

Building Consistent Transactions with Inconsistent Replication

IRENE ZHANG, Microsoft Research

NAVEEN KR. SHARMA, ADRIANA SZEKERES, and ARVIND KRISHNAMURTHY,
University of Washington

DAN R. K. PORTS, Microsoft Research

Application programmers increasingly prefer distributed storage systems with strong consistency and distributed transactions (e.g., Google’s Spanner) for their strong guarantees and ease of use. Unfortunately, existing transactional storage systems are expensive to use—in part, because they require costly replication protocols, like Paxos, for fault tolerance. In this article, we present a new approach that makes transactional storage systems more affordable: We eliminate consistency from the replication protocol, while still providing distributed transactions with strong consistency to applications.

We present the Transactional Application Protocol for Inconsistent Replication (TAPIR), the first transaction protocol to use a novel replication protocol, called *inconsistent replication*, that provides fault tolerance without consistency. By enforcing strong consistency only in the transaction protocol, TAPIR can commit transactions in a single round-trip and order distributed transactions without centralized coordination. We

This work is an extended version of the paper by the same title that appeared in SOSP 2015 (Zhang et al. 2015a). The additional content includes:

- (1) A complete description of the IR view change protocol (Section 3.2.2)
- (2) A description of the IR client recovery protocol (Section 3.2.3)
- (3) A full proof of correctness for IR (Section 3.3)
- (4) Additional pseudocode for TAPIR-EXEC-CONSENSUS, which executes TAPIR’s `Prepare` operation, and TAPIR-SYNC, which synchronizes replicas with missed IR operations and consensus results. (Section 5)
- (5) The complete coordinator recovery protocol for TAPIR (Section 5.2.3)
- (6) A full proof of correctness for for TAPIR on IR (Section 5.3)
- (7) Extensions to the TAPIR protocol for:
 - (a) Supporting read-only transactions at a consistent timestamp and Spanner-style linearizable read-only transactions (Section 6.1)
 - (b) Relaxing from linearizable transaction ordering to serializable (Section 6.2)
 - (c) Optimizing retry timestamp selections for greater transaction success rates (Section 6.3).
 - (d) Coping with very high clock skews (Section 6.4)
- (8) A full latency and clock skew profile of our Google Compute Engine testbed (Section 7.1.1)
- (9) An evaluation of performance during node failures and recovery (Section 7.6)

This work was supported by the National Science Foundation under Grants No. CNS-0963754, No. CNS-1217597, No. CNS-1318396, No. CNS-1420703, No. CNS-1518702, and No. CNS-1615102, by NSF GRFP and IBM Ph.D. fellowships, and by gifts from Google and VMware.

Authors’ addresses: I. Zhang and D. R. K. Ports, Microsoft Research, 14820 NE 36th St. Redmond, WA 98052; emails: Irene.Zhang@microsoft.com, dan@drkp.net; N. K. Sharma, A. Szekeres, and A. Krishnamurthy, Department of Computer Science & Engineering, Box 352350, University of Washington, Seattle, WA 98195; emails: {naveenks, aaasz, arvind}@cs.washington.edu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

0734-2071/2018/12-ART12 \$15.00

<https://doi.org/10.1145/3269981>

demonstrate the use of TAPIR in a transactional key-value store, TAPIR-KV. Compared to conventional systems, TAPIR-KV provides better latency *and* better throughput.

CCS Concepts: • **Information systems** → **Distributed database transactions**; • **Computer systems organization** → **Distributed architectures**; **Reliability**;

Additional Key Words and Phrases: Distributed transactional storage, inconsistent replication, strict serializability

ACM Reference format:

Irene Zhang, Naveen Kr. Sharma, Adriana Szekeres, Arvind Krishnamurthy, and Dan R. K. Ports. 2018. Building Consistent Transactions with Inconsistent Replication. *ACM Trans. Comput. Syst.* 35, 4, Article 12 (December 2018), 37 pages.

<https://doi.org/10.1145/3269981>

1 INTRODUCTION

Distributed storage systems provide fault tolerance and availability for large-scale web applications. Increasingly, application programmers prefer systems that support distributed transactions with strong consistency to help them manage application complexity and concurrency in a distributed environment. Several recent systems (Kraska et al. 2013; Baker et al. 2011; Escriva et al. 2013; Cooper et al. 2008) reflect this trend, notably Google’s Spanner system (Corbett et al. 2012), which guarantees linearizable transaction ordering.¹

For application programmers, distributed transactional storage with strong consistency comes at a price. These systems commonly use replication for fault-tolerance, and replication protocols with strong consistency, like Paxos, impose a high performance cost, while more efficient, weak consistency protocols fail to provide strong system guarantees.

Significant prior work has addressed improving the performance of transactional storage systems—including systems that optimize for read-only transactions (Baker et al. 2011; Corbett et al. 2012), more restrictive transaction models (Kraska et al. 2013; Aguilera et al. 2007; Cowling and Liskov 2012), or weaker consistency guarantees (Lloyd et al. 2011; Sovran et al. 2011; Bailis et al. 2014). However, none of these systems have addressed both latency *and* throughput for general-purpose, replicated, read-write transactions with strong consistency.

In this article, we use a new approach to reduce the cost of replicated, read-write transactions and make transactional storage more affordable for programmers. Our key insight is that existing transactional storage systems waste work and performance by incorporating a distributed transaction protocol and a replication protocol that *both* enforce strong consistency. Instead, we show that it is possible to provide distributed transactions with better performance and the same transaction and consistency model using replication with *no consistency*.

To demonstrate our approach, we designed the Transactional Application Protocol for Inconsistent Replication (TAPIR), which uses a new replication technique, *inconsistent replication* (IR), that provides fault tolerance without consistency. Rather than an ordered operation log, IR presents an *unordered operation set* to applications. Successful operations execute at a majority of the replicas and survive failures, but replicas can execute them in any order. Thus, IR needs no cross-replica coordination or designated leader for operation processing. However, unlike eventual consistency, IR allows applications to enforce higher-level invariants when needed.

Thus, despite IR’s weak consistency guarantees, TAPIR provides *linearizable read-write transactions* and supports globally-consistent reads across the database at a timestamp—the same

¹Spanner’s linearizable transaction ordering is also referred to as strict serializable isolation or external consistency.

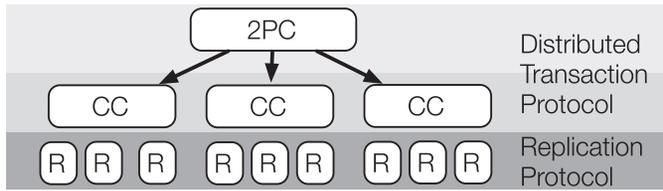


Fig. 1. A common architecture for distributed transactional storage systems today. The distributed transaction protocol consists of an atomic commitment protocol, commonly Two-Phase Commit (2PC), and a concurrency control (CC) mechanism. This runs atop a replication (R) protocol, like Paxos.

guarantees as Spanner. TAPIR efficiently leverages IR to distribute read-write transactions in a *single round-trip* and order transactions globally across shards and replicas *with no centralized coordination*.

We implemented TAPIR in a new distributed transactional key-value store called TAPIR-KV, which supports linearizable transactions over a sharded set of keys. Our experiments found that TAPIR-KV had: (1) 50% lower commit latency, (2) more than 3× better throughput compared to systems using conventional transaction protocols, including an implementation of Spanner’s transaction protocol, and (3) comparable performance to MongoDB (MongoDB 2013) and Redis (Redis 2013), widely used eventual consistency systems.

This article makes the following contributions to the design of distributed, replicated transaction systems:

- We define *inconsistent replication*, a new replication technique that provides fault tolerance without consistency.
- We design *TAPIR*, a new distributed transaction protocol that provides strict serializable transactions using inconsistent replication for fault tolerance.
- We build and evaluate TAPIR-KV, a key-value store that combines inconsistent replication and TAPIR to achieve high-performance transactional storage.

2 OVER-COORDINATION IN TRANSACTION SYSTEMS

Replication protocols have become an important component in distributed storage systems. Modern storage systems commonly partition data into *shards* for scalability and then replicate each shard for fault-tolerance and availability (Baker et al. 2011; Chang et al. 2008; Corbett et al. 2012; Mahmoud et al. 2013). To support transactions with strong consistency, they must implement both a *distributed transaction protocol*—to ensure atomicity and consistency for transactions across shards—and a *replication protocol*—to ensure transactions are not lost (provided that no more than half of the replicas in each shard fail at once). As shown in Figure 1, these systems typically place the transaction protocol, which combines an atomic commitment protocol and a concurrency control mechanism, on top of the replication protocol (although alternative architectures have also occasionally been proposed (Mahmoud et al. 2013)).

Distributed transaction protocols assume the availability of an *ordered, fault-tolerant log*. This ordered log abstraction is easily and efficiently implemented with a spinning disk but becomes more complicated and expensive with replication. To enforce the serial ordering of log operations, transactional storage systems must integrate a costly replication protocol with strong consistency (e.g., Paxos (Lamport 2001), Viewstamped Replication (Oki and Liskov 1988) or virtual synchrony (Birman and Joseph 1987)) rather than a more efficient, weak consistency protocol (Ladin et al. 1992; Saito and Shapiro 2005).

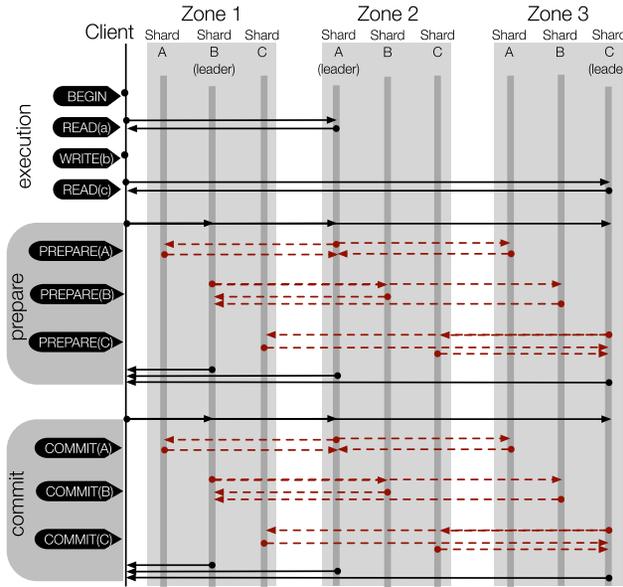


Fig. 2. Example read-write transaction using two-phase commit, viewstamped replication, and strict two-phase locking. Availability zones represent a cluster, datacenter, or geographic region. Each *shard* holds a partition of the data stored in the system and has replicas in each zone for fault tolerance. The red, dashed lines represent redundant coordination in the replication layer.

The traditional log abstraction imposes a serious performance penalty on replicated transactional storage systems, because it enforces strict serial ordering using expensive distributed coordination *in two places*: the replication protocol enforces a serial ordering of operations across replicas in each shard, while the distributed transaction protocol enforces a serial ordering of transactions across shards. This redundancy impairs latency and throughput for systems that integrate both protocols. The replication protocol must coordinate across replicas on every operation to enforce strong consistency; as a result, it takes *at least two round-trips* to order any read-write transaction. Further, to efficiently order operations, these protocols typically rely on a replica leader, which can introduce a throughput bottleneck to the system.

As an example, Figure 2 shows the redundant coordination required for a single read-write transaction in a system like Spanner. Within the transaction, Read operations go to the shard leaders (which may be in other datacenters), because the operations must be ordered across replicas, even though they are not replicated. To Prepare a transaction for commit, the transaction protocol must coordinate transaction ordering across shards, and then the replication protocol coordinates the Prepare operation ordering across replicas. As a result, it takes at least two round-trips to commit the transaction.

In the TAPIR and IR design, we eliminate the redundancy of strict serial ordering over the two layers and its associated performance costs. IR is the first replication protocol to provide *pure fault tolerance* without consistency. Instead of an ordered operation log, IR presents the abstraction of an *unordered operation set*. Existing transaction protocols cannot efficiently use IR, so TAPIR is the first transaction protocol designed to provide linearizable transactions on IR.

3 INCONSISTENT REPLICATION

IR is an efficient replication protocol designed to be used with a higher-level protocol, like a distributed transaction protocol. IR provides fault-tolerance without enforcing any consistency

guarantees of its own. Instead, it allows the higher-level protocol, which we refer to as the *application protocol*, to decide the outcome of conflicting operations and recover those decisions through IR's fault-tolerant, unordered operation set.

3.1 IR Overview

IR is similar to other state machine replication-based protocols: application clients invoke application-provided operations on replicas through IR for fault-tolerance. Note that the applications for IR are a higher-level protocol, like TAPIR. For example, TAPIR uses IR to invoke prepare across replicas in a shard to ensure that prepare is fault-tolerant. Application protocols invoke operations through IR in one of two modes:

- **inconsistent**—IR can execute application operations in any order at replicas. IR guarantees that successful operations persist across failures. On recovery, IR may re-execute operations in any order.
- **consensus**—IR can execute application operations in any order at replicas, but the operation must return a single *consensus result*. The consensus result is decided by either a majority of the replicas returning the same result to the application operation or running an application provided *decide* function. Successful operations and their consensus results persist across failures. On recovery, IR will re-execute operations with the provided consensus result.

inconsistent operations are similar to operations in weak consistency replication protocols: they can execute in different orders at each replica, and the application protocol must resolve conflicts afterwards. **inconsistent** operations do not always return results, and indeed none of the **inconsistent** operations in TAPIR do (save for the optional read-only extension in Section 6.1). If they do, then each replica independently produces its own result, and the client receives a *set* of results—the result returned by each of the replicas that executed the operation. IR does not maintain or enforce any consistency for these results; it is up to the application protocol to interpret them.

In contrast, **consensus** operations either require either: (1) the application operation returns the same result after unordered execution at the replicas, or (2) the application protocol resolves conflicts with a provided *decide* function. The *decide* function can choose result that it wants, even one that is not returned by any of the replicas's execution of the operation. IR ensures that the application can recover the decision after failure by preserving the consensus result. In this way, **consensus** operations can serve as the basic building block for the higher-level guarantees of application protocols. For example, TAPIR decides which conflicting transaction will commit and which will abort, while IR will ensure that decision persists across failures.

3.1.1 IR Application Protocol Interface. Figure 3 summarizes the IR interfaces at clients and replicas. Application protocols invoke operations through a client-side IR library using `InvokeInconsistent` and `InvokeConsensus`, and then IR runs operations using the `ExecInconsistent` and `ExecConsensus` upcalls into application functions at the replicas.

If replicas return conflicting/non-matching results for a **consensus** operation, then IR allows the application protocol to decide the operation's outcome by invoking the *decide* function—passed in by the application protocol to `InvokeConsensus`—in the client-side library. The *decide* function takes the list of returned results (the candidate results) and returns a single result, which IR ensures will persist as the *consensus result*. The application protocol can later recover the consensus result to find out its decision to conflicting operations.

