

Rolis: A software approach to efficiently replicating multi-core transactions

Anonymous Author(s)

Submission Id: <63>

Abstract

This paper presents Rolis, a new speedy and fault-tolerant replicated multi-core transactional database system. Rolis’s aim is to mask the high cost of replication by ensuring that cores are always doing useful work and not waiting for each other or for other replicas. Rolis achieves this by not mixing the multi-core concurrency control with multi-machine replication, as is traditionally done by systems that use Paxos to replicate the transaction commit protocol. Instead, Rolis takes an “execute-replicate-replay” model. Rolis first speculatively executes the transaction on the leader machine, and then replicates the per-thread transaction log to the followers using a novel protocol that leverages independent Paxos instances to avoid coordination, while still allowing followers to safely replay. The execution, replication, and replay are carefully designed to be scalable and have nearly zero coordination overhead across cores. Our evaluation shows that Rolis can achieve 1.24M transactions per second on the TPC-C workload, using a 3-replica setup where each server has 32 cores. This throughput result is orders of magnitude higher than traditional software approaches we tested (e.g., 2PL), and is comparable to state-of-the-art, fault-tolerant, in-memory storage systems built using kernel bypass and advanced networking hardware, even though Rolis runs on commodity machines.

1 Introduction

Transactional storage systems are a key backend component in large-scale online services. They typically come in two flavors: single-machine multi-core databases [22, 23, 26, 41], and distributed and replicated databases [7, 32, 40, 49]. The former tend to be faster and achieve higher throughput owing to the lack of coordination across machines, while the latter achieve important properties such as the ability to continue operating in the presence of crash failures. Recent works [5, 11, 17, 20, 38, 39] attempt to bridge this performance gap by combining multiple machines to achieve fault tolerance, with multiple cores per machine, each of which operates on a different database partition, to achieve high throughput. Despite this progress, we find that they still fall short of the performance achieved by single-machine multi-core databases, owing to the cross-replica coordination needed by distributed transactions.

This paper explores the following research questions: can we significantly improve the throughput of replicated transactions (the bottleneck in making distributed databases perform as fast as multi-core ones) with a more clever coordination protocol between a multi-core leader and its multi-core replicas? And can this protocol finally close the gap between single-machine multi-core and fault-tolerant distributed databases? Our proposed system, *Rolis*, answers both questions affirmatively. Perhaps surprisingly, Rolis’s approach to managing and delegating work between a multi-core leader and its multi-core replicas achieves higher throughput than recent works that rely on advanced networking hardware and kernel bypass [5, 11, 20, 39, 45], even though Rolis runs on commodity servers.

Our main idea is to start with a well-known principle: if one maximizes the pipeline of transaction processing and replication by having more transactions outstanding, the long latency of distributed coordination protocols is roughly masked and the system can achieve similar throughput as a non-replicated transactional system. Of course, just because a principle is well-known does not mean that it can be applied easily. Many designs of distributed transaction systems cannot simply increase the number of outstanding transactions, as this significantly degrades performance (the opposite of our aim). The main reason for this is that the replication protocol is closely tied to the transaction execution and commit protocols. For example, Google’s Spanner [7] applies Paxos to replicate the critical steps in two-phase locking and two-phase commit; increasing the number of outstanding transactions will not improve the system’s performance, but will instead cause more aborts by increasing the chance of conflicting accesses. More recent protocols like Calvin [40] and Janus [32] partially overcome this limitation, but still fall short of achieving the performance of a single multi-core transactional store.

Rolis’s key contribution finds a way out of this conundrum. Rolis uses a combination of speculation and deterministic replay to ensure that we can increase the number of outstanding transactions—maximizing the pipeline and masking the cost of distributed coordination protocols—without causing more aborts. In more detail, Rolis first speculatively executes all transactions on the leader. Then, each thread of the leader creates an independent Paxos stream and uses it to replicate its transaction logs, including the serialization order, to a corresponding thread at each follower. Finally, the followers deterministically replay the logs to arrive at the same state.

