

Indicus: Unchaining Byzantine Databases

Your N. Here
Your Institution

Second Name
Second Institution

Abstract

Your abstract text goes here. Just a few facts. Whet our appetites. Not more than 200 words, if possible, and preferably closer to 150.

1 Introduction

There is an increasing tension between the desire to share data online, and the security concerns it entails. This paper asks the question: how can we enable *mutually distrustful parties* to consistently and reliably share data, while minimizing centralization?

The ability to share data online offers exciting opportunities. In banking, systems like SWIFT enable financial institutions to quickly and accurately receive information such as money transfer instructions; and in manufacturing, online data sharing can improve accountability and auditing amongst the globally distributed supply chain.

Increased data sharing, however, raises questions of how to *decentralize trust*. Banking institutions must currently place their trust in the centralized SWIFT's network to issue payment orders. Sometimes there is even no identifiable source of trust. Consider the supply chain for the latest iPhone: it spans three continents, and hundreds of different contractors []; neither Apple nor these contractors trust each other, yet all must be willing to agree and share information about the construction of the same product.

Recognizing this challenge by both the research and industry communities, much effort has focused on enabling shared computation between mutually distrustful parties, in the context of byzantine fault tolerance (BFT), and blockchains. Systems proposed in the literature of BFT[][] provide the abstraction of a totally ordered log; the log is agreed upon by the n participants in the system, of which at most f can misbehave. Each participant executes operations that may touch one to multiple objects in the log. In the blockchain world, Bitcoin and Ethereum have become popular distributed computing platforms providing the same log abstraction and aiming

for decentralizing trust. Furthermore, Microsoft Azure has launched projects[] that leverage these blockchain platforms and extend the digital transformation beyond the companies' four walls in a supply chain.

This paper argues that there exists a fundamental mismatch between the implementation of a totally ordered log and the reality of much large-scale distributed processing. Many large-scale distributed systems consist primarily of unordered and unrelated operations. For example, a product supply chain consists of many concurrent steps that do not require ordering. Imposing an ordering on non-conflicting operations is not only often unnecessary, but costly: participants in the shared computation must vote to order operations, store the full state of the system, and replay the full log for auditing.

While there exists work on mitigating this scalability bottleneck through sharding [], the latent total order requirement introduces unnecessary coordination overhead, as coordination is performed twice, at the level of individual shards, and across shards. Callinicos [] and Omniledger [], for instance, runs a full BFT protocol for every operation. This is especially problematic when workloads are geo-replicated [], or when, as in BFT, the replication factor is high. Further, these systems support transactions under the assumption that their read and write operations are known a priori, which limits the set of applications that they can support.

As another research trend of mitigating the scalability bottleneck, EPaxos[], TAPIR[] and CURP[] only consider the ordering between potentially conflicting operations, instead of commutative operations. However, these systems assume the crash-failure model and are non-trivial to be extended to the Byzantine model, so that they cannot directly solve the problem of data sharing among mutually distrustful parties.

Existing research, in essence, is either attempting to build concurrency control and sharding functionalities over BFT replication, or integrating these functionalities into a crash-failure replication protocol. In this paper, we will show how to build these desiring functionality inside a BFT replication protocol. Specifically, our goal is to *provide the illusion of a centralized shared log, rather than the non-scalable reality of*

a totally ordered log.

2 Indicus

In the following we outline the Architecture and Protocols of Indicus, the first highly scalable database that tolerates both byzantine clients and replicas.

2.1 Architecture

Indicus follows a traditional optimistic concurrency control architecture as shown in Figure 1. Clients speculatively execute their Transactions, issuing remote reads when necessary and buffering writes locally. As Clients execute their own Transactions they need not declare their read/write keys or values preemptively, but rather may conduct interactive Transactions, the most general Transaction model. Since execution is speculative and Clients are unaware of potential concurrency, they must validate their Transactions for Isolation correctness in order to be able to Commit. Intuitively, if there are no conflicting concurrent Transactions and a Client observed a consistent snapshot of the distributed Database state then it may commit, and otherwise it must abort or retry. Lastly, since Transactions may span multiple shards, the Commit/Abort decisions of each shard must be aggregated in a Two Phase Commit manner, finalizing the result. Most notably, in Indicus not just execution is Client driven but the entire Transaction life cycle, maximizing scalability, fairness and putting each Client in charge of its own liveness. In the following we describe our Fault Model and define corresponding Transactional guarantees. We then outline the protocols for Execution, Validation and Writeback respectively.

2.2 Model and Definitions

2.2.1 System Model

We assume a partitioned Database where data objects are spread uniformly across shards (randomly or based on locality). For fault tolerance we assume that each Shard contains $5f+1$ replicas where f is the number of (static) faulty or compromised replicas. A faulty replica may behave in Byzantine fashion, i.e. in addition to crashing or omitting messages it may send arbitrary messages and deviate from prescribed protocols in any way. We denote any participant (replica or client) that follows the protocol as *honest*, while faulty participants are dubbed *byzantine*. We assume there may exist a finite but unbounded number of byzantine Clients that may deviate from the protocol arbitrarily. We further assume a strong, but static adversary that can coordinate the faulty participants. We do, however, assume the existence of sufficiently hard cryptographic primitives that allow for private/public key signatures and collision-resistant hashes that cannot be compromised by

byzantine participants. We denote a signed message m , signed by principal p as $\langle m \rangle_{\sigma p}$.