Because IR replicas can be inconsistent, the application protocol at each replica may need to “fix” its state periodically. These inconsistencies appear in two ways. First, at each replica,

Client Interface

- $\text{InvokeInconsistent}(op) \rightarrow \text{results}$ - Invoke op across replicas as an **inconsistent** operation and return. Once returned, IR guarantees fault-tolerance for op . If the operations have results, a set of $f + 1$ results (one from each replica) is returned.
- $\text{InvokeConsensus}(op, \text{decide}(\text{results})) \rightarrow \text{result}$ - Invoke op across replicas as a **consensus** operation. If replicas return same results, return as result ; otherwise, invoke decide with different returned results and return result from decide . Once returned, IR guarantees fault-tolerance for op and returned result .

Replica Upcalls

- $\text{ExecInconsistent}(op) \rightarrow \text{result}$ - Upcall to application to execute inconsistent op . The result is optional, as many operations do not return results.
- $\text{ExecConsensus}(op) \rightarrow \text{result}$ - Upcall to application to execute op . Application returns result to be collected for deciding consensus result.
- $\text{Sync}(R)$ - Upcall to application with finalized results to let application fix any inconsistencies.
- $\text{Merge}(d, u) \rightarrow \text{record}$ - Upcall to application to make merge decision for d – **inconsistent** operations and finalized **consensus** operations – and u – tentative consensus operations. Only needed for unfinished consensus operations during synchronization and recovery.

Client State

- client id - unique identifier for the client
- operation counter - # of sent operations

Replica State

- state - current replica state; either NORMAL or VIEW-CHANGING
- record - unordered set of inconsistent and consensus operations. For each operation, the record contains:
 - id - The operation identifier, which consists of the issuing client's id and operation counter
 - operation type - The type of operation, either **inconsistent** or **consensus**
 - operation state - The state of the operation, either TENTATIVE or FINALIZED
 - result - The application returned result for the operation. If the operation is a **consensus** operation, then, if the operation state is TENTATIVE, this field holds the result from this replica, and, if the operation state is FINALIZED, then this field holds the consensus result.

Fig. 3. Summary of IR interfaces and client/replica state.

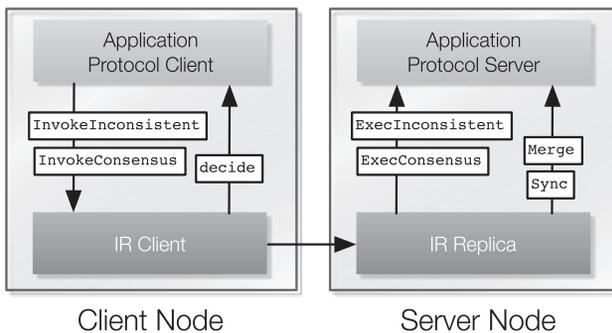


Fig. 4. IR call flow.

operations in the record progress through two states. Operations in the `TENTATIVE` state may not persist (i.e., they may be lost later due to failure) or their results may differ from the consensus result. Operations in the `FINALIZED` will never be lost and their result matches the consensus result. Next, IR replica periodically *synchronize* to ensure that their records converge and all operations eventually move to the `FINALIZED` state (e.g., in case of client failure). Similar to eventual consistency protocols, IR relies on the application to reconcile inconsistent replicas; however, we expect the application to be a higher-level protocol, like TAPIR, not an end-user application. On synchronization, a designated IR replica, which we call the leader for simplicity, collects records from a quorum of inconsistent replicas. The leader then upcalls into the application protocol with the `Merge`, which merges the collected records into a *master record* of successful operations and consensus results. The `Merge` function does not have to be deterministic, because only one leader will call it at a time. The leader sends the master replica back to the other replicas, which each upcall into the application protocol with `Sync`. `Sync` gives the *master record* to the application protocol at each replica to make the replica consistent with the chosen consensus results.

3.1.2 IR Guarantees. We define a *successful operation* to be one that returns to the application protocol and a *persistent operation* to be one that is guaranteed to survive failures. IR will never return a successful operation back to the application protocol unless it is also persistent. However, persistent operations may not have been returned successfully back to the application protocol (e.g., if the application client fails).

The *operation set* of any IR group is guaranteed to include all persistent operations. We define an operation X as being *visible* to an operation Y if one of the replicas executing Y has previously executed X . IR ensures the following properties for the operation set:

- P1. [Fault tolerance].* At any time, every operation in the operation set is in the record of at least one replica in any quorum of $f + 1$ non-failed replicas.
- P2. [Visibility].* For any two operations in the operation set, at least one is visible to the other.
- P3. [Consensus results].* At any time, the result returned by a successful **consensus** operation is in the record of at least one replica in any quorum. The only exception is if the consensus result has been explicitly modified by the application protocol through `Merge`, after which the outcome of `Merge` will be recorded instead.

IR ensures guarantees are met for up to f simultaneous failures out of $2f + 1$ replicas and any number of client failures. Replicas must be fail-stop, without Byzantine behavior. We assume an asynchronous network where messages can be lost or delivered out of order. IR *does not* require synchronous disk writes during operation execution, ensuring guarantees are maintained even if clients or replicas lose disks on failure. However, if more than f replicas fail (e.g., in the case of a power failure), IR can lose data. IR makes progress (operations will eventually become successful) provided that messages that are repeatedly resent are eventually delivered before the recipients time out.

3.1.3 Application Protocol Example: Fault-Tolerant Lock Server. As an example, we show how to build a simple lock server using IR. The lock server's guarantee is mutual exclusion: a lock cannot be held by two clients at once. We replicate `Lock` as a **consensus** operation and `Unlock` as an **inconsistent** operation. A client application acquires the lock only if `Lock` successfully returns `OK` as a consensus result.

Due to operations executing in different orders on different replicas, clients must use unique ids (e.g., a tuple of client id and a sequence number) to identify corresponding `Lock` and `Unlock`

operations and call `Unlock` if `Lock` first succeeds. Replicas will therefore be able to later match up `Lock` and `Unlock` operations, regardless of order, and determine the lock's status.

Note that **inconsistent** operations *are not commutative*, because they can have side-effects that affect the outcome of **consensus** operations. If an `Unlock` and `Lock` execute in different orders at different replicas, then some replicas might have the lock free, while others might not. If replicas return different results from `Lock`, then IR invokes the lock server's `decide` function, which returns `OK` if $f + 1$ replicas returned `OK` and `NO` otherwise. IR invokes `Merge` and `Sync` only during synchronization and recovery, so we defer their discussion until Section 3.2.2.

IR's guarantees ensure correctness for our lock server. P1 ensures that held locks are persistent: a `Lock` operation persists at one or more replicas in any quorum. P2 ensures mutual exclusion: for any two conflicting `Lock` operations, one is visible to the other in any quorum, so IR will never receive $f + 1$ matching `OK` results. The `decide` function is designed to only return `OK` if at $f + 1$ replicas return `OK`, so two conflicting `Lock` operations will never both return `OK`. P3 ensures that once the client application receives `OK` from a `Lock`, the result will not change. The lock server's `Merge` function will not change it, as we will show later, and IR ensures that the `OK` will persist in the record of at least one replica out of any quorum.

3.2 IR Protocol

Figure 3 shows the IR state at the clients and replicas. Each IR client keeps an *operation counter*, which, combined with the *client id*, uniquely identifies operations. Each replica keeps an unordered *record* of executed operations and results for **consensus** operations. Replicas add **inconsistent** operations to their record as `TENTATIVE` and then mark them as `FINALIZED` once they execute. **consensus** operations are first marked `TENTATIVE` with the result of locally executing the operation, then `FINALIZED` once the record has the consensus result.

IR uses four sub-protocols—*operation processing*, *replica recovery/synchronization*, *client recovery*, and *group membership change*. We discuss the first three here; the last is identical to that of Viewstamped Replication (Liskov and Cowling 2012).

3.2.1 Operation Processing. We begin by describing IR's normal-case **inconsistent** operation processing protocol without failures:

- (1) The client sends $\langle \text{PROPOSE}, id, op \rangle$ to all replicas, where id is the operation id and op is the operation.
- (2) Each replica writes id and op to its record as `TENTATIVE`, then responds to the client with $\langle \text{REPLY}, id \rangle$.
- (3) Once the client receives $f + 1$ responses from replicas (retrying if necessary), it returns to the application protocol and asynchronously sends $\langle \text{FINALIZE}, id \rangle$ to all replicas. (`FINALIZE` can also be piggy-backed on the client's next message.)
- (4) On `FINALIZE`, replicas upcall into the application protocol with `ExecInconsistent(op)` and mark op as `FINALIZED`.

Due to the lack of consistency, IR can successfully complete an **inconsistent** operation with a *single round-trip* to $f + 1$ replicas and no coordination across replicas. If the IR client does not receive a response to its `PREPARE` from $f + 1$ replicas, then it will retry until it does.

Next, we describe the normal-case **consensus** operation processing protocol, which has both a *fast path* and a *slow path*. IR uses the fast path when it can achieve a *fast quorum* of $\lceil \frac{3}{2}f \rceil + 1$ replicas that *return matching results* to the operation. Similar to Fast Paxos (Lamport 2006a) and Speculative Paxos (Ports et al. 2015), IR requires a fast quorum to ensure that a majority of the replicas in any regular majority quorum agrees to the consensus result. This quorum size is necessary to execute

operations in a single round-trip when using a replica group of size $2f + 1$ (Lamport 2006b); an alternative would be to use quorums of size $2f + 1$ in a system with $3f + 1$ replicas.

When IR cannot achieve a fast quorum, either because replicas did not return enough matching results (e.g., if there are conflicting concurrent operations) or because not enough replicas responded (e.g., if more than $\frac{f}{2}$ are down), then it must take the slow path. We describe both below:

- (1) The client sends $\langle \text{PROPOSE}, id, op \rangle$ to all replicas.
- (2) Each replica calls into the application protocol with $\text{ExecConsensus}(op)$ and writes id, op , and $result$ to its record as `TENTATIVE`. The replica responds to the client with $\langle \text{REPLY}, id, result \rangle$.
- (3) If the client receives at least $\lceil \frac{3}{2}f \rceil + 1$ matching $results$ (within a timeout), then it takes the *fast path*: the client returns $result$ to the application protocol and asynchronously sends $\langle \text{FINALIZE}, id, result \rangle$ to all replicas.
- (4) Otherwise, the client takes the *slow path*: once it receives $f + 1$ responses (retrying if necessary), then it sends $\langle \text{FINALIZE}, id, result \rangle$ to all replicas, where $result$ is obtained from executing the *decide* function.
- (5) On receiving `FINALIZE`, each replica marks the operation as `FINALIZED`, updating its record if the received $result$ is different, and sends $\langle \text{CONFIRM}, id \rangle$ to the client.
- (6) On the slow path, the client returns $result$ to the application protocol once it has received $f + 1$ `CONFIRM` responses.

The fast path for **consensus** operations takes a single round-trip to $\lceil \frac{3}{2}f \rceil + 1$ replicas, while the slow path requires two round-trips to at least $f + 1$ replicas. Note that IR replicas can execute operations in different orders and *still* return matching responses, so IR can use the fast path without a strict serial ordering of operations across replicas. IR can also run the fast path and slow path in parallel as an optimization.

3.2.2 Replica Recovery and Synchronization. IR uses a single protocol for recovering failed replicas and running periodic synchronizations. On recovery, we must ensure that the failed replica applies all operations it may have lost or missed in the operation set, so we use the same protocol to periodically bring all replicas up-to-date.

To handle recovery and synchronization, we introduce *view changes* into the IR protocol, similar to Viewstamped Replication (VR) (Oki and Liskov 1988). These maintain IR's correctness guarantees across failures. Each IR view change is run by a leader, which is one of the replicas. Leaders coordinate only view changes, *not* operation processing. During a view change, the leader has just one task: to make at least $f + 1$ replicas up-to-date (i.e., they have applied all operations in the operation set) and consistent with each other (i.e., they have applied the same consensus results). IR view changes require a leader, because polling inconsistent replicas can lead to conflicting sets of operations and consensus results. Thus, the leader must decide on a *master record* that replicas can then use to synchronize with each other.

To support view changes, each IR replica maintains a current *view*, which consists of the identity of the leader, a list of the replicas in the group, and a (monotonically increasing) *view number* uniquely identifying the view. Each IR replica can be in one of the three states: `NORMAL`, `VIEW-CHANGING`, or `RECOVERING`. Replicas process operations only in the `NORMAL` state. We make four additions to IR's operation processing protocol:

- (1) IR replicas send their current view number in every response to clients. For an operation to be considered successful, the IR client must receive responses with matching view

numbers. For **consensus** operations, the view numbers in REPLY and CONFIRM must match as well. If a client receives responses with different view numbers, then it notifies the replicas in the older view.

- (2) On receiving a message with a view number that is higher than its current view, a replica moves to the VIEW-CHANGING state and requests the master record from any replica in the higher view. It replaces its own record with the master record and upcalls into the application protocol with Sync before returning to NORMAL state.
- (3) On PROPOSE, each replica first checks whether the operation was already FINALIZED by a view change. If so, then the replica responds with $\langle \text{REPLY}, id, \text{FINALIZED}, v, [result] \rangle$, where v is the replica's current view number and $result$ is the consensus result for **consensus** operations.
- (4) If the client receives REPLY with a FINALIZED status for **consensus** operations, then it sends $\langle \text{FINALIZE}, id, result \rangle$ with the received $result$ and waits until it receives $f + 1$ CONFIRM responses in the same view before returning $result$ to the application protocol.

IR's view change protocol is similar to VR's. Each view change is coordinated by a leader, which is unique per view and deterministically chosen. There are three key differences. First, in IR the leader *merges* records during a view change rather than simply taking the longest log from the latest view. The reason for this is that, with inconsistent replicas and unordered operations, any single record could be incomplete. Second, in VR, the leader is used to process operations in the normal case, but IR uses the leader *only* for performing view changes. Finally, on recovery, an IR replica performs a view change, rather than simply interrogating a single replica. This makes sure that the recovering replica either receives all operations it might have sent a reply for, or prevents them from completing.