The main technical difficulties that Rolis must overcome are: (1) how to design a replication protocol that allows each Paxos stream to work independently without having to coordinate with each other, and without having more outstanding transactions cause aborts; (2) how to replay transactions at the followers concurrently to ensure that they can keep pace with the leader; and (3) how to ensure that leader failures do not leave the different Paxos streams in an inconsistent state.

This paper addresses the above challenges and makes the following contributions:

- Identifies the performance gap between distributed transaction systems and multi-core transaction systems that cannot be remedied by merely increasing the number of outstanding requests in traditional approaches.
- Introduces a novel watermark tracking method to eliminate coordination between independent Paxos streams. This watermark is also used for visibility control and safe deterministic replay at followers.
- Describes an implementation and evaluation of Rolis running on a 3-replica setup with 32-core Azure servers that can process 1.24M TPC-C transactions per second. This throughput is an order of magnitude higher than prior distributed transaction protocols on commodity machines, and it is competitive (and even higher) than recent systems that rely on kernel bypass and RDMA NICs.

The major limitation of Rolis is that it does not (yet) support sharding. While we would like to support a system that scales linearly with more shards, we find that Rolis’s throughput, with a single shard, is already equivalent to a system with ~1,000 shards, and should suffice for many applications.

2 Overview

2.1 Background

Setup. In this paper we study the replication of multi-core transactions. We describe our design and implementation in the context of a key-value storage system, but our design can also apply to other settings, e.g., relational databases and software transactional memory systems. In our setup, a transaction includes multiple read, write, and range query operations that may access different keys. Concurrent transactions that access overlapping keys are isolated, replicated consistently, and obey their chronological order, i.e., transactions are *strictly serializable* [16, 36].

Consensus and replication. We assume an *asynchronous* network: messages can be arbitrarily delayed and there is no perfect failure detector that can detect or force a failure in the system. The standard way to achieve consistent replication in an asynchronous network is through consensus protocols such as Paxos [24] or Raft [34]. The standard interface of consensus-based replication is state machine replication (SMR): replicate a sequence of operation logs

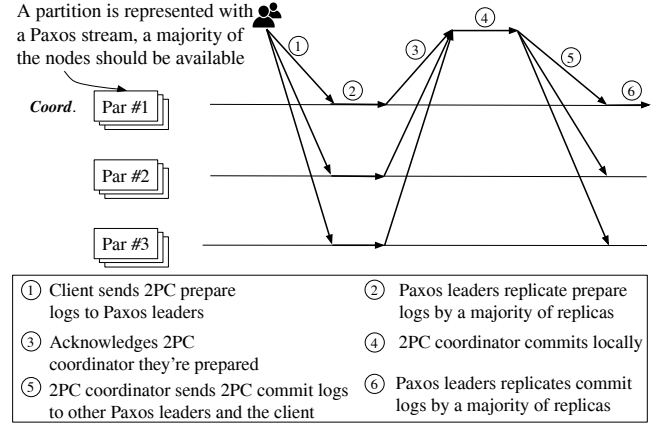


Figure 1. Spanner: two-phase commit along with replication on the Paxos replicated partitions. There are three partitions in this example, one partition serves as a coordinator, the others are the participants.

to all replicas and then have each replica deterministically apply those logs in the same order.

Distributed transactions. The de facto approach for distributing and replicating a transaction across different servers is to use consensus protocols to replicate the critical steps in the transaction execution and commit protocols. Taking Google Spanner for example, it uses two-phase commit (2PC) as a transaction commit protocol and uses Paxos to replicate the critical steps of 2PC, as shown in Figure 1. Other systems may have differences in the protocols they choose, but they all share an important similarity with Spanner: the replication happens side-by-side with the transaction execution and/or commit. That is, replication happens before the transaction’s serialization order is determined by the transaction execution/commit protocol.