We make no assumption on network synchrony in order to maintain safety, but in some cases may provide liveness only when the network is synchronous and messages are delayed by no more than a fixed but potentially unknown window. Unlike traditional State Machine Replication protocols in which the liveness of all Clients is correlated with the fate of the system (or often more specifically a leader), our system guarantees liveness not on a system basis, but on a per client basis. Concretely, we only guarantee liveness to clients that follow the protocol and conversely an honest client only loses liveness (even in an asynchronous setting) when it intertwines its fate with byzantine clients.

Application services may restrict the influence of Byzantine Clients on the Database state by authentication and enforcing access control. Beyond these measures, byzantine Clients that follow the protocol are indistinguishable from honest clients and may read and write at their leisure. While potential damage to the Database state cannot be avoided, it can be re-traced by auditing the Transaction logs.

2.2.2 System properties

We offer Clients an interactive Transaction interface that implements ACID transactions. While all ACID guarantees may be violated by individual byzantine Replicas, the system maintains the illusion of an ACID compliant state to all honest Clients. Moreover, we guarantee to Clients that the Database is *byzantine Serializable* as defined below. Intuitively this Isolation level guarantees that all honest Clients experience the Database as if there serializable. In order to formally capture this we lay some ground work:

- history - projection -

This definition captures the requirements for any byzantine tolerant protocol that strives to maintain Isolation level I.

We further define some progress properties to limit the influence Byzantine Clients have on system throughput. These will be motivating factors that we will return to in order to justify our choice of $5f+1$ replicas as well as the use of an Optimistic Concurrency Control over a Locking based schema. We later outline an alternative design using $3f+1$ replicas that does not maintain this property.

2.3 Execution

2.4 Validation

2.5 Concurrency Control

- Optimization: retries - Dependency resolution

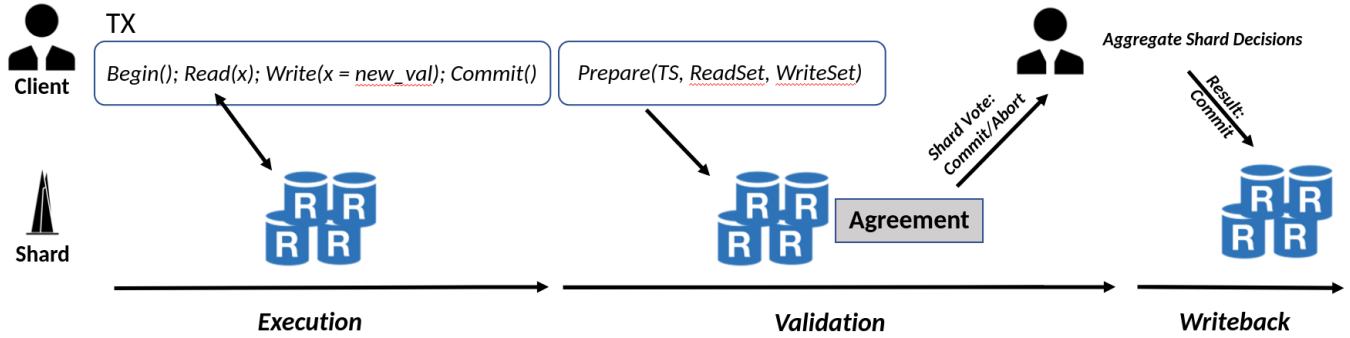


Figure 1: Transaction Lifecycle

2.6 Consistent logging.

2.7 Writeback and Multi-shard 2pc

- Optimization: Single shard logging

2.8 Failures

- Fallback: election, resolution, subtleties with mvts, necessity even without dependencies.

2.9 Low Cost mode

3f+1 if not defending against byz colluders OCC if not worried about reads aborting

3 Implementation and Evaluation

4 Limitations

As is inherent to any Optimistic Execution and Concurrency Control, Indicus is vulnerable to highly congested workloads. When contention on select objects is high, concurrent execution of Transactions must yield in the abort of some Transactions during Validation in order to maintain the Database Isolation guarantees. Note however, that when clients are in charge of execution, a pessimistic concurrency control solution such as two-phase-locking would incur an equal amount of deadlocks which would require resolution. The observation to make is that any system that conducts execution at the client application side speculates on concurrency. This however we stipulate, is unavoidable when trying to scale a system to the number of users rather than replica processing power. The traditional ways to avoid the abort rate conundrum is to either restrict the transaction model, which in turn weakens the general applicability of the protocol, or to delegate execution to replicas and utilize State Machine Replication to serialize Transactions. Indicus does not make these concessions in order to offer interactive Transactions and remain

scalable. A workload that exhibits low commutativity and high contention should therefore refrain from adopting our system.