The full view change protocol follows:

- (1) A replica that notices the need for a view change advances its view number and sets its status to either VIEW-CHANGING or RECOVERING—if the replica just started a recovery. A replica notices the need for a view change either based on a timeout, because it is a recovering replica, or because it received a DO-VIEW-CHANGE message for a view with a larger number than its own current view-number. It records the new view number to disk.
- (2) The replica then sends a $\langle \text{DO-VIEW-CHANGE}, rec, v, v' \rangle$ message to the new leader, *except when the sending replica is a recovering replica*. It also sends the same message, without the rec field, to the other replicas. Here, v identifies the new view, v' is the latest view in which the replica's status was NORMAL, and rec is its unordered record of executed operations.
- (3) Once the new leader receives f records from f other replicas, it considers all records with the highest value of v' . It uses a merge function, shown in Figure 5, to join these into a master record R .
- (4) The leader updates its view number to v_{new} , where v_{new} is the view number from the received messages, and its status to NORMAL. It then informs the other replicas of the completion of the view change by sending a $\langle \text{START-VIEW}, v_{new}, R \rangle$, where R is the master record.
- (5) When a replica receives a START-VIEW message with v_{new} greater than or equal to its current view number, it replaces its own record with R and upcalls into the application protocol with Sync.
- (6) Once Sync is complete, the replica updates its current view number to v_{new} , records this to disk, and enters the NORMAL state.

```

IR-MERGE-RECORDS(records)
1  R, d, u = ∅
2  for ∀op ∈ records
3    if op.type == inconsistent
4      R = R ∪ op
5    elseif op.type == consensus and op.status == FINALIZED
6      R = R ∪ op
7    elseif op.type == consensus and op.status == TENTATIVE
8      if op.result in more than  $\frac{f}{2} + 1$  records
9        d = d ∪ op
10     else
11       u = u ∪ op
12  Sync(R)
13  return R ∪ Merge(d, u)

```

Fig. 5. Merge function for the master record. The IR leader merges records from all replicas in the latest view, which is always a strict superset of the records from replicas in lower views.

Merging Records. The IR-MERGE-RECORDS function is used by the new leader to merge the set of received records. This function is shown in Figure 5. IR-MERGE-RECORDS starts by adding all **inconsistent** operations and **consensus** operations marked **FINALIZED** to R and calling **Sync** into the application protocol. These operations must persist in the next view, so we first apply them to the leader, ensuring that they are visible to any operations for which the leader will decide consensus results next in **Merge**. As an example, **Sync** for the lock server matches up all corresponding **Lock** and **Unlock** by id; if there are unmatched **Locks**, it sets $locked = \text{TRUE}$; otherwise, $locked = \text{FALSE}$.

IR asks the application protocol to decide the consensus result for the remaining **TENTATIVE consensus** operations, which either (1) have a matching result, which we define as the *majority result*, in at least $\lceil \frac{f}{2} \rceil + 1$ records, or (2) do not. IR places these operations in d and u , respectively, and calls **Merge**(d, u) into the application protocol, which must return a consensus result for every operation in d and u .

IR must rely on the application protocol to decide consensus results for several reasons. For operations in d , IR cannot tell whether the operation succeeded with the majority result on the fast path, or whether it took the slow path and the application protocol *decide*d a different result that was later lost. In some cases, it is not safe for IR to keep the majority result, because it would violate application protocol invariants. For example, in the lock server, **OK** could be the majority result if only $\lceil \frac{f}{2} \rceil + 1$ replicas replied **OK**, but the other replicas might have accepted a conflicting lock request. However, it is also possible that the other replicas *did* respond **OK**, in which case **OK** would have been a successful response on the fast-path.

The need to resolve this ambiguity is the reason for the caveat in IR's consensus property (P3) that consensus results can be changed in **Merge**. Fortunately, the application protocol can ensure that successful consensus results *do not change* in **Merge**, simply by maintaining the majority results in d on **Merge** *unless they violate invariants*. The merge function for the lock server, therefore, does not change a majority response of **OK**, *unless* another client holds the lock. In that case, the operation in d could not have returned a successful consensus result to the client (either through the fast or the slow path), so it is safe to change its result.

For operations in u , IR needs to invoke *decide* but cannot without at least $f + 1$ results, so uses **Merge** instead. The application protocol can decide consensus results in **Merge** without $f + 1$ replica results and still preserve IR's visibility property, because IR has already applied all of the operations in R and d , which are the only operations definitely in the operation set, at this point.

The leader adds all operations returned from **Merge** and their consensus results to R , then sends R to the other replicas, which call **Sync**(R) into the application protocol and *replace their own*

records with R . The view change is complete after at least $f + 1$ replicas have exchanged and merged records and SYNC'd with the master record. A replica can only process requests in the new view (in the NORMAL state) *after* it completes the view change protocol. At this point, any recovering replicas can also be considered recovered. If the leader of the view change does not finish the view change by some timeout, then the group will elect a new leader to complete the protocol by starting a new view change with a larger view number.

Periodic Synchronization. Synchronizations need to be performed on view changes. It can also be beneficial to perform them periodically even during normal operation. This reduces the cost of future synchronizations, as the set of operations they need to process is smaller. This compensates for an inherent tradeoff in the unordered operation set used by IR. In systems that use a replicated log, a replica's state can be entirely characterized by the index of the last executed operation; this is not true for IR's operation set. As a result, extra work is needed for Merge and Sync operations, the latter of which must revisit every unsynchronized operation in the operation set. Periodic synchronizations help mitigate this cost.

3.2.3 Client Recovery. We assume that clients can lose some or all of their state on failure. On recovery, a client must ensure that (1) it recovers its latest operation counter, and (2) any operations that it started but did not finish are FINALIZED. To do so, the recovering client requests the *id* for its latest operation from a majority of the replicas. This poll gets the client the largest *id* that the group has seen from it, so the client takes the largest returned *id* and increments it to use as its new operation counter.

A view change finalizes all TENTATIVE operation on the next synchronization, so the client does not need to finish previously started operations and IR does not have to worry about clients failing to recover after failure.

3.3 Correctness

For correctness, we show that IR provides the following properties for operations in the *operation set*:

P1. [Fault tolerance]. At any time, every operation in the operation set is in the record of at least one replica in any quorum of $f + 1$ non-failed replicas.

P2. [Visibility]. For any two **consensus** operations in the operation set, at least one is visible to the other.

P3. [Consensus results]. At any time, every successful consensus result is in the record of at least one replica in any quorum. Again, the only exception being that the application protocol modified the result through Merge.

IR also provides the following *eventual consistency* property, which is not necessary for correctness, but is useful for application protocols. As this is a liveness property, it holds only during periods of synchrony, when messages that are repeatedly resent are eventually delivered before the recipient times out (Fischer et al. 1985):

P4. [Eventual Consistency]. Given a sufficiently long period of synchrony, any operation in the operation set (and its consensus result, if applicable) will eventually have executed or Synced at every non-faulty replica.

In addition to the proof of correctness below, we have also formally specified both IR and TAPIR in the TLA+ language (Lamport 1994), and model-checked its correctness. The TLA+ specification is available in a technical report (Zhang et al. 2015b).

We begin our proof of correctness by defining the following terms:

- D1.* An operation is *applied* at a replica if that replica has executed (through `ExecInconsistent` or `ExecConsensus`) or synchronized (through `Sync`) the operation.
- D2.* An operation X is *visible* to a **consensus** operation Y if one of the replicas providing candidate results for Y has previously applied X .
- D3.* The *persistent operation set* is the set of operations applied at at least one replica in any quorum of $f + 1$ non-failed replicas.

We first prove a number of invariants about the persistent operation set. Given these invariants, we can show that the IR properties hold.

I1. The size of persistent operation set is monotonically increasing.

I1 holds at every replica during normal operation, because replicas never roll back executed operations. I1 also hold across view changes. The leader merges all operations from the records of $f + 1$ non-faulty replicas into the master record, so by quorum intersection, the master record contains every operation in the persistent operation set. Then, at least $f + 1$ non-faulty replicas replace their record with the master record and applies the master record (through `Sync`), so any persistent operation before the view change will continue to persist after the view change.

*I2. All operations in the persistent operation set are visible to any **consensus** operation added to the set.*

consensus operations are added to the persistent set by either (1) executing at at least a quorum of $f + 1$ replicas or (2) being merged by the leader into the master record. In case 1, by definition, every operation already in the persistent operation set must be applied at at least 1 replica out of the quorum and will be visible to the added **consensus** operation. In case 2, the leader applies all operations in the persistent operation set (through `Sync`) before running `Merge`, ensuring that every operation already in the persistent operation set is visible to operations added to the persistent operation set through `Merge`.

*I3. The result of any **consensus** operation in the persistent operation set is either the successful consensus result or the Merge result.*

The result of any **consensus** operations in the persistent set is (1) a matching result from executing the operation (through `ExecConsensus`) at a fast quorum of $\lceil \frac{3}{2}f \rceil + 1$ replicas, (2) a result from executing the application protocol-specific *decide* function in the client-side library, or (3) a result from executing `Merge` at the leader during a view change. In case 1, the matching result will be both the result in the persistent operation set and the successful consensus result. The same holds for the result returned from *decide* in case 2. During a view change, the leader may get an operation that has already fulfilled either case 1 or case 2, and change the result in `Merge`. The result from `Merge` will be in the record and applied to at least $f + 1$ replicas. Thus, either the successful consensus result or, if the application protocol changed the result in `Merge`, the `Merge` result, will continue to persist in the persistent operation set.

I4. All operations and consensus results in the persistent operation set in all previous views must be applied at a replica before it executes any operations in the new view.

IR clients require that all responses come from replicas in the same view. Thus, if any replica is in view v and at least $f + 1$ other replicas are in a higher view $V > v$, that replica cannot successfully complete an operation until it joins the higher view. To join the higher view, the replica in the lower view must obtain the master record from a replica in the higher view, and `Sync` with that master record. The master record contains all operations in the persistent operation set, so the

replica will apply all operations from the persistent operation set before processing operations in the new view.

Given these four invariants for the persistent operation set, we can show that the four properties of IR hold. Any operation in the operation set must have executed at (and received matching responses from) $f + 1$ of $2f + 1$ replicas, so by quorum intersection, all operations in the operation set must be in the persistent operation set. Thus, I1 directly implies P1, as any operation in the persistent operation set will continue to be in the set. I1 and I2 imply P2, because for any **consensus** operation X , all operations added to the persistent operation set before X are visible to X and X will be visible to all operations added to the persistent operation set after it. I1 and I3 implies P3, because either the successful consensus result will remain in the persistent operation set or the Merge result will. I4 implies P4, because if all replicas are non-faulty for long enough, they will eventually all attempt to participate in processing operations, which will cause them to Sync all operations in the persistent operation set.

4 BUILDING ATOP IR

IR obtains performance benefits because it offers weak consistency guarantees and relies on application protocols to resolve inconsistencies, similar to eventual consistency protocols such as Dynamo (DeCandia et al. 2007) and Bayou (Terry et al. 1995). However, IR provides an important benefit over eventual consistency systems, which expect applications to resolve conflicts *after they happen*. In addition to **inconsistent** operations, IR provides **consensus** operations, which allow applications to *detect and resolve conflicts as they happen*. Using **consensus** operations, application protocols can enforce higher-level guarantees (e.g., TAPIR’s linearizable transaction ordering) across replicas despite IR’s weak consistency.

However, building strong guarantees on IR requires careful application protocol design. IR cannot support certain application protocol invariants. Moreover, if misapplied, IR can even provide applications with *worse* performance than a strongly consistent replication protocol. In this section, we discuss the properties that application protocols need to have to correctly and efficiently enforce higher-level guarantees with IR and TAPIR’s techniques for efficiently providing linearizable transactions.

4.1 IR Application Protocol Requirement: *Invariant checks must be performed pairwise*

Application protocols can enforce certain types of invariants with IR, but not others. IR guarantees that in any pair of **consensus** operations, at least one will be visible to the other (P2). Thus, IR readily supports invariants that can be safely checked by examining *pairs* of operations for conflicts. For example, our lock server example can enforce mutual exclusion. However, application protocols cannot check invariants that require the entire history, because each IR replica may have an incomplete history of operations. For example, tracking bank account balances and allowing withdrawals only if the balance remains positive is problematic, because the invariant check must consider the entire history of deposits and withdrawals.

Despite this seemingly restrictive limitation, application protocols can still use IR to enforce useful invariants, including lock-based concurrency control, like Strict Two-Phase Locking (S2PL). As a result, distributed transaction protocols like Spanner (Corbett et al. 2012) or Replicated Commit (Mahmoud et al. 2013) would work with IR. IR can also support optimistic concurrency control (OCC) (Kung and Robinson 1981), because OCC checks are pairwise as well: each committing transaction is checked against every previously committed transaction, so **consensus** operations suffice to ensure that *at least one replica sees any conflicting transaction* and aborts the transaction being checked.

4.2 IR Application Protocol Requirement: *Application protocols must be able to change consensus operation results*

Inconsistent replicas could execute **consensus** operations with one result and later find the group agreed to a different consensus result. For example, if the group in our lock server agrees to reject a Lock operation that one replica accepted, the replica must later free the lock, and vice versa. As noted in the section on pairwise invariant checks, the group as a whole continues to enforce mutual exclusion, so these temporary inconsistencies are tolerable and are always resolved by the end of synchronization.

In TAPIR, we take the same approach with distributed transaction protocols. 2PC-based protocols are always prepared to abort transactions, so they can easily accommodate a Prepare result changing from PREPARE-OK to ABORT. If ABORT changes to PREPARE-OK, then it might temporarily cause a conflict at the replica, which can be correctly resolved, because the group as a whole could not have agreed to PREPARE-OK for two conflicting transactions.

Note that IR replicas are not directly informed (via an upcall) if their consensus result changes as a result of finalization. We use a simple strategy in TAPIR to ensure that replicas, whether they were part of the Prepare quorum or not, learn the transaction outcome. The Prepare message is followed by a Commit or Abort operation, executed as an IR **inconsistent** operation. This idiom, in which a client follows a **consensus** operation with an **inconsistent** operation indicating the outcome, is generally useful with IR.

Changing Prepare results does sometimes cause unnecessary aborts. To reduce these, TAPIR introduces two Prepare results in addition to PREPARE-OK and ABORT: ABSTAIN and RETRY. ABSTAIN helps TAPIR distinguish between conflicts with *committed* transactions, which will not abort, and conflicts with *prepared* transactions, which may later abort. Replicas return RETRY if the transaction has a chance of committing later. The client can retry the Prepare *without* re-executing the transaction.

4.3 IR Performance Principle: *Application protocols should not expect operations to execute in the same order*

To efficiently achieve agreement on consensus results, application protocols should not rely on operation ordering for application ordering. For example, many transaction protocols (Gray and Lamport 2006; Baker et al. 2011; Kraska et al. 2013) use Paxos operation ordering to determine transaction ordering. They would perform worse with IR, because replicas are unlikely to agree on which transaction should be next in the transaction ordering.

In TAPIR, we use *optimistic timestamp ordering* to ensure that replicas agree on a single transaction ordering despite executing operations in different orders. Like Spanner (Corbett et al. 2012), every committed transaction has a timestamp, and committed transaction timestamps reflect a linearizable ordering. However, TAPIR clients, not servers, propose a timestamp for their transaction; thus, if TAPIR replicas agree to commit a transaction, they have all agreed to the same transaction ordering.

TAPIR replicas use these timestamps to order their transaction logs and multi-versioned stores. Therefore, replicas can execute Commit in different orders but still converge to the same application state. TAPIR leverages loosely synchronized clocks at the clients for picking transaction timestamps, which improves performance but is not necessary for correctness.