2.2 Problem statement and strawman

In this paper we wish to answer the question: can we significantly improve the throughput of a replicated transactional key-value store that leverages the many cores available in today’s servers, with a more clever coordination protocol between a leader and its replicas?

To answer this question, we start by analyzing a simple strawman that separates replication and transaction execution/commit. In this strawman, a multi-core leader first executes and serializes all transaction requests into a log, and then uses MultiPaxos to replicate the log to all replicas. When the replication is complete, the leader responds to clients. Each replica, which is optimized to use multiple threads to replay the log (we discuss how to do this in Section 3.4) then deterministically executes the transactions in the log.

Figure 2 shows the performance of this strawman replicated database on the TPC-C benchmark. It is able to achieve

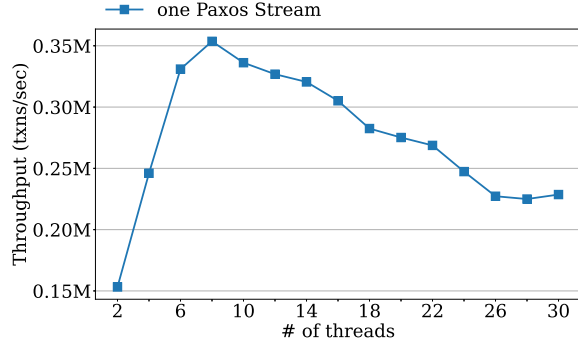


Figure 2. Throughput of TPC-C benchmark in a 3-replica setting with single Paxos stream replication.

over 0.35 million transactions per second, which is already substantially higher than a single shard in traditional replicate-before-commit approaches. However, the performance degrades after the leader and replicas use more than ~ 8 threads. This is due to the bottleneck in the system: a single log is replicated by MultiPaxos, which we refer to as a single *Paxos stream*. All threads on the leader need to write to this Paxos stream and will be bottlenecked by the thread synchronization cost and the replication rate of a single stream. Some recent works [1, 37, 42] proposed RDMA-based Paxos protocol that aims to be scalable to client connections. In theory, it is possible that those systems can be scalable to the higher thread number (i.e., 20 threads), but this is not our goal.

2.3 Challenges and high-level idea

A potential solution to the aforementioned bottleneck is to use more than one Paxos stream. For example, each thread could get its own stream. There are, however, two challenges with using multiple Paxos streams. First, with a single stream we can rely on the natural order of the requests in the stream as the serialization order of transactions, and all replicas can simply follow that order. With more than one stream, requests across streams are not ordered relative to each other, so we need to find another way to ensure that all replicas execute transactions in the same order.

To address this issue, Rolis adopts an execute-replicate-replay model. First, we have one replica (the leader) execute and “commit” the transaction locally as if no replication were to take place, except that the commit is held back from writing real values into the database. Rolis uses the commit operation to determine a serialization order among transactions, which will be captured through a monotonically increasing hardware counter that is available on commodity machines. This timestamp will be encoded into each log and will be used to order requests across different Paxos streams. As we will see, this timestamp can be used to determine

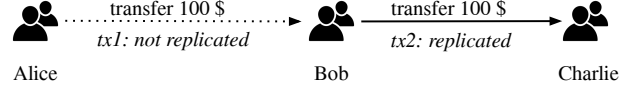


Figure 3. A transitive transaction: $tx1$ fails to be replicated and $tx2$ succeeds being replicated while $tx2$ depends on $tx1$.

dependencies between transactions without the need for coordination among threads.

The second problem is more subtle and difficult to solve; it relates to the failure recovery of the Paxos streams. Suppose that each Paxos stream is independent and therefore recovers independently of the other streams. This could lead to a situation where transactions that are ordered early in the serialization order are lost in recovery, but transactions that are ordered later in the serialization order are preserved. We give an example of this behavior in Figure 3.