Similarly, as is the case in any transaction protocol, Indicus is vulnerable to ddos attacks by byzantine participants. A byzantine clients only opportunity at subverting progress for honest users is to artificially increase congestion. When such a client has unrestricted access control it may do so strategically iff it has control over the network. If it does not, it cannot reliably gain knowledge about concurrent transactions before they pass the validation step and must resort to flooding based attacks. Defense against such attacks is out of scope in our work, but is disincentivised as participants can be held accountable for their actions in a closed membership setting.

5 Related Work/Comparison

A lot of recent effort has gone into designing high throughput and low latency databases that leverage synergies between transaction and replication layer to squeeze out any last performance. The recently proposed TAPIR transaction protocol leverages redundancy between transaction ordering and replication ordering to reduce total roundtrips, thus reducing latency in wide area networks. TAPIR shares several similarities with our system, most notably the absence of a leader and the resulting unordered validation structure. TAPIR too, leverages optimistic concurrency control to allow for concurrency among commutative transactions. When congestion is high however, throughput and tail latency worsen as abort rate grows. Janus avoids transaction aborts by dynamically re-ordering conflicting transactions. This however, is made possible by assuming one-shot transactions, i.e. fixed read/write keys, and thus reduces the application generality. Another DB, Carousel similarly assumes a restricted Transaction model in order to parallelize execution and validation. While all of these databases offer low latency replication, they are fundamentally limited to tolerating crash-failures. This strong failure assumption makes them less secure and not suitable for storing mission-critical, financial or highly sensitive data.

With the surge in Blockchain interest Byzantine Fault Tolerant (BFT) protocols are experiencing a second Spring. While originally developed to tolerate arbitrary bugs, these protocols find increasing importance in settings where participants are untrusted or malicious. Permissioned Blockchains, organized by a consortium of registered participants can use traditional BFT State Machine Replication (SMR) protocols in order to achieve agreement. Starting with PBFT, there have been numerous adaptations such as FaB or Generalized byzantine Paxos, and most notably Zyzzyva which leverages speculative execution and a semi-client driven protocol to reduce latency. SBFT modifies PBFT to scale to large replication degrees by utilizing collectors, threshold signatures and a fast path akin to Zyzzyva. Aardvark states the importance of robustness against byzantine failures and takes measures to increase the rate of leader rotation. Nevertheless, all protocols derived from the PBFT family suffer from the leader bottleneck as well as enforcing a total order even on commutative operations. Biblos achieves leaderless SMR by leveraging a non-skipping timestamp protocol. It furthermore allows for commutative Transactions to be executed in parallel at the cost of requiring read/write key sets to be known in advance. In order to preserve liveness however Bilbos falls back to a PBFT resolution. Q/U too, offers leaderless agreement via Quorums for a limited read/write interface, but fails to terminate under contention. H/Q improves upon Q/U by adding PBFT fallback path under contention. Liskov et Al further explore byzantine Quorum protocols for a Read/Write interface, giving special attention to bounding the effects of byzantine clients. While the literature on BFT state machine replication is extensive, the efforts to offer a transactional interface for BFT is scarce. SMR itself can be utilized as a straw-man system to implement pre-defined Transactions (one-shot or stored procedures) by enforcing a common total order in which replicas will execute the transactions. HRDB offers a dedicated BFT Database, but assumes a trusted shepherd layer, thus not being truly BF resilient. Byzantium offers Snapshot Isolation for Transactions by re-purposing PBFT as atomic broadcast. It executes requests only at a primary and uses replicas to validate results in the total order defined by the SMR protocol. Augustus implements mini-transactions, a limited TX model that declares all operations before execution. It offers

scalability via partitioning and achieves consistency within partitions by utilizing atomic broadcast. While Augustus assumes an optimistic execution model that allows for aborts under concurrency, its follow up work Callinicos implements a locking scheme in order to avoid Transaction aborts. To do so efficiently it resorts to a limited transaction model that requires knowledge of read/write sets and otherwise locks an entire partition in order to guarantee mutual exclusion, thus voiding any concurrency. BFT Deferred update replication adopts an interactive OCC transaction model comparable to ours, allowing execution to be speculative at clients. However, it uses PBFT atomic broadcast to enforce SMR on validation, thus resorting to a leader, and totally ordering all requests. It moreover does not extend to multiple shards. Chainspace and Omniledger implement blockchain sharding by layering 2PC and atomic broadcast (within shards) for UTXO transactions. Rapidchain offers efficient sharding for a permissionless system.

[1]

6 Conclusion

Acknowledgments

The USENIX latex style is old and very tired, which is why there's no \acks command for you to use when acknowledging. Sorry.

Availability

USENIX program committees give extra points to submissions that are backed by artifacts that are publicly available. If you made your code or data available, it's worth mentioning this fact in a dedicated section.

References

- [1] Ramakrishna Kotla, Lorenzo Alvisi, Mike Dahlin, Allen Clement, and Edmund Wong. Zyzzyva: speculative byzantine fault tolerance. *ACM SIGOPS Operating Systems Review*, 41(6):45–58, 2007.