4.4 IR Performance Principle: *Application protocols should use cheaper inconsistent operations whenever possible rather than consensus operations*

By concentrating invariant checks in a few operations, application protocols can reduce **consensus** operations and improve their performance. For example, in a transaction protocol, any

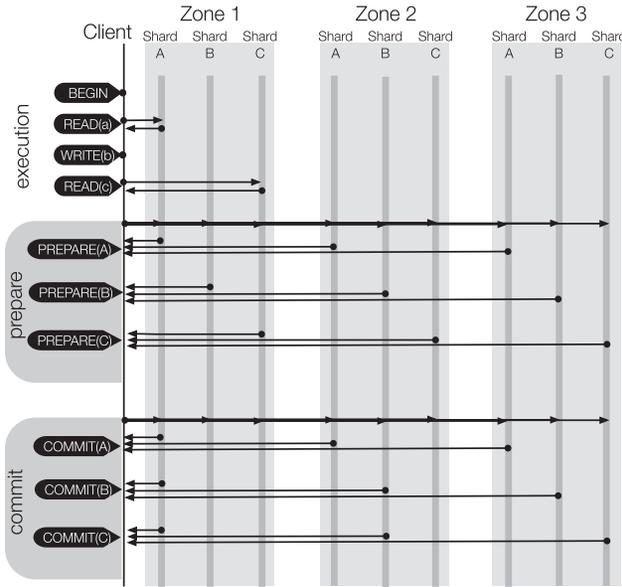


Fig. 6. *Example read-write transaction in TAPIR.* TAPIR executes the same transaction pictured in Figure 2 with less redundant coordination. Reads go to the closest replica and `Prepare` takes a single round-trip to all replicas in all shards. TAPIR sends `Commit` asynchronously, because participants will always come to the same commit decision (see Section 5.2.3). Overall, TAPIR is able to commit transactions in a single round-trip to all participating replicas.

operation that decides transaction ordering must be a **consensus** operation to ensure that replicas agree to the same transaction ordering. For locking-based transaction protocols, this is any operation that acquires a lock. Thus, every `Read` and `Write` must be replicated as a **consensus** operation.

TAPIR improves on this by using optimistic transaction ordering and OCC, which reduces **consensus** operations by concentrating all ordering decisions into a single set of validation checks at the proposed transaction timestamp. These checks execute in `Prepare`, which is TAPIR's only **consensus** operation. `Commit` and `Abort` are **inconsistent** operations, while `Read` and `Write` are not replicated.

5 TAPIR

This section details TAPIR. As noted, TAPIR is designed to efficiently leverage IR's weak guarantees to provide high-performance linearizable transactions. Using IR, TAPIR can order a transaction in a *single round-trip* to all replicas in all participant shards without *any centralized coordination*.

TAPIR is designed to be layered atop IR in a replicated, transactional storage system. Together, TAPIR and IR eliminate the redundancy in the replicated transactional system, as shown in Figure 2. As a comparison, Figure 6 shows the coordination required for the same read-write transaction in TAPIR with the following benefits: (1) TAPIR does not have any leaders or centralized coordination; (2) TAPIR Reads always go to the closest replica; and (3) TAPIR `Commit` takes a single round-trip to the participants in the common case.

5.1 Overview

TAPIR is designed to provide distributed transactions for a scalable storage architecture. This architecture partitions data into shards and replicates each shard across a set of storage servers for

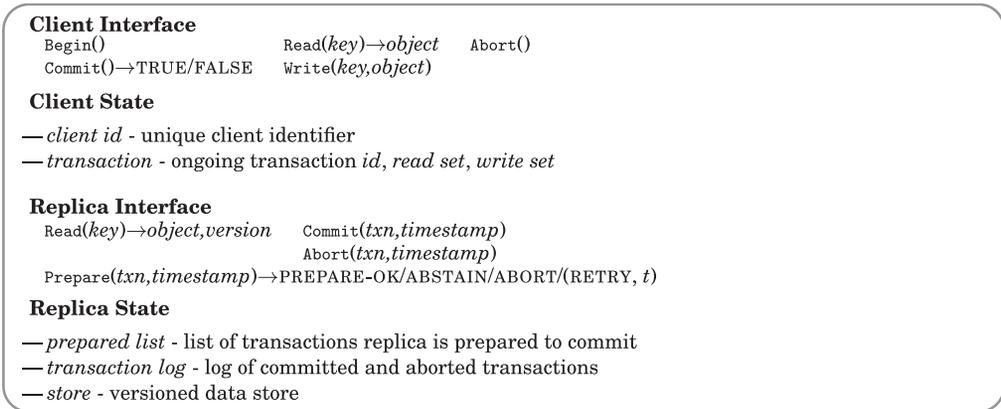


Fig. 7. Summary of TAPIR interfaces and client and replica state.

availability and fault tolerance. Clients are front-end application servers, located in the same or another datacenter as the storage servers, not end-hosts or user machines. They have access to a directory of storage servers using a service like Chubby (Burrows 2006) or ZooKeeper (Hunt et al. 2010) and directly map data to servers using a technique like consistent hashing (Karger et al. 1997).

TAPIR provides a general storage and transaction interface for applications via a client-side library. Note that TAPIR is the application protocol for IR; applications using TAPIR do not interact with IR directly.

A TAPIR application Begins a transaction, then executes Reads and Writes during the transaction’s *execution period*. During this period, the application can Abort the transaction. Once it finishes execution, the application Commits the transaction. Once the application calls Commit, it can no longer abort the transaction. The 2PC protocol will run to completion, committing or aborting the transaction based entirely on the decision of the participants. Since TAPIR’s 2PC coordinators are not allowed to alter the decision taken by the participants (e.g., decide to abort the transaction even if all participants decided to commit), they do not have to be fault-tolerant. This property allows TAPIR to use clients as 2PC coordinators, as in MDCC (Kraska et al. 2013), to reduce the number of round-trips to storage servers.

TAPIR provides the traditional ACID guarantees with the strictest level of isolation: strict serializability (or linearizability) of committed transactions.

5.2 Protocol

TAPIR provides transaction guarantees using a *transaction processing protocol*, *IR functions*, and a *coordinator recovery protocol*.

Figure 7 shows TAPIR’s interfaces and state at clients and replicas. Replicas keep committed and aborted transactions in a *transaction log* in timestamp order; they also maintain a multi-versioned *data store*, where each version of an object is identified by the timestamp of the transaction that wrote the version. TAPIR replicas serve reads from the versioned data store and maintain the transaction log for synchronization and checkpointing. Like other 2PC-based protocols, each TAPIR replica also maintains a *prepared list* of transactions that it has agreed to commit.

Each TAPIR client supports one ongoing transaction at a time. In addition to its *client id*, the client stores the state for the ongoing *transaction*, including the *transaction id* and *read and write sets*. The transaction id must be unique, so the client uses a tuple of its client id and *transaction*

```

TAPIR-EXEC-CONSENSUS(op)
1  txn = op.args.txn
2  timestamp = op.args.timestamp
3  if txn.id ∈ txn-log
4      if txn-log[txn.id].status == COMMITTED
5          return PREPARE-OK
6      else
7          return ABORT
8      elseif txn.id ∈ prepared-list
9          return PREPARE-OK
10     else
11         return TAPIR-OCC-CHECK(txn, timestamp)

```

Fig. 8. Since Prepare is TAPIR’s only **consensus** operations, TAPIR-EXEC-CONSENSUS simply runs TAPIR’s prepare algorithm at replicas.

counter, similar to IR. TAPIR does not require synchronous disk writes at the client or the replicas, as clients do not have to be fault-tolerant and replicas use IR.

5.2.1 Transaction Processing. We begin with TAPIR’s protocol for executing transactions.

- (1) For Write(*key*, *object*), the client buffers *key* and *object* in the write set until commit and returns immediately.
- (2) For Read(*key*), if *key* is in the transaction’s write set, the client returns *object* from the write set. If the transaction has already read *key*, then it returns a cached copy. Otherwise, the client sends Read(*key*) to the replica.
- (3) On receiving Read, the replica returns *object* and *version*, where *object* is the latest committed version of *key* and *version* is the timestamp of the transaction that wrote that version.
- (4) On response, the client puts (*key*, *version*) into the transaction’s read set and returns *object* to the application.

Once the application calls Commit or Abort, the execution phase finishes. To commit, the TAPIR client coordinates across all *participants*—the shards that are responsible for the keys in the read or write set—to find a single timestamp, consistent with the strict serial order of transactions, to assign the transaction’s reads and writes, as follows:

- (1) The TAPIR client selects a *proposed timestamp*. Proposed timestamps must be unique, so clients use a tuple of their local time and their *client id*.
- (2) The TAPIR client invokes Prepare(*txn*, *timestamp*) on all participants as an IR **consensus** operation, where *timestamp* is the proposed timestamp and *txn* includes the transaction id (*txn.id*) and the transaction read (*txn.read_set*) and write sets (*txn.write_set*).
- (3) Each TAPIR replica that receives Prepare (invoked by IR through ExecConsensus) first checks its transaction log for *txn.id*. If found, then it returns PREPARE-OK if the transaction committed or ABORT if the transaction aborted.
- (4) Otherwise, the replica checks if *txn.id* is already in its *prepared list*. If found, then it returns PREPARE-OK.
- (5) Otherwise, the replica runs TAPIR’s *OCC validation checks*, which check for conflicts with the transaction’s read and write sets at *timestamp*, shown in Figure 9.
- (6) Once the TAPIR client receives results from all shards, the client sends Commit(*txn*, *timestamp*) if all shards replied PREPARE-OK or Abort(*txn*, *timestamp*) if any shards replied ABORT or ABSTAIN. If any shards replied RETRY, then the client retries with a new proposed timestamp (up to a set limit of retries).

```

TAPIR-OCC-CHECK(txn, timestamp)
1  for  $\forall key, version \in txn.read\text{-}set$ 
2    if version < store[key].latestVersion
3      return ABORT
4    elseif timestamp > MIN(prepared-writes[key])
5      return ABSTAIN
6  for  $\forall key \in txn.write\text{-}set$ 
7    if timestamp < MAX(PREPARED-READS(key))
8      return RETRY, MAX(PREPARED-READS(key))
9    elseif timestamp < store[key].latestVersion
10     return RETRY, store[key].latestVersion
11  prepared-list[txn.id] = timestamp
12  return PREPARE-OK

```

Fig. 9. Validation function for checking for OCC conflicts on Prepare. PREPARED-READS and PREPARED-WRITES get the proposed timestamps for all transactions that the replica has prepared and read or write to *key*, respectively.

```

TAPIR-DECIDE(results)
1  if ABORT  $\in results$ 
2    return ABORT
3  if count(PREPARE-OK, results)  $\geq f + 1$ 
4    return PREPARE-OK
5  if count(ABSTAIN, results)  $\geq f + 1$ 
6    return ABORT
7  if RETRY  $\in results$ 
8    return RETRY, max(results.retry-timestamp)
9  return ABORT

```

Fig. 10. TAPIR's *decide* function. IR runs this if replicas return different results on Prepare.

- (7) On receiving a Commit, the TAPIR replica: (1) commits the transaction to its transaction log, (2) updates its versioned store with *w*, (3) removes the transaction from its prepared list (if it is there), and (4) responds to the client.
- (8) On receiving a Abort, the TAPIR replica: (1) logs the abort, (2) removes the transaction from its prepared list (if it is there), and (3) responds to the client.

Like other 2PC-based protocols, TAPIR can return the outcome of the transaction to the application as soon as Prepare returns from all shards (in Step 6) and send the Commit operations asynchronously. As a result, using IR, TAPIR can commit a transaction with a *single round-trip* to all replicas in all shards.

5.2.2 IR Support. Because TAPIR's Prepare is an IR **consensus** operation, TAPIR must implement a client-side *decide* function, shown in Figure 10, which merges inconsistent Prepare results from replicas in a shard into a single result. TAPIR-DECIDE is simple: if a majority of the replicas (i.e., at least $f + 1$ out of the total of $2f + 1$ replicas) replied PREPARE-OK, then it commits the transaction. This is safe because no conflicting transaction could also get a majority of the replicas to return PREPARE-OK.

TAPIR also supports Merge, shown in Figure 11, and Sync, shown in Figure 12, at replicas. TAPIR-MERGE first removes any prepared transactions from the leader where the Prepare operation is TENTATIVE (i.e., all operations passed as arguments by the IR merge procedure). This step removes any inconsistencies that the leader may have, because it executed a Prepare differently—out-of-order or missed—than the rest of the group.

```

TAPIR-MERGE( $d, u$ )
1  for  $\forall op \in d \cup u$ 
2     $txn = op.args.txn$ 
3    if  $txn.id \in prepared-list$ 
4      DELETE( $prepared-list, txn.id$ )
5  for  $op \in d$ 
6     $txn = op.args.txn$ 
7     $timestamp = op.args.timestamp$ 
8    if  $txn.id \notin txn-log$  and  $op.result == PREPARE-OK$ 
9       $R[op].result = TAPIR-OCC-CHECK(txn, timestamp)$ 
10   else
11      $R[op].result = op.result$ 
12  for  $op \in u$ 
13     $txn = op.args.txn$ 
14     $timestamp = op.args.timestamp$ 
15     $R[op].result = TAPIR-OCC-CHECK(txn, timestamp)$ 
16  return  $R$ 

```

Fig. 11. TAPIR's merge function. IR runs this function at the leader on synchronization and recovery.

```

TAPIR-SYNC( $R$ )
1  for  $\forall op \in R$ 
2    if  $op \notin r$  or  $op.result \neq r[op].result$ 
3       $txn = op.args.txn$ 
4       $timestamp = op.args.timestamp$ 
5      if  $op.func == Prepare$ 
6        if  $op.result == PREPARE-OK$ 
7          if  $txn.id \notin prepared-list$  and  $txn.id \notin txn-log$ 
8             $prepared-list[txn.id] = timestamp$ 
9          elseif  $txn.id \in prepared-list$ 
10             DELETE( $prepared-list, txn.id$ )
11         else
12            $txn-log[txn.id].txn = txn$ 
13            $txn-log[txn.id].timestamp = timestamp$ 
14           if  $op.func == Commit$ 
15              $txn-log[txn.id].status = COMMITTED$ 
16           else
17              $txn-log[txn.id].status = ABORTED$ 
18           if  $txn.id \in prepared-list$ 
19             DELETE( $prepared-list, txn.id$ )

```

Fig. 12. TAPIR's function for synchronizing inconsistent replica state. IR runs this on each replica except the leader during synchronization. r is the replica's local record.