Consider two transactions $tx1$ and $tx2$. In $tx1$, Alice transfers \$100 to Bob; in $tx2$, Bob transfers \$100 to Charlie. $tx1$ happens before $tx2$ in the leader’s execution. $tx1$ is replicated by thread-1 in stream-1 and $tx2$ by thread-2 in stream-2. If there is a failure, a possible outcome for the two independent Paxos streams is that stream-1 recovers with $tx1$ not replicated, and stream-2 recovers with $tx2$ replicated. In this case, if replicas only replay $tx2$ ’s transaction, the system will be in an incorrect state, where Alice never transferred her money out but Charlie got an extra \$100.

To address this issue, we need some way to track the dependencies between transactions to make sure that, during failure recovery, transactions with missing dependencies are not replayed—e.g., in the above example, $tx2$ should not be replayed. This is unfortunate, since our initial goal was to avoid all coordination between threads. Nevertheless, we devise a low overhead mechanism for tracking dependencies that preserves the system’s performance. In the context of multi-core transactional systems, low-overhead coordination is well known to be a challenging problem [26, 39]. Rolis’s innovation is the use of a lightweight mechanism to track dependencies based on the idea of keeping a *watermark* across all Paxos streams (threads). The watermark is effectively a boundary for replay visibility. All threads will synchronize periodically (with a frequency that does not affect performance) to advance the watermark. All transactions that fall within the watermark are safe to replay, and those beyond the watermark are unsafe to replay. We discuss the details of this and other components of Rolis in the following sections.

3 Main System Design

This section describes the main design of Rolis. First, we overview the architecture of Rolis including its major components and workflow. Then, we describe each stage of the workflow in more detail.

3.1 Architecture

Rolis has two major building blocks that we use (almost) as opaque boxes: a multi-core in-memory high-speed transactional database, and a consensus-based replication layer. Our design is generic and is not tied to specific choices for these building blocks. In our implementation, we choose to use Silo [41] as the local database; other options (e.g., ERMA [22] and Cicada [26]) are applicable as well.¹ Silo is a speedy in-memory database that is fast and has good multi-core scalability. It uses an optimized optimistic concurrency control scheme to execute and commit transactions. We adopt Silo’s design and interface that the client (application) and database are in the same OS process. For the consensus-based replication, we implemented a standard MultiPaxos algorithm.

Figure 4 shows the architecture of our Rolis. Each replica can run either as a *leader* or as a *follower*. Only a leader can accept new transaction requests. A transaction is first executed on the leader replica, where it is guaranteed isolation from other transactions. The execution will generate a log for the transaction, which is assigned a globally unique timestamp that represents the serialization order of the transaction (§3.2). Then, the log is replicated to the followers via the replication layer (§3.3). To avoid any scalability bottleneck, the replication layer has multiple Paxos streams, each of which is dedicated to a worker thread in the database. The Paxos streams are entirely independent from each other in their replication, thus avoiding any coordination overhead between them.

When the transaction log is replicated to enough followers and is considered durable by the replication layer, the followers will replay the log. A follower sees multiple Paxos streams, and replays the logs in a scalable manner using multiple database threads, ensuring that they are applied in the same order as on the leader (§3.4). We also design a novel scheme (§4.1) to ensure that replay is safe against failures, so that data inconsistencies such as the one described earlier cannot occur. Once the leader confirms that the transaction’s log *will* be replayed—that is, the leader does not need to wait for the actual replay—the leader can release the transaction’s results to the client.

3.2 Executing transactions on the leader

Transactions enter the system following the original database’s (Silo’s) paths, and they are executed up to their commit point in Silo. For completeness, we briefly discuss how Silo works. In Silo, all worker threads are egalitarian and work in a shared memory. A thread starts a transaction and executes it until it finishes. Silo has the standard interfaces of transaction start, read, write, and end. Figure 5 shows an example of these interfaces.