The next step checks d for any PREPARE-OK results that might have succeeded on the IR fast path and need to be preserved. If the transaction has not committed or aborted already, then we re-run TAPIR-OCC-CHECK to check for conflicts with other previously prepared or committed transactions. If the transaction *conflicts*, then we know that its PREPARE-OK did not succeed at a fast quorum and is yet to be decided through slow quorum agreement, so we can safely propose an ABORT; otherwise, for correctness, we must preserve the PREPARE-OK, because TAPIR may have moved on to the commit phase of 2PC. Further, we know that it is safe to preserve these PREPARE-OK results, because if they conflicted with another transaction, then the conflicting transaction *must* have gotten its consensus result on the IR slow path, so if TAPIR-OCC-CHECK did not find a conflict, then the conflicting transaction's Prepare must not have succeeded.

Finally, for the operations in u , we simply decide a result for each operation and preserve it. We know that the leader is now consistent with $f + 1$ replicas (which stopped processing new

requests since the synchronization started), so it can make decisions on consensus result for the majority.

TAPIR's sync function, shown in Figure 12, runs at the other replicas to reconcile TAPIR state with the master records, correcting missed operations or consensus results where the replica did not agree with the group. It simply applies operations and consensus results to the replica's state: it logs aborts and commits, and it prepares uncommitted transactions where the group responded PREPARE-OK.

5.2.3 Coordinator Recovery. If a client fails while in the process of committing a transaction, then TAPIR ensures that the transaction runs to completion (either commits or aborts). Further, the client may have returned the commit or abort to the application, so we must ensure that the client's commit decision is preserved. For this purpose, TAPIR uses the *cooperative termination protocol* defined by Bernstein (Bernstein et al. 1987) for coordinator recovery and used by MDCC (Kraska et al. 2013). TAPIR designates one of the participant shards as a *backup shard*, the replicas in which can serve as a backup coordinator if the client fails. As observed by MDCC, because coordinators cannot unilaterally abort transactions (i.e., if a client receives $f + 1$ PREPARE-OK responses from each participant, it must commit the transaction), a backup coordinator can safely complete the protocol without blocking. However, we must ensure that no two coordinators for a transaction are active at the same time.

Coordinator Changes. We use a coordinator change protocol, similar to IR's view change protocol to ensure that only one coordinator is active at a time.² For each transaction, we designate one of the participant shards as a *backup shard*. The initial coordinator for every transaction is the client. In every subsequent view, the currently active backup coordinator is a replica from the backup shard.

For every transaction in its *prepared – list*, each TAPIR replica keeps the transaction's backup shard and a current *coordinator view*. Replicas only process and respond to Prepare, Commit and Abort operations from the active coordinator designated by the current view, identified by indexing into the list of backup shard replicas with the coordinator view number. Replicas also keep a *no-vote list* with transactions that the replica knows a backup coordinator may abort.

If the current coordinator is suspected to have failed, then any of the participants can initiate a coordinator change. In doing so, it keeps the client or any previous backup coordinator from sending operations to the participating replicas. The new coordinator can then poll the participant using Prepare, and make a commit decision without interference from other coordinators. The election protocol for a new backup coordinator progresses as follows:

- (1) Any replica in any participant shard calls `CoordinatorChange` through IR as a **consensus** operation on the backup shard.
- (2) Each replica that executes `CoordinatorChange` through IR, increments and returns its current coordinator view number v . If the replica is not already in the `COORDINATOR-VIEW-CHANGE` state, then it sets its state to `COORDINATOR-VIEW-CHANGE` and stops responding to operations for that transaction.
- (3) The *decide* function for `CoordinatorChange` returns the highest v returned by the replicas.
- (4) Once `CoordinatorChange` returns successfully, the replica sends `StartCoordinatorView(v_{new})`, where v_{new} is the returned view number from `CoordinatorChange`, as an IR **inconsistent** operation to *all* participant shards, including its own.

²Other possible ways to achieve this goal include logging the currently active backup coordinator to a service like Chubby (Burrows 2006) or ZooKeeper (Hunt et al. 2010), or giving each backup coordinator a lease in turn.

```

TAPIR-RECOVERY-DECIDE(results)
1  if ABORT  $\in$  results
2      return ABORT if  $\text{count}(\text{NO-VOTE}, \text{results}) \geq f + 1$ 
3      return ABORT
4  if  $\text{count}(\text{PREPARE-OK}, \text{results}) \geq f + 1$ 
5      return PREPARE-OK
6  return RETRY

```

Fig. 13. TAPIR’s *decide* function for *Prepare* on coordinator recovery. IR runs this if replicas return different results on *Prepare*. This *decide* function differs from the normal case execution *decide*, because it is not safe to return ABORT unless it is sure the original coordinator did not receive PREPARE-OK.

- (5) Any replica that receives `StartCoordinatorView` checks if v_{new} is higher or equal to its current view. If so, then the replica updates its current view number and begins accepting *Prepare*, *Commit* and *Abort* from the active backup coordinator designated by the new view. If the replica is in the backup shard, then it can set its state back to `NORMAL`.
- (6) When a replica executes `StartCoordinatorView` for the view where it is the designated backup coordinator, it begins the cooperative termination protocol.

The *Merge* function for `CoordinatorChange` preserves the consensus result if it is greater than or equal to the current view number at the leader during synchronization. The *Sync* function for `CoordinatorChange` sets the replica state to `COORDINATOR-VIEW-CHANGE` if the consensus result is larger than the replica’s current view number. The *Sync* function for `StartCoordinatorView` just executes the function: it updates the replica’s current view number if v_{new} is greater than or equal to it and sets the state back to `NORMAL` if the replica is in the backup shard.

Cooperative Termination. The backup coordination protocol executed by the active coordinator is based on the cooperative termination protocol described by Bernstein (Bernstein et al. 1987), with several changes to accommodate IR and TAPIR. The most notable changes are that the backup coordinators do not propose timestamps. If the client successfully prepared the transaction at a timestamp t (i.e., achieved at least $f + 1$ `PREPARE-OK` in every participant shard), then the transaction will commit at t . Otherwise, the backup coordinator will eventually abort the transaction.

Next, in Bernstein’s algorithm, any single participant can abort the transaction if they have not yet voted (i.e., replied to a coordinator). However, with IR, no single replica can abort the transaction without information about the state of the other replicas in the shard. As a result, replicas return a `NO-VOTE` response and add the transaction to their *no-vote-list*. Once a replica adds a transaction to the `NO-VOTE-LIST`, it will always respond `NO-VOTE` to *Prepare* operations. Eventually, all replicas in the shard will either converge to a response (i.e., `PREPARE-OK`, `ABORT`) to the original coordinator’s *Prepare* or to a `NO-VOTE` response. TAPIR’s modified cooperative termination protocol proceeds as follows:

- (1) The backup coordinator polls the participants with *Prepare* with no proposed timestamp by invoking *Prepare* as a **consensus** operation in IR with the *decide* function outlined in Figure 13.
- (2) Any replica that receives *Prepare* with no propose timestamp, returns `PREPARE-OK` if it has committed or prepared the transaction, `ABORT` if it has received an `Abort` for the transaction or committed a conflicting transaction and `NO-VOTE` if it does not have the transaction in its *prepared-list* or *txn-log*. If the replica returns `NO-VOTE`, then it adds the transaction to its *no-vote-list*.

```

TAPIR-MERGE( $d, u$ )
1  for  $\forall op \in d \cup u$ 
2     $txn = op.args.txn$ 
3    if  $txn.id \in prepared-list$ 
4      DELETE( $prepared-list, txn.id$ )
5    if  $txn.id \in no-vote-list$ 
6      DELETE( $no-vote-list, txn.id$ )
7  for  $op \in d$ 
8     $txn = op.args.txn$ 
9     $timestamp = op.args.timestamp$ 
10   if  $txn.id \in no-vote-list$ 
11      $R[op].result = NO-VOTE$ 
12   elseif  $txn.id \notin txn-log$  and  $op.result == PREPARE-OK$ 
13      $R[op].result = TAPIR-OCC-CHECK(txn, timestamp)$ 
14   else
15      $R[op].result = op.result$ 
16  for  $op \in u$ 
17     $txn = op.args.txn$ 
18    if  $txn.id \in no-vote-list$ 
19       $R[op].result = NO-VOTE$ 
20    else
21       $timestamp = op.args.timestamp$ 
22       $R[op].result = TAPIR-OCC-CHECK(txn, timestamp)$ 
23  return  $R$ 

```

Fig. 14. TAPIR's merge function. IR runs this function at the leader on synchronization and recovery. This version handles NO-VOTE results.

- (3) The coordinator continues to send Prepare as an IR operation until it either receives a ABORT or PREPARE-OK from all participant shards. Note that the result will be ABORT if a majority of replicas respond NO-VOTE.
- (4) If all participant shards return PREPARE-OK, then the coordinator sends Commit; otherwise, it sends Abort.

Assuming at least $f + 1$ replicas are up in each participant shard and shards are able to communicate, this process will eventually terminate with a backup coordinator sending Commit or Abort to all participants.

We must also incorporate the NO-VOTE into our Merge and Sync handlers for Prepare. We make the following changes to Merge for the final function shown in Figure 14: (lines 5 and 6) delete any tentative NO-VOTES from the *no-vote-list* at the leader for consistency, (lines 10 and 11) return NO-VOTE without running TAPIR-OCC-CHECK if the transaction is already in the *no-vote-list*, because any result to the original Prepare could not have succeeded (i.e., was not returned to the coordinator that invoked the Prepare operation), (lines 18 and 19) do the same for operations without *majority result* (i.e., the result does not match with the results of at least $\lceil \frac{f}{2} \rceil + 1$ of the merged records), where the original coordinator's Prepare definitely did not succeed. If the consensus result to the Prepare is NO-VOTE in Sync, then we add transactions to the *no-vote-list* and remove it from the *prepared-list*, as shown in lines 11 and 12 of Figure 15.

5.3 Correctness

To prove correctness, we show that TAPIR maintains the following properties³ given up to f failures in each replica group and any number of client failures:

³We do not prove database consistency, as it depends on application invariants; however, strict serializability is sufficient to enforce consistency.

```

TAPIR-SYNC( $R$ )
1  for  $\forall op \in R$ 
2    if  $op \notin r$  or  $op.result \neq r[op].result$ 
3       $txn = op.args.txn$ 
4       $timestamp = op.args.timestamp$ 
5      if  $op.func == Prepare$ 
6        if  $op.result == PREPARE-OK$ 
7          if  $txn.id \notin prepared-list$  and  $txn.id \notin txn-log$ 
8             $prepared-list[txn.id] = timestamp$ 
9          elseif  $txn.id \in prepared-list$ 
10              $DELETE(prepared-list, txn.id)$ 
11            if  $op.result == NO-VOTE$  and  $txn.id \notin txn-log$ 
12               $no-vote-list[txn.id] = timestamp$ 
13          else
14             $txn-log[txn.id].txn = txn$ 
15             $txn-log[txn.id].timestamp = timestamp$ 
16            if  $op.func == Commit$ 
17               $txn-log[txn.id].status = COMMITTED$ 
18            else
19               $txn-log[txn.id].status = ABORTED$ 
20            if  $txn.id \in prepared-list$ 
21               $DELETE(prepared-list, txn.id)$ 

```

Fig. 15. TAPIR's function for synchronizing inconsistent replica state. IR runs this on each replica except the leader during synchronization. r is the replica's local record.

- **Isolation.** There exists a global linearizable ordering of committed transactions.
- **Atomicity.** If a transaction commits at any participating shard, then it commits at them all.
- **Durability.** Committed transactions stay committed, maintaining the original linearizable order.

A model-checked TLA+ specification of TAPIR and its correctness properties is available in a technical report (Zhang et al. 2015b).

5.3.1 Isolation. For correctness, we must show that any two conflicting transactions, A and B , that violate the linearizable transaction ordering cannot both commit. If A and B have a conflict, then there must be at least one common shard that is participating in both A and B . We show that, in that shard, $Prepare(A)$ and $Prepare(B)$ cannot both return $PREPARE-OK$, so one transaction must abort.

In the common shard, IR's visibility property (P2) guarantees that $Prepare(A)$ must be *visible* to $Prepare(B)$ (i.e., executes first at one replica out of every $f + 1$ quorum) *or* $Prepare(B)$ is visible to $Prepare(A)$. Without loss of generality, suppose that $Prepare(A)$ is visible to $Prepare(B)$ and the group returns $PREPARE-OK$ to $Prepare(A)$. Any replica that executes TAPIR-OCC-CHECK for both A and B will not return $PREPARE-OK$ for both, so at least one replica out of any $f + 1$ quorum will not return $PREPARE-OK$ to $Prepare(B)$. IR will not get a fast quorum of matching $PREPARE-OK$ results for $Prepare(B)$, and TAPIR's *decide* function will not return $PREPARE-OK$, because it will never get the $f + 1$ matching $PREPARE-OK$ results that it needs. Thus, IR will never return a consensus result of $PREPARE-OK$ for $Prepare(B)$. The same holds if $Prepare(B)$ is visible to $Prepare(A)$ and the group returns $PREPARE-OK$ to $Prepare(B)$. Thus, IR will never return a successful consensus result of $PREPARE-OK$ to executing both $Prepare(A)$ and $Prepare(B)$ in the common participant shard and TAPIR will not be able to commit both transactions.

Further, once decided, the successful consensus results for $Prepare(A)$ and $Prepare(B)$ will persist in the record of at least one replica out of every quorum, unless it has been modified by the application through Merge. TAPIR will never change another result to a $PREPARE-OK$, so the shard

will never respond `PREPARE-OK` to both transactions. IR will ensure that the successful consensus result is eventually Sync'd at all replicas. Once a TAPIR replica prepared a transaction, it will continue to return `PREPARE-OK` until it receives a `Commit` or `Abort` for the transaction. As a result, if the shard returned `PREPARE-OK` as a successful consensus result to `Prepare(A)`, then it will never allow `Prepare(B)` to also return `PREPARE-OK` (unless *A* aborts), ensuring that *B* is never able to commit. The opposite also holds true.

5.3.2 Atomicity. If a transaction commits at any participating shard, then the TAPIR client must have received a successful `PREPARE-OK` from every participating shard on `Prepare`. Barring failures, it will ensure that `Commit` eventually executes successfully at every participant. TAPIR replicas always execute `Commit`, even if they did not prepare the transaction, so `Commit` will eventually commit the transaction at every participant if it executes at one participant.

If the coordinator fails, then at least one replica in a participant shard will detect the failure and initiate the coordinator recovery protocol. Assuming no more than *f* simultaneous failures in the backup shard, the coordinator change protocol will eventually pick a new active backup coordinator from the backup shard. At this point, the participants will have stopped processing operations from the client and any previous backup coordinators.