¹When using non-strictly serializable databases such as Cicada [26], the isolation guarantee of Rolis downgrades to that of the local database.

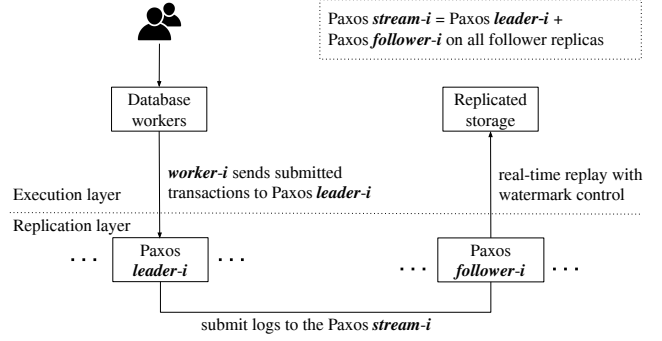


Figure 4. There are multiple database workers on the leader replica and each worker has its separate Paxos stream. There are two replicas (one leader plus one follower replica) and each replica runs a local instance of the Rolis system.

```
void *txn = db->new_txn();           // start a transaction
void *table = db->table_instance();
int account_alice = table->get("Alice"); // read a key
int account_bob = table->get("Bob");
assert(account_alice >= 100);
table->put("Alice", account_alice - 100); // write a key
table->put("Bob", account_bob + 100);
db->commit_txn(txn);                 // end a transaction
```

Figure 5. An example of transaction interfaces: Alice transfers \$100 to Bob using transaction *start*, *read*, *write* and *end* interfaces.

Silo executes transactions using OCC. During transaction execution, all the reads will record the current versions, and all the writes are buffered in a thread-local workspace. During transaction commit, the worker thread will first (spin-)lock all the keys in the write-set, and then it will certify that the keys in the read-set have not been updated by conflicting transactions, by comparing the most recent versions with the recorded versions. If the certification passes, the transaction will commit and all locks are released. If the certification fails, the transaction will abort and retry by starting over from the execution.

We made two modifications to Silo. First, if the certification passes, the worker thread will call the timestamp counter instruction, `rdtscp`, to obtain a monotonically increasing timestamp, before it unlocks the keys in the write-set. The `rdtscp` counter is not a scaling bottleneck for OLTP workloads these days [38], and the overhead is negligible compared to the serialization of transactions. We use this instruction to generate ordered timestamps across threads on the same machine. The timestamp represents a serialization order between transactions, which the system will utilize to detect dependencies in later replays. After unlocking the

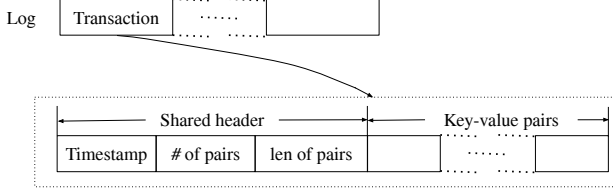


Figure 6. Decomposition of batched logs: each log contains several transactions and each transaction maintains the shared header and key-value pairs.

write set, the worker thread generates the log for replication, which contains the keys and values of the write set.

To reduce the overhead of network processing (e.g., interrupts), Rolis batches many transactions into a log (e.g., 1000 transactions in our implementation), which significantly reduces the network transmission and communication cost between two replicas. Rolis borrows ideas from the TCP protocol to reduce the size of serialized logs. In Figure 6, instead of encoding a timestamp into each key-value pair, we provide a piece of header information, which includes a timestamp, count of key-value pairs, and length of key-value pairs in the same transaction. We use the timestamp from the last transaction in a log to represent the timestamp of the log; this timestamp is compared against the watermark during replay (§3.4). For each key-value pair, we include the key and value information, along with the necessary metadata, to facilitate replay in the follower replicas.