Backup coordinators do not propose timestamps, so if any replica in a participant shard received a `Commit`, then the client's `Prepare` must have made it into the operation set of every participant shard with `PREPARE-OK` as the consensus result. IR's consensus result and eventual consistency properties (P3 and P4) ensure that the `PREPARE-OK` will eventually be applied at all replicas in every participant shard and TAPIR ensures that successful `PREPARE-OK` results are not changed in `Merge` (as shown above). Once a TAPIR replica applies `PREPARE-OK`, it will continue to return `PREPARE-OK`, so once replicas in participant groups have stopped processing operations from previous coordinators, all non-failed replicas in all shards will eventually return `PREPARE-OK`. As a result, the backup coordinator must eventually receive `PREPARE-OK` as well from all participants.

In the meantime, the backup coordinator is guaranteed to not receive an `ABORT` from a participant shard. A participant shard will only return an `ABORT` if: (1) a conflicting transaction committed, (2) a majority of the replicas return `NO-VOTE`, because they did not have a record of the transaction, or (3) the transaction was aborted on the shard. Case (1) is not possible, because the conflicting transaction could not have also received a successful consensus result of `PREPARE-OK` (based on our isolation proof) and IR's consensus result property (P3) ensures that the conflicting transaction could never get a `PREPARE-OK` consensus result, so the conflicting transaction cannot commit. Case (2) is not possible, because the client could not have received `PREPARE-OK` as a consensus result if a majority of the replicas do not have the transaction in their *prepared-list* and IR's P3 and P4 ensures the transaction eventually makes its way into the *prepared-list* of every replica. Case (3) is not possible, because the client could not have sent `Abort` if it got `PREPARE-OK` from all participant shards and no previous backup coordinator could have sent `Abort`, because cases (1) and (2) will never happen. As a result, the backup coordinator will not abort the transaction.

5.3.3 Durability. For all committed transactions, either the client or a backup coordinator will eventually execute `Commit` successfully as an IR **inconsistent** operation. IR guarantees that the `Commit` will never be lost (P1) and every replica will eventually execute or synchronize it. On `Commit`, TAPIR replicas use the transaction timestamp included in `Commit` to order the transaction in their log, regardless of when they execute it, thus maintaining the original linearizable ordering. If there are no coordinator failures, then a transaction would eventually be finalized through an IR inconsistent operation (`Commit/Abort`), which ensures that the decision would never be lost. As

described above, for coordinator failures, the coordinator recovery protocol ensures that a backup coordinator would eventually send `Commit` or `Abort` to all participants.

6 TAPIR EXTENSIONS

We now describe four useful extensions to TAPIR.

6.1 Read-Only Transactions

Since it uses a multi-versioned store, TAPIR easily supports globally consistent read-only transactions at a given snapshot timestamp. However, since TAPIR replicas are inconsistent, it is important to ensure that: (1) reads are up-to-date and (2) later write transactions do not invalidate the reads. To achieve these properties, TAPIR replicas keep a read timestamp for each object.

TAPIR's read-only transactions have a single round-trip fast path that sends the `Read` to only one replica. If that replica has a *validated version* of the object—where the write timestamp precedes the snapshot timestamp and the read timestamp follows the snapshot timestamp—then we know that the returned object is valid, because it is up-to-date, and will remain valid, because it will not be overwritten later. If the replica lacks a validated version, then TAPIR uses the slow path and executes a `QuorumRead` through IR as an **inconsistent** operation. Note that this is an instance of a **inconsistent** operation with return values, which we have not used previously. Recall that such an operation returns a result set containing results from at least $f + 1$ replicas.

A `QuorumRead` operation returns the state of an object at the replica. It is performed with respect to a snapshot timestamp. In addition, it updates the read timestamp stored at each replica. As a result, it ensures that at least $f + 1$ replicas will not subsequently accept writes that would invalidate the read.

More precisely, the protocol for read-only transactions follows:

- (1) The TAPIR client chooses a *snapshot timestamp* for the transaction; for example, the client's local time.
- (2) The client sends `Read(key,version)`, where *key* is what the application wants to read and *version* is the snapshot timestamp.
- (3) If the replica has a validated version of the object, then it returns it. Otherwise, it returns a failure.
- (4) If the client could not get the value from the replica, then it executes a `QuorumRead(key,version)` through IR as an inconsistent operation.
- (5) Any replica that receives `QuorumRead` returns the latest version of the object from the data store. It also writes the `Read` to the transaction log and updates the data store to ensure that it will not prepare for transactions that would invalidate the `Read`.
- (6) The client returns the object with the highest timestamp to the application.

As a brief sketch of correctness, it is always safe to read a version of the key that is *validated* at the snapshot timestamp. The version will always be valid at the snapshot timestamp, because the write timestamp for the version is earlier than the snapshot timestamp and the read timestamp is after the snapshot timestamp. If the replica does not have a validated version, then the replicated `QuorumRead` ensures that (1) the client gets the latest version of the object (because at least 1 of any $f + 1$ replicas must have it), and (2) a later write transaction cannot overwrite the version (because at least 1 of the $f + 1$ `QuorumRead` replicas will block it).

Since TAPIR also uses loosely synchronized clocks, it could be combined with Spanner's algorithm for providing externally consistent read-only transactions as well. This combination would require Spanner's TrueTime technology and waits at the client for the TrueTime uncertainty bound. Note that while TAPIR itself provides external consistency for read-write transactions

regardless of clock skew, this read-only protocol would provide linearizability guarantees only if the clock skew did not exceed the TrueTime bound, like Spanner (Corbett et al. 2012).

6.2 Serializability

TAPIR is restricted in its ability to accept transactions out of order, because it provides linearizability, i.e., strict serializability. Thus, TAPIR replicas cannot accept writes that are older than the last write for the same key, and they cannot accept reads of older versions of the same key.

However, if TAPIR's guarantees were weakened to (non-strict) *serializability*, then it can then accept proposed timestamps any time in the past as long as they respect the serializable transaction ordering. This increases the number of updates that can be accepted in highly concurrent workloads. Implementing this optimization requires tracking the timestamp of the transaction that last read and wrote each version.

With this optimization, TAPIR can now accept (1) reads of past versions, as long as the read timestamp precedes the write timestamp of the next version, and (2) writes in the past (per the Thomas Write Rule (Thomas 1979)), as long as the write timestamp follows the read timestamp of the previous version and precedes the write timestamp of the next version.

6.3 Retry Timestamp Selection

The proposed timestamp for a given transaction determines whether the participant replicas will accept the transaction: It will be rejected if a transaction with lower timestamp has already been processed. Normally, clients use their local clocks to select timestamps; assuming that all clients have loosely synchronized clocks, this provides fairness.

Alternatively, a client can increase the likelihood that participant replicas will accept its proposed timestamp by proposing a very large timestamp; this decreases the likelihood that the participant replicas have already accepted a higher timestamp. The downside to this approach is that other clients' transactions will be rejected until their clocks pass that higher timestamp, potentially creating a starvation problem. As a balance, we propose that clients increment their proposed timestamp by an exponentially weighted amount on each retry. This decreases the chances of any individual transaction being forced to retry forever, while balancing against starvation risk for other transactions.

6.4 Tolerating Very High Skew

If there is significant clock skew between servers and clients, then TAPIR can employ a batching protocol of sorts at the participant replicas to decrease the likelihood that transactions will arrive out of timestamp order. On receiving each Prepare message, the participant replica can wait (for the estimated clock error-bound period) to see if any other transactions with smaller timestamps will arrive. After the wait, the replica can process transactions in timestamp order. This wait increases the chances that the participant replica can process transactions in timestamp order and decreases the number of transactions that it will have to reject for arriving out of order.

7 EVALUATION

In this section, our experiments demonstrate the following:

- TAPIR provides better latency *and* throughput than conventional transaction protocols in both the datacenter and wide-area environments.
- TAPIR's abort rate scales similarly to other OCC-based transaction protocols as contention increases.

Table 1. Measured RTTs and Clock Skews between Google Compute Engine VMs

	Latency (ms)			Clock Skew (ms)		
	U.S.	Europe	Asia	U.S.	Europe	Asia
U.S.	1.2	111.3	166.5	3.4	1.3	1.86
Europe	–	0.8	261.8	–	0.1	1.9
Asia	–	–	10.8	–	–	.3

- Clock synchronization sufficient for TAPIR’s needs is widely available in both datacenter and wide-area environments.
- TAPIR provides performance comparable to systems with weak consistency guarantees and no transactions.

7.1 Experimental Setup

We ran our experiments on Google Compute Engine (GCE) with VMs spread across three geographical regions—U.S., Europe, and Asia—and placed in different availability zones within each geographical region. Each server has a virtualized, single core 2.6GHz Intel Xeon, 8GB of RAM, and 1Gb NIC.

7.1.1 Testbed Measurements. As TAPIR’s performance depends on clock synchronization and round-trip times, we first present latency and clock skew measurements of our test environment. As clock skew increases, TAPIR’s latency increases and throughput decreases, because clients may have to retry more `Prepare` operations. It is important to note that TAPIR’s performance depends on the *actual* clock skew, not a worst-case bound like Spanner (Corbett et al. 2012).

We measured the clock skew by sending a ping message with timestamps taken on either end. We calculate skew by comparing the timestamp taken at the destination to the one taken at the source plus half the round-trip time (assuming that network latency is symmetric). Table 1 reports the average skew and latency between the three geographic regions. Within each region, we average over the availability zones. Our VMs benefit from Google’s reliable wide-area network infrastructure; although we use UDP for RPCs over the wide-area, we saw negligible packet drops and little variation in round-trip times.

The average RTT between U.S.-Europe was 110ms; U.S.-Asia was 165ms; Europe-Asia was 260ms. We found the clock skew to be low, averaging between 0.1 and 3.4ms, demonstrating the feasibility of synchronizing clocks in the wide area. However, there was a long tail to the clock skew, with the worst case clock skew being around 27ms—making it significant that TAPIR’s performance depends on actual rather than worst-case clock skew. As our measurements show, the skew in this environment is low enough to achieve good performance.

7.1.2 Implementation. We implemented TAPIR in a transactional key-value storage system, called TAPIR-KV. Our prototype consists of 9,094 lines of C++ code, not including the testing framework.

We also built two comparison systems. The first, OCC-STORE, is a “standard” implementation of 2PC and OCC, combined with an implementation of Multi-Paxos (Lamport 2001). Like TAPIR, OCC-STORE accumulates a read and write set with read versions at the client during execution and then runs 2PC with OCC checks to commit the transaction. OCC-STORE uses a centralized timestamp server to generate transaction timestamps, which we use to version data in the multi-versioned storage system. We verified that this timestamp server was not a bottleneck in our experiments.

Table 2. Transaction Profile for Retwis Workload

Transaction Type	# gets	# puts	workload %
Add User	1	3	5%
Follow/Unfollow	2	2	15%
Post Tweet	3	5	30%
Load Timeline	rand(1,10)	0	50%

Our second system, LOCK-STORE, is based on the Spanner protocol (Corbett et al. 2012). Like Spanner, it uses 2PC with S2PL and Multi-Paxos. The client acquires read locks during execution at the Multi-Paxos leaders and buffers writes. On Prepare, the leader replicates these locks and acquires write locks. We use loosely synchronized clocks at the leaders to pick transaction timestamps, from which the coordinator chooses the largest as the commit timestamp. We use the client as the coordinator, rather than one of the Multi-Paxos leaders in a participant shard, for a more fair comparison with TAPIR-KV. Lacking access to TrueTime, we set the TrueTime error bound to 0, eliminating the need to wait out clock uncertainty and thereby giving the benefit to this protocol.

7.1.3 Workload. We use two workloads for our experiments. We first test using a synthetic workload based on the Retwis application (Leau 2013). Retwis is an open-source Twitter clone designed to use the Redis key-value storage system (Redis 2013). Retwis has a number of Twitter functions (e.g., add user, post tweet, get timeline, follow user) that perform Puts and Gets on Redis. We treat each function as a transaction, and generate a synthetic workload based on the Retwis functions as shown in Table 2.

Our second experimental workload is YCSB+T (Dey et al. 2014), an extension of YCSB (Cooper et al. 2010)—a commonly used benchmark for key-value storage systems. YCSB+T wraps database operations inside simple transactions such as read, insert or read-modify-write. However, we use our Retwis benchmark for many experiments, because it is more sophisticated: transactions are more complex—each touches 2.5 shards on average—and longer—each executes 4–10 operations.

7.2 Single Datacenter Experiments

We begin by presenting TAPIR-KV’s performance within a single datacenter. We deploy TAPIR-KV and the comparison systems over 10 shards, all in the U.S. geographic region, with three replicas for each shard in different availability zones. We populate the systems with 10M keys and make transaction requests with a Zipf distribution (coefficient 0.75) using an increasing number of closed-loop clients.

Figure 16 shows the average latency for a transaction in our Retwis workload at different throughputs. At low offered load, TAPIR-KV has lower latency, because it is able to commit transactions in a single round-trip to all replicas, whereas the other systems need two; its commit latency is thus reduced by 50%. However, Retwis transactions are relatively long, so the difference in *transaction* latency is relatively small.

Compared to the other systems, TAPIR-KV is able to provide roughly 3× the peak throughput, which stems directly from IR’s weak guarantees: it has no leader and does not require cross-replica coordination. Even with moderately high contention (Zipf coefficient 0.75), TAPIR-KV replicas are able to inconsistently execute operations and still agree on ordering for transactions at a high rate.

7.3 Wide-Area Latency

For wide-area experiments, we placed one replica from each shard in each geographic region. For systems with leader-based replication, we fix the leader’s location in the U.S. and move the client

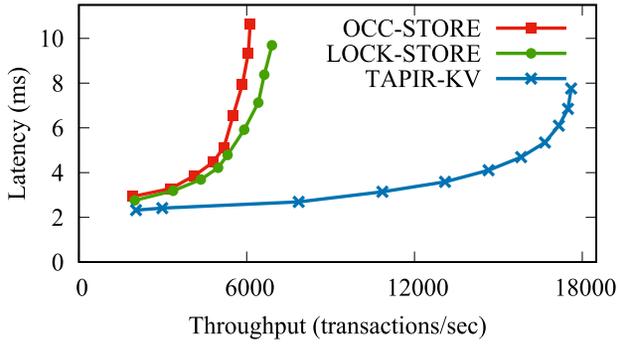


Fig. 16. Average Retwis transaction Latency (Zipf coefficient 0.75) versus throughput within a datacenter.

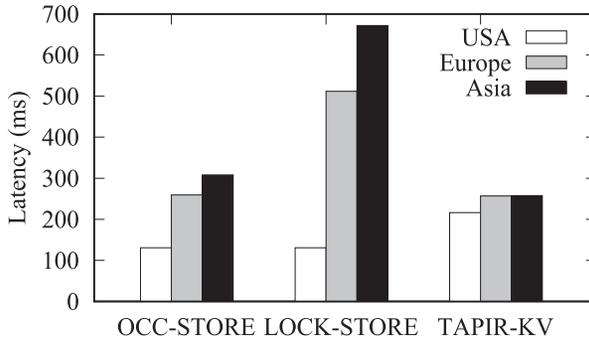


Fig. 17. Average wide-area latency for Retwis transactions, with leader located in the U.S. and client in U.S., Europe, or Asia.

between the U.S., Europe, and Asia. Figure 17 gives the average latency for Retwis transactions using the same workload as in previous section.

When the client shares a datacenter with the leader, the comparison systems are faster than TAPIR-KV, because TAPIR-KV must wait for responses from all replicas, which takes 160ms to Asia, while OCC-STORE and LOCK-STORE can commit with a round-trip to the local leader and one other replica, which is 115ms to Europe.

When the leader is in a different datacenter, LOCK-STORE suffers, because it must go to the leader on Read for locks, which takes up to 160ms from Asia to the U.S., while OCC-STORE can go to a local replica on Read like TAPIR-KV. In our setup TAPIR-KV takes longer to Commit, as it has to contact the *furthest* replica, and the RTT between Europe and Asia is more expensive than two round-trips between U.S. to Europe (likely because Google’s traffic goes through the U.S.). In fact, in this setup, IR’s slow path, at two RTT to the two closest replicas, is *faster* than its fast path, at one RTT to the furthest replica. We do not implement the optimization of running the fast and slow paths in parallel, which could provide better latency in this case.

7.4 Abort and Retry Rates

TAPIR is an optimistic protocol, so transactions can abort due to conflicts, as in other OCC systems. Moreover, TAPIR transactions can also be forced to abort or retry when conflicting timestamps are chosen due to clock skew. We measure the abort rate of TAPIR-KV compared to OCC-STORE, a conventional OCC design, for varying levels of contention (varying Zipf coefficients). These

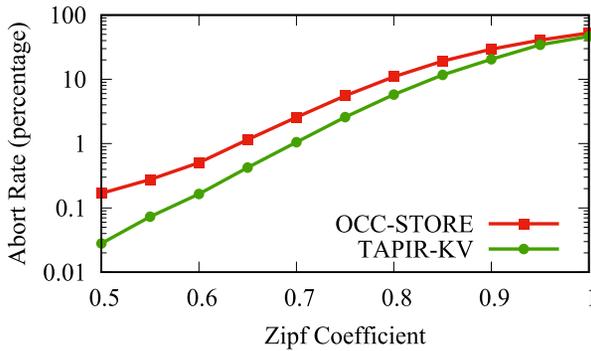


Fig. 18. Abort rates at varying Zipf co-efficients with a constant load of 5,000 transactions/second in a single datacenter. This graph excludes LOCK-STORE, which does not abort.

experiments run in a single region with replicas in three availability zones. We supply a constant load of 5,000 transactions/second.

With a uniform distribution, both TAPIR-KV and OCC-STORE have very low abort rates: 0.005% and 0.04%, respectively. Figure 18 gives the abort rate for Zipf co-efficients from 0.5 to 1.0. At lower Zipf co-efficients, TAPIR-KV has abort rates that are roughly an order of magnitude lower than OCC-STORE. TAPIR’s lower commit latency and use of optimistic timestamp ordering reduce the time between Prepare and Commit or Abort to a single round-trip, making transactions less likely to abort.

Under heavy contention (Zipf coefficient 0.95), both TAPIR-KV and OCC-STORE have moderately high abort rates: 36% and 40%, respectively, comparable to other OCC-based systems like MDCC (Kraska et al. 2013). These aborts are primarily due to the most popular keys being accessed very frequently. For these workloads, locking-based systems like LOCK-STORE would make better progress; however, clients would have to wait for extended periods to acquire locks.

TAPIR rarely needs to retry transactions due to clock skew. Even at moderate contention rates, and with simulated clock skew of up to 50ms, we saw less than 1% TAPIR retries and negligible increase in abort rates, demonstrating that commodity clock synchronization infrastructure is sufficient.

7.5 Comparison with Weakly Consistent Systems

We also compare TAPIR-KV with three widely used eventually consistent storage systems, MongoDB (MongoDB 2013), Cassandra (Lakshman and Malik 2010), and Redis (Redis 2013). For these experiments, we used YCSB+T (Dey et al. 2014), with a single shard with three replicas and 1M keys. MongoDB and Redis support master-slave replication; we set them to use synchronous replication by setting WriteConcern to REPLICAS_SAFE in MongoDB and the WAIT command (Sanfilippo 2013) for Redis. Cassandra uses REPLICATION_FACTOR=2 to store copies of each item at any 2 replicas. We chose to use synchronous replication to give similar fault-tolerance guarantees to TAPIR-KV, i.e., that an update is recorded to multiple replicas before being considered committed. Nevertheless, it provides weaker consistency guarantees.

Figure 19 demonstrates that the latency and throughput of TAPIR-KV is comparable to these systems. We do not claim this to be an entirely fair comparison; these systems have features that TAPIR-KV does not. At the same time, the other systems do not support distributed transactions—in some cases, not even single-node transactions—while TAPIR-KV runs a distributed transaction protocol that ensures strict serializability. Despite this, TAPIR-KV’s performance remains

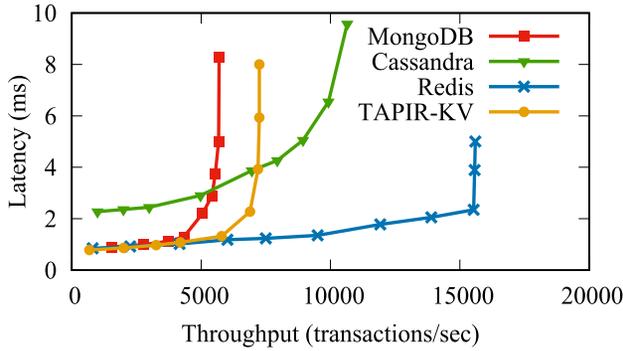


Fig. 19. Comparison with weakly consistent storage systems using the YCSB+T benchmark with one shard, three replicas, and 1M keys.

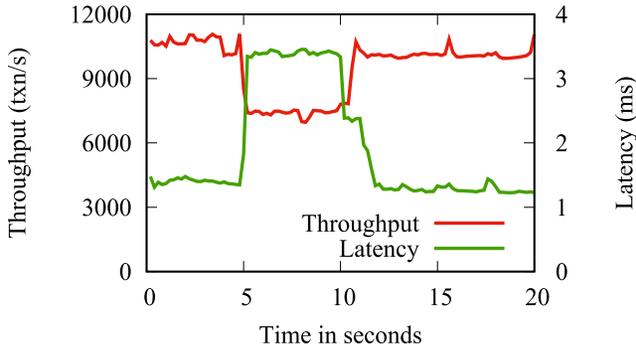


Fig. 20. Performance of TAPIR recovery protocol. One replica is failed at $t = 5s$, and restarted at $t = 10s$.

competitive: it outperforms MongoDB, and has throughput within a factor of 2 of Cassandra and Redis, demonstrating that strongly consistent distributed transactions are not incompatible with high performance.

7.6 Performance under Recovery

We studied the effect of TAPIR’s recovery mechanism on the protocol performance by measuring transaction throughput and commit latency while a replica has crashed and is recovering. For this experiment, we use a simple setup of a single shard with three replicas and a workload consisting of single key read-modify-write transactions. Each replica contains 1M keys and every transaction chooses a key at random uniformly. We kill one of the replicas 5s into the experiment and restart it 5s later at $t = 10s$. The restarting replica reloads all keys back into memory and runs the recovery protocol before beginning to process transactions again.

Figure 20 shows the drop in throughput and increase in latency when a single replica becomes unresponsive. The overall throughput drops, because each transaction takes the slow path, requiring one extra message to be processed to commit every transaction at each replica. The latency increases for the same reason, because the slow path requires every transaction to take an extra round-trip to commit. However, the main cause of the latency increase shown in Figure 20 is that each TAPIR client only begins the slow path protocol after a fixed timeout—here, 2ms. The increase in latency is directly dependent on this timeout. Reducing it would lead to a lower latency penalty during failures, but setting it too low would cause unnecessary messages to be transmitted by the

client and processed by each replica, as the slow path requires an additional exchange of messages. In the extreme, it may be possible to start the fast path and slow path in parallel. This would offer benefits in environments with highly skewed pairwise latency, where executing a multiround protocol among a simple majority of replicas may still offer lower latency than a single round to a supermajority of replicas, a classic problem for Fast Paxos-like protocols (Junqueira et al. 2007).

8 RELATED WORK

Inconsistent replication shares the same principle as past work on commutativity, causal consistency and eventual consistency: Operations that do not require ordering are more efficient. TAPIR leverages IR's weak guarantees, in combination with optimistic timestamp ordering and optimistic concurrency control, to provide semantics similar to past work on distributed transaction protocols but with both lower latency and higher throughput.

8.1 Replication

Transactional storage systems currently rely on strict consistency protocols, like Paxos (Lamport 2001) and VR (Oki and Liskov 1988). These protocols enforce a strict serial ordering of operations and no divergence of replicas. IR is most closely related to quorum systems (Gifford 1979; Malkhi and Reiter 1998), which achieve durability by making sufficient copies of a data. Like IR's **inconsistent** operations, these systems can make operations durable more rapidly than consensus protocols, but do not guarantee strict serial ordering. IR's **consensus** operations achieve single-round consensus in the absence of contention, like Fast Paxos (Lamport 2006a) or Bosco (Song and van Renesse 2008). When replicas disagree on the consensus result, it falls back to a two-round slow path; a similar hybrid approach was used for the Byzantine environment by HQ Replication (Cowling et al. 2006).

IR is also closely related to eventually consistent replication protocols, like Bayou (Terry et al. 1995), Dynamo (DeCandia et al. 2007), and others (Ladin et al. 1992; Saito and Shapiro 2005; Lakshman and Malik 2010). The key difference is that applications resolve conflicts after they happen with eventually consistent protocols, whereas IR **consensus** operations allow applications to decide conflicts and recover that decision later. As a result, applications can enforce higher-level guarantees (e.g., mutual exclusion, strict serializability) that they cannot with eventual consistency.

IR is also related to replication protocols that avoid coordination for *commutative operations* (e.g., Generalized Paxos (Lamport 2005), EPaxos (Moraru et al. 2013)). These protocols are more general than IR, because they do not require application invariants to be pairwise. For example, EPaxos could support invariants on bank account balances, while IR cannot. However, these protocols consider two operations to commute if their order does not matter when applied to *any* state, whereas IR requires only that they produce the same results *in a particular execution*. This is a form of state-dependent commutativity similar to SIM-commutativity (Clements et al. 2013). As a result, in the example from Section 3.1.3, EPaxos would consider any operations on the same lock to conflict, whereas IR would allow two unsuccessful Lock operations to the same lock to commute.

8.2 Distributed Transactions

A technique similar to optimistic timestamp ordering was first explored by Thomas (Thomas 1979), while CLOCC (Adya et al. 1995) was the first to combine it with loosely synchronized clocks. We extend Thomas's algorithm to: (1) support multiple shards, (2) eliminate synchronous disk writes, and (3) ensure availability across coordinator failures. Spanner (Corbett et al. 2012) and Granola (Cowling and Liskov 2012) are two recent systems that use loosely synchronized clocks to improve performance for read-only transactions and independent transactions, respectively. TAPIR's use

Table 3. Comparison of Read-write Transaction Protocols in Replicated Transactional Storage Systems

Transaction System	Replication Protocol	Read Latency	Commit Latency	Msg At Bottleneck	Isolation Level	Transaction Model
Spanner	Multi-Paxos	2 (leader)	4	$2n + \text{reads}$	Strict Serializable	Interactive
MDCC	Gen. Paxos	2 (any)	3	$2n$	Read-Committed	Interactive
Repl. Commit	Paxos	$2n$	4	2	Serializable	Interactive
CLOCC	VR	2 (any)	4	$2n$	Serializable	Interactive
Lynx	Chain Repl.	–	$2n$	2	Serializable	Stored procedure
TAPIR	IR	2 (to any)	2	2	Strict Serializable	Interactive

of loosely synchronized clocks differs from Spanner’s in two key ways: (1) TAPIR depends on clock synchronization only for performance, not correctness, and (2) TAPIR’s performance is tied to the *actual* clock skew, not TrueTime’s worst-case estimated bound. Spanner’s approach for read-only transactions complements TAPIR’s high-performance read-write transactions, and the two could be easily combined. Like TAPIR and CLOCC, Clock-SI (Du et al. 2013) uses loosely synchronized clocks to implement an optimistic concurrency control protocol, and its performance depends on the actual rather than worst-case clock skew. However, it provides a weaker isolation level (snapshot isolation (Berenson et al. 1995) rather than strict serializability); it also does not incorporate replication.

CLOCC and Granola were both combined with VR (Liskov et al. 1999) to replace synchronous disk writes with in-memory replication. These combinations still suffer from the same redundancy—enforcing ordering both at the distributed transaction and replication level—that we discussed in Section 2. Other layered protocols, like the examples shown in Table 3, have a similar performance limitation.

Recent work has applied new hardware technologies to support faster transaction processing. In particular, FaRM (Dragojević et al. 2014, 2015) uses high-performance RDMA networks to build a strict serializable transaction processing protocol. It optimizes for the local cluster network, where latency is far lower; in this environment, minimizing the number of message exchanges (as TAPIR does) is less critical. DrTM (Wei et al. 2015; Chen et al. 2016) combines RDMA networks with hardware transactional memory to accelerate concurrency control checks on multicore machines; this approach could be used to extend TAPIR to the multicore setting. Finally, Eris (Li et al. 2017) uses an in-network sequencer (Li et al. 2016) co-designed with a concurrency control to provide coordination-free execution of independent transactions. Inspired by TAPIR, Eris combines both replication and concurrency control into a single protocol to minimize the overhead of both.

Some previous work included in Table 3 improves throughput (e.g., Warp (Escriva et al. 2013), Transaction Chains (Zhang et al. 2013), Tango (Balakrishnan et al. 2013)), while others improve performance for read-only transactions (e.g., MegaStore (Baker et al. 2011), Spanner (Corbett et al. 2012)) or other limited transaction types (e.g., Sinfonia’s mini-transactions (Aguilera et al. 2007), Granola’s independent transactions (Cowling and Liskov 2012), Lynx’s transaction chains (Zhang et al. 2013), and MDCC’s commutative transactions (Kraska et al. 2013)), or weaker consistency guarantees (Lloyd et al. 2011; Sovran et al. 2011). In comparison, TAPIR is the first transaction protocol to provide better performance (both throughput and latency) for general-purpose, read-write transactions using replication.

A design goal for TAPIR was to support fully general, interactive transactions. Other optimizations are possible for specific transaction models. Subsequent to the initial publication of this work, other systems have applied TAPIR’s idea of integrating concurrency control and replication layers

to more restricted transaction models. As mentioned above, Eris (Li et al. 2017) uses a network-level sequencer and optimizes for independent transactions; it also supports fully general transactions through an additional protocol. Janus (Mu et al. 2016) considers a one-shot transaction model where write sets are known at transaction start; this allows it to employ transaction reordering (Mu et al. 2014) to increase the commit rate under highly contented workloads.