The second change we made is holding the worker thread from releasing commit results to the clients until the replication finishes. We use a watermark to control this. Only transactions whose timestamp is smaller than the watermark are safe to release to the clients; transactions whose timestamp is greater are still speculative. Section 3.4 explains how the watermark is designed and how it grows so that it is not a bottleneck in the system.

Note that while the worker thread holds off on releasing results to the clients, the thread can still process new requests. In this way, the throughput of the system is preserved and only latency is affected. This introduces another problem in our system design. Since Rolis runs at a high throughput (over a million transactions per second), there will be (tens of) thousands of transactions whose execution is speculative. In the common case, this is okay, because the leader is stable for most of the time and the speculative executions will eventually commit. In the worst case, the speculatively executed transactions need to be abandoned. Standard speculative approaches use undo logs to roll back, but in our case, logging undo logs and rolling back thousands of transactions is too costly. Instead, we take a simpler approach when an old leader rejoins after a crash, it will drop the inconsistent state and join as a new replica (see more discussion in Section 4.3).

3.3 Replicating with multiple Paxos streams

Rolis utilizes multiple Paxos streams, based on the intuition that associating each database worker thread with a separate Paxos stream will allow us to maximize scalability while minimizing cross-core synchronization. In general, state machine replication (e.g., Paxos [24] and Raft [34]) is a poor match for multi-core scalability, because SMR usually requires a sequential ordering of all operations. In our implementation, we use the same number of Paxos streams and database workers, and then pair each Paxos stream to a database worker. Each Paxos stream only needs to deal with transactions from its associated database worker.

Each Paxos stream receives a log from its database worker and replicates the log to at least a majority of replicas. After that, the log is durable. i.e., the system can recover the log from any healthy replica despite a minority of replica failures. Each Paxos stream reaches agreements independently of the other streams, and thus avoids the cost of serialization between multiple cores. The leader replica executes transactions speculatively, which are then replicated using each worker’s corresponding Paxos stream. The per-Paxos stream logs extends progressively with strictly increasing timestamps. Since each Paxos stream on the follower replica has a correct partial order, we can progressively mark consistent snapshots on the followers and replay them asynchronously. In the following section, we describe this replay mechanism in detail.

In order to make it possible to compare timestamps generated by different leaders after a failover, we make our Paxos implementation use a monotonically increasing *epoch* number to distinguish a new leader when there is a leader change. The epoch number is encoded into each transaction’s log. This epoch-based method for distinguishing leaders is widely used in consensus algorithms, e.g., it is similar to the term number in Raft. Together, $\langle \text{epoch}, \text{timestamp} \rangle$ forms a pair that we can use to serialize all transactions.

3.4 Replaying transactions on followers

When a new log is durable in the Paxos stream, a follower cannot replay the log just yet, since this might lead to the inconsistencies discussed in Section 2.3. Instead, we use a watermark scheme to enable safe replay. We avoid the use of explicit dependency tracking [2, 3, 12] for two reasons: (1) explicit dependency tracking (i.e., maintaining a dependency graph) typically involves more expensive and complicated protocols compared to our simple watermark tracking; and (2) explicit dependency tracking may result in cycles, which ultimately result in higher tail latencies [3].

The watermark tracking works as follows. Suppose we have n Paxos streams $1, 2, \dots, n$, for a certain epoch e , and the most recent durable timestamp for this epoch is ts_1, ts_2, \dots, ts_n on a replica. Then the watermark W_e on this replica is:

$$W_e = \min(ts_1, ts_2, \dots, ts_n)$$

| | | | | | | | | | | | |
|----|--|----|----|----|----|----|----|----|----|----|----|
| P1 | 1 | 12 | 24 | 34 | 42 | 59 | 60 | 82 | 83 | 84 | 85 |
| P2 | 3 | 7 | 27 | 44 | 46 | 57 | 61 | 78 | 80 | | |
| P3 | 4 | 8 | 26 | 41 | 45 | 55 | 62 | 74 | 75 | | |
| P4 | 2 | 9 | 21 | 47 | 48 | 53 | 63 | 73 | | | |
| P5 | 5 | 11 | 23 | 50 | 52 | 67 | 69 | 70 | | | |
| | <div style="text-align: center;"> durability committed logs </div> | | | | | | | | | | |

Figure 7. There are five Paxos streams and logs in the same color are computed within the same calculation interval (0.5ms), each entry contains timestamp of a log.