9 CONCLUSION

This article demonstrates that it is possible to build distributed transactions with better performance and strong consistency semantics by building on a replication protocol with *no* consistency. We present inconsistent replication, a new replication protocol that provides fault tolerance without consistency, and TAPIR, a new distributed transaction protocol that provides linearizable transactions using IR. We combined IR and TAPIR in TAPIR-KV, a distributed transactional key-value storage system. Our experiments demonstrate that TAPIR-KV lowers commit latency by 50% and increases throughput by 3× relative to conventional transactional storage systems. In many cases, it matches the performance of weakly consistent systems while providing much stronger guarantees.

ACKNOWLEDGMENTS

We thank the anonymous reviewers from both SOSP and TOCS and our SOSP shepherd Miguel Castro for their helpful feedback. We thank Michael Whittaker for implementing the IR recovery protocol, and Jialin Li, Neha Narula, and Xi Wang for early feedback on the article. We also appreciate the support of our local zoo tapirs, Ulan and Bintang.

REFERENCES

- Atul Adya, Robert Gruber, Barbara Liskov, and Umesh Maheshwari. 1995. Efficient optimistic concurrency control using loosely synchronized clocks. In *Proceedings of the ACM International Conference on Management of Data (SIGMOD'95)*.
- Marcos K. Aguilera, Arif Merchant, Mehul Shah, Alistair Veitch, and Christos Karamanolis. 2007. Sinfonia: A new paradigm for building scalable distributed systems. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP'07)*.
- Peter Bailis, Aaron Davidson, Alan Fekete, Ali Ghodsi, Joseph M. Hellerstein, and Ion Stoica. 2014. Highly available transactions: Virtues and limitations. In *Proceedings of the Conference on Very Large Databases (VLDB'14)*.
- Jason Baker, Chris Bond, James Corbett, J. J. Furman, Andrey Khorlin, James Larson, Jean-Michel Léon, Yawei Li, Alexander Lloyd, and Vadim Yushprakh. 2011. Megastore: Providing scalable, highly available storage for interactive services. In *Proceedings of the Conference on Innovative Data Systems Research (CIDR'11)*.
- Mahesh Balakrishnan, Dahlia Malkhi, Ted Wobber, Ming Wu, Vijayan Prabhakaran, Michael Wei, John D. Davis, Sriram Rao, Tao Zou, and Avi Zuck. 2013. Tango: Distributed data structures over a shared log. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP'13)*.
- Hal Berenson, Phil Bernstein, Jim Gray, Jim Melton, Elizabeth O'Neil, and Patrick O'Neil. 1995. A critique of ANSI SQL isolation levels. In *Proceedings of the 1995 ACM SIGMOD International Conference on Management of Data*. ACM.
- Philip A. Bernstein, Vassos Hadzilacos, and Nathan Goodman. 1987. *Concurrency Control and Recovery in Database Systems*. Addison Wesley.
- Ken Birman and Thomas A. Joseph. 1987. Exploiting virtual synchrony in distributed systems. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP'87)*.
- Mike Burrows. 2006. The Chubby lock service for loosely coupled distributed systems. In *Proceedings of the USENIX Symposium on Operating Systems Design and Implementation (OSDI'06)*.
- Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach, Mike Burrows, Tushar Chandra, Andrew Fikes, and Robert E. Gruber. 2008. Bigtable: A distributed storage system for structured data. *ACM Trans. Comput. Syst.* 26, 2, Article 4 (June 2008), 26.
- Yanzhe Chen, Xinda Wei, Jiaxin Shi, Rong Chen, and Haibo Chen. 2016. Fast and general distributed transactions using RDMA and HTM. In *Proceedings of the 11th ACM SIGOPS EuroSys (EuroSys'16)*. ACM.
- Austin T. Clements, M. Frans Kaashoek, Nikolai Zeldovich, Robert T. Morris, and Eddie Kohler. 2013. The scalable commutativity rule: Designing scalable software for multicore processors. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP'13)*.

- Brian F. Cooper, Raghu Ramakrishnan, Utkarsh Srivastava, Adam Silberstein, Philip Bohannon, Hans-Arno Jacobsen, Nick Puz, Daniel Weaver, and Ramana Yerneni. 2008. Pnuts: Yahoo!'s hosted data serving platform. *Proceedings of the Conference on Very Large Databases (VLDB'08)*.
- Brian F. Cooper, Adam Silberstein, Erwin Tam, Raghu Ramakrishnan, and Russell Sears. 2010. Benchmarking cloud serving systems with YCSB. In *Proceedings of the ACM Symposium on Cloud Computing (SOCC'10)*.
- James C. Corbett et al. 2012. Spanner: Google's globally distributed database. In *Proceedings of the USENIX Symposium on Operating Systems Design and Implementation (OSDI'12)*.
- James Cowling and Barbara Liskov. 2012. Granola: Low-overhead distributed transaction coordination. In *Proceedings of the USENIX Annual Technical Conference (ATC'12)*.
- James Cowling, Daniel Myers, Barbara Liskov, Rodrigo Rodrigues, and Liuba Shrira. 2006. HQ replication: A hybrid quorum protocol for Byzantine fault tolerance. In *Proceedings of the USENIX Symposium on Operating Systems Design and Implementation (OSDI'06)*.
- Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, Gunavardhan Kakulapati, Avinash Lakshman, Alex Pilchin, Swaminathan Sivasubramanian, Peter Vosshall, and Werner Vogels. 2007. Dynamo: Amazon's highly available key-value store. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP'07)*.
- Akon Dey, Alan Fekete, Raghunath Nambiar, and Uwe Rohm. 2014. YCSB+T: Benchmarking web-scale transactional databases. In *Proceedings of the International Conference on Data Engineering Workshops (ICDEW'14)*.
- Aleksandar Dragojević, Dushyanth Narayanan, Orion Hodson, and Miguel Castro. 2014. FaRM: Fast remote memory. In *Proceedings of the 11th USENIX Symposium on Networked Systems Design and Implementation (NSDI'14)*. USENIX.
- Aleksandar Dragojević, Dushyanth Narayanan, Edmund B. Nightingale, Matthew Renzelmann, Alex Shamis, Anirudh Badam, and Miguel Castro. 2015. No compromises: Distributed transactions with consistency, availability, and performance. In *Proceedings of the 25th ACM Symposium on Operating Systems Principles (SOSP'15)*. ACM.
- Jiaqing Du, Sameh Elnikety, and Willy Zwaenepoel. 2013. Clock-SI: Snapshot isolation for partitioned data stores using loosely synchronized clocks. In *Proceedings of the 32nd IEEE Symposium on Reliable Distributed Systems (SRDS'13)*. IEEE.
- Robert Escriva, Bernard Wong, and Emin Gun Sirer. 2013. *Warp: Multi-Key Transactions for Key-Value Stores*. Technical Report. Cornell.
- Michael J. Fischer, Nancy A. Lynch, and Michael S. Patterson. 1985. Impossibility of distributed consensus with one faulty process. *J. ACM* 32, 2 (Apr. 1985), 374–382.
- David K. Gifford. 1979. Weighted voting for replicated data. In *Proceedings of the USENIX Symposium on Operating Systems Design and Implementation (OSDI'79)*.
- Jim Gray and Leslie Lamport. 2006. Consensus on transaction commit. *ACM Trans. Database Syst.* 31, 1 (Mar. 2006), 133–160.
- Patrick Hunt, Mahadev Konar, Flavio Paiva Junqueira, and Benjamin Reed. 2010. ZooKeeper: Wait-free coordination for internet-scale systems. In *Proceedings of the USENIX Annual Technical Conference (ATC'10)*.
- Flavio Junqueira, Yanhua Mao, and Keith Marzullo. 2007. Classic Paxos vs Fast Paxos: Caveat emptor. In *Proceedings of the 3rd Workshop on Hot Topics in System Dependability (HotDep'07)*. USENIX.
- David Karger, Eric Lehman, Tom Leighton, Rina Panigrahy, Matthew Levine, and Daniel Lewin. 1997. Consistent hashing and random trees: Distributed caching protocols for relieving hot spots on the world wide web. In *Proceedings of the ACM Symposium on Theory of Computing (STOC'97)*.
- Tim Kraska, Gene Pang, Michael J. Franklin, Samuel Madden, and Alan Fekete. 2013. MDCC: Multi-data center consistency. In *Proceedings of the ACM SIGOPS EuroSys (EuroSys'13)*.
- Hsiang-Tsung Kung and John T. Robinson. 1981. On optimistic methods for concurrency control. *ACM Trans. Database Syst.* 6, 2 (June 1981), 213–226.
- Rivka Ladin, Barbara Liskov, Liuba Shrira, and Sanjay Ghemawat. 1992. Providing high availability using lazy replication. *ACM Trans. Comput. Syst.* 10, 4 (Nov. 1992), 360–391.
- Avinash Lakshman and Prashant Malik. 2010. Cassandra: A decentralized structured storage system. *ACM SIGOPS Operat. Syst. Rev.* 44, 2 (Apr. 2010), 35–40.
- Leslie Lamport. 1994. *ACM Trans. Prog. Lang. Syst.* 16, 3 (May 1994), 872–923.
- Leslie Lamport. 2001. Paxos made simple. *ACM SIGACT News* 32, 4 (Dec. 2001), 51–58.
- Leslie Lamport. 2005. *Generalized Consensus and Paxos*. Technical Report 2005-33. Microsoft Research.
- Leslie Lamport. 2006a. Fast Paxos. *Distrib. Comput.* 19, 2 (2006).
- Leslie Lamport. 2006b. Lower bounds for asynchronous consensus. *Distrib. Comput.* 19, 2 (Oct. 2006), 104–125.
- Costin Leau. 2013. Spring Data Redis–Retwis-J. Retrieved from <http://docs.spring.io/spring-data/data-keyvalue/examples/retwisj/current/>.
- Jialin Li, Ellis Michael, and Dan R. K. Ports. 2017. Eris: Coordination-free consistent transactions using network multi-sequencing. In *Proceedings of the 26th ACM Symposium on Operating Systems Principles (SOSP'17)*. ACM.
- Jialin Li, Ellis Michael, Adriana Szekeres, Naveen Kr. Sharma, and Dan R. K. Ports. 2016. Just say no to Paxos overhead: Replacing consensus with network ordering. In *Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI'16)*. USENIX.

- Barbara Liskov, Miguel Castro, Liuba Shrira, and Atul Adya. 1999. Providing persistent objects in distributed systems. In *Proceedings of the European Conference on Object-Oriented Programming (ECOOP'99)*.
- Barbara Liskov and James Cowling. 2012. Viewstamped replication revisited. Technical report MIT-CSAIL-TR-2012-021. MIT.
- Wyatt Lloyd, Michael J. Freedman, Michael Kaminsky, and David G. Andersen. 2011. Don't settle for eventual: Scalable causal consistency for wide-area storage with COPS. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP'11)*.
- Hatem Mahmoud, Faisal Nawab, Alexander Pucher, Divyakant Agrawal, and Amr El Abbadi. 2013. Low-latency multi-database databases using replicated commit. *Proceedings of the Conference on Very Large Databases (VLDB'13)*.
- Dahlia Malkhi and Michael Reiter. 1998. Byzantine quorum systems. *Distrib. Comput.* 11 (1998), 203–213.
- MongoDB. 2013. MongoDB: A open-source document database. Retrieved from <http://www.mongodb.org/>.
- Iulian Moraru, David G. Andersen, and Michael Kaminsky. 2013. There is more consensus in Egalitarian parliaments. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP'13)*.
- Shuai Mu, Yang Cui, Yang Zhang, Wyatt Lloyd, and Jinyang Li. 2014. Extracting more concurrency from distributed transactions. In *Proceedings of the USENIX Symposium on Operating Systems Design and Implementation (OSDI'14)*.
- Shuai Mu, Lamont Nelson, Wyatt Lloyd, and Jinyang Li. 2016. Consolidating concurrency control and consensus for commits under conflicts. In *Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI'16)*. USENIX.
- Brian M. Oki and Barbara H. Liskov. 1988. Viewstamped replication: A new primary copy method to support highly available distributed systems. In *Proceedings of the ACM SIGACT-SIGOPS Symposium on Principles of Distributed Computing (PODC'88)*.
- Dan R. K. Ports, Jialin Li, Vincent Liu, Naveen Kr. Sharma, and Arvind Krishnamurthy. 2015. Designing distributed systems using approximate synchrony in data center networks. In *Proceedings of the USENIX Symposium on Networked Systems Design and Implementation (NSDI'15)*.
- Redis. 2013. Redis: Open Source Data Structure Server. Retrieved from <http://redis.io/>.
- Yasushi Saito and Marc Shapiro. 2005. Optimistic replication. *Comput. Surveys* 37, 1 (Mar. 2005), 42–81.
- Salvatore Sanfilippo. 2013. WAIT: Synchronous replication for Redis. Retrieved from <http://antirez.com/news/66>.
- Yee Jiun Song and Robbert van Renesse. 2008. Bosco: One-step Byzantine asynchronous consensus. In *Proceedings of the International Symposium on Distributed Computing (DISC'08)*.
- Yair Sovran, Russell Power, Marcos K. Aguilera, and Jinyang Li. 2011. Transactional storage for geo-replicated systems. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP'11)*.
- Douglas B. Terry, Marvin M. Theimer, Karin Petersen, Alan J. Demers, Mike J. Spreitzer, and Carl H. Hauser. 1995. Managing update conflicts in Bayou, a weakly connected replicated storage system. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP'95)*.
- Robert H. Thomas. 1979. A majority consensus approach to concurrency control for multiple copy databases. *ACM Trans. Database Syst.* 4, 2 (June 1979), 180–209.
- Xingda Wei, Jiaxin Shi, Yanzhe Chen, Rong Chen, and Haibo Chen. 2015. Fast in-memory transaction processing using RDMA and HTM. In *Proceedings of the 25th ACM Symposium on Operating Systems Principles (SOSP'15)*. ACM.
- Irene Zhang, Naveen Kr. Sharma, Adriana Szekeres, Arvind Krishnamurthy, and Dan R. K. Ports. 2015a. Building consistent transactions with inconsistent replication. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP'15)*.
- Irene Zhang, Naveen Kr. Sharma, Adriana Szekeres, Arvind Krishnamurthy, and Dan R. K. Ports. 2015b. *Building Consistent Transactions with Inconsistent Replication (extended version)*. Technical Report 2014-12-01 v2. University of Washington. Retrieved from <http://syslab.cs.washington.edu/papers/tapir-tr-v2.pdf>.
- Yang Zhang, Russell Power, Siyuan Zhou, Yair Sovran, Marcos K. Aguilera, and Jinyang Li. 2013. Transaction chains: Achieving serializability with low latency in geo-distributed storage systems. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP'13)*.

Received October 2016; revised July 2018; accepted August 2018