All transaction logs whose timestamps are smaller than W_e in epoch e are safe against failures. For a leader replica, these transactions are safe to release results to the clients. For a follower replica, these transactions are safe to replay. The replay happens as follows. Similar to transaction execution on the leader, each Paxos stream has a corresponding replay thread. The thread will replay the logs below the watermark from the stream sequentially; all threads can replay concurrently. For each key/value pair in the log, the worker thread will do a “compare-and-swap” to determine if the database should be updated. If the key in the database has a smaller $\langle \text{epoch}, \text{timestamp} \rangle$ than that in the log, then the thread will update the database to the newer one. Otherwise, the thread does nothing and moves on to the next key. In this way, all worker threads can replay their logs concurrently, while ensuring that the latest state is reflected in the database. In our design we assume the compare-and-swap is atomic; how this can be done is discussed in Section 5.

The watermark calculation can be done asynchronously with respect to the log replication. That is, an outdated watermark does not affect safety, because the watermark is always growing and a transaction is always safe to replay once the watermark is beyond the transaction’s timestamp. This enables two aspects of our design that make the implementation easier. First, each replica can calculate its watermark on its own without any external communication. Second, the system does not need to calculate W_e instantly each time it is accessed, which would incur synchronization across threads. Instead, we calculate W_e periodically, e.g., every 0.5ms as in our evaluation.

Example. Figure 7 shows the logs durability committed by five Paxos streams in the first 1.5ms, when W_e is updated at 0.5ms intervals. Rolis advances the watermark from 8, to 44, to 70 at the corresponding intervals. Note that the watermark calculation always happens within the same epoch; i.e., each epoch has its own watermark tracking and the system is not allowed to advance the watermark across two different epochs. The reader may notice that during an epoch (leader)

change, the transactions in the previous epoch may not be safe to replay due to failures. We discuss this further in Section 4.1.

3.5 A review of stages in Rolis

One way to look at Rolis is that it separates the “commit points” of transaction processing and consensus processing. We would like to clarify these different commit points, because the building blocks we rely on have their own definitions of commits, which can be confusing when comparing Rolis to its peers. In fact, in Rolis, there are three separate commit points: in the in-memory database, in the replication layer, and in the holistic system.

The first commit point is in the in-memory database (Silo), where the serialization point of the transaction is speculatively determined, which we can refer to as *execution commit*.

The second commit point is in a Paxos stream, where a log is accepted by a majority of replicas and considered durable, which we can refer to as *durability commit*.

The third commit point is when the watermark grows beyond a transaction, so it is safe to release the transaction’s results and replay it, which we can refer to as *release commit*.

These commit points occur at the execution stage, replication stage, and replay stage, respectively. Note that the release commit is the real commit point of a transaction; the execution commit and durability commit are not. A transaction can be both execution and durability committed, but aborted by the system in the end.

4 Availability

In this section, we discuss how our protocol guarantees fault tolerance without sacrificing either correctness or availability. Rolis needs $2f + 1$ replicas to tolerate f failures. Rolis’s Paxos implementation is standard and resembles traditional Paxos systems (e.g., Chubby [4]); we give a brief review here.

Rolis runs an election module on all replicas, which periodically send and receive heartbeats to maintain the leader’s liveness in case of machine failures or network partitions. Each replica maintains a single epoch number that indicates the round of election. The follower replicas wait for a random time interval, and if no heartbeat is received within this interval, they declare the leader replica as failed and trigger a new round of election that increments the epoch number. When a new leader is elected, it can start receiving requests from the application and tries to commit them as normal.

In our Paxos implementation, we use an optimization that prevents a log entry from committing until all previous log entries in the same stream have committed. That is, each stream is always growing sequentially. This simplifies our design because the replay process does not need to deal with any “holes” in the stream where a later log entry commits before an earlier entry.

| | old epoch | | | | | new epoch | | |
|----|-----------|----|-----------|-----------|--------|-----------|----|----|
| P1 | 1 | 12 | <u>24</u> | <u>34</u> | no-ops | 1 | 11 | 23 |
| P2 | 3 | 7 | <u>27</u> | | no-ops | 3 | 13 | 20 |
| P3 | 4 | 8 | <u>26</u> | <u>41</u> | no-ops | 7 | 21 | 26 |
| P4 | 2 | 9 | 21 | | no-ops | 4 | 8 | 33 |
| P5 | 5 | 11 | <u>23</u> | | no-ops | 5 | 9 | 31 |

22
 the lost log

Figure 8. An example of failover: durability committed logs in five Paxos streams during a failure. Each entry shows the timestamp of the log.

4.1 Replay under failures

Using multiple Paxos streams to support high throughput raises two challenges for correctness. First, as transactions have dependency relationships, Rolis needs to determine when and if a transaction can be replayed safely. Second, the Paxos streams on the new leader replica might not have the latest durability-committed transactions locally.

To discuss this in more detail, once a new leader replica has been elected, all Paxos streams on the new leader will start bringing themselves up to date by retrieving any missing durability-committed transactions from other replicas. Then, the new leader commits a no-op in each Paxos stream to end the old epoch, after which it can calculate the latest watermark W_e based on timestamps from the old epoch. Then, Rolis can safely replay all transactions up to W_e . Importantly, Rolis can skip transactions above W_e in the old epoch, which is safe for two reasons. First, those transactions have never been executed on other follower replicas due to our watermark scheme. Second, Rolis never releases the results of those transactions to clients on the old leader, or any replica. In effect, the watermark represents the visibility boundary to clients. After replaying transactions in the old epoch, the new leader re-initializes the watermark to continue processing in the new epoch.

Example. Figure 8 shows an example of 5 Paxos streams, where all transactions have been durability committed. The timestamps are monotonically increasing in the same epoch within each Paxos stream. Whenever a replay thread sees a no-op, it immediately realizes that membership has changed and then waits for all other replay threads to receive no-ops as well. Once all replay threads observe no-ops, Rolis advances the latest watermark W_e to 21 (the smallest one among timestamps: 34, 27, 41, 21, 23) in the old epoch. All replay threads have to replay transactions up to 21, and must skip all remaining transactions (the underlined entries) as they might depend on other transactions that have not durability committed yet. For example, the log with timestamp 22

might not have durability committed yet in Paxos stream P4, and will be lost when a failure occurs. If the system replays logs with timestamps 26 and 41 from Paxos stream P3, the system may become inconsistent because logs with timestamp 26 and 41 might depend on the lost log. Once all replay threads succeed in replaying and skipping transactions in the old epoch, the new leader proceeds as normal in the new epoch (§3.4).

4.2 Correctness

We have explained in-place why our mechanisms work. Here, we give a brief review and explain why Rolis guarantees strict serializability across failovers. We separate the discussion into two cases, based on whether the system has replicated all proposed log entries in an epoch before the failover.

We first consider the simpler case when Paxos replication has replicated all log entries in an epoch. Within the same epoch (before a failure occurs), the leader generates a strictly serializable execution whose serial order is captured by the timestamps. This follows directly from Silo’s strict serializability guarantees and the way Rolis assigns timestamps. Thus, **transactions with smaller timestamps must be serialized before those with larger timestamps.**

When a failure occurs, the Paxos algorithm guarantees that all replicas see the same Paxos streams, a new leader always has a higher epoch number than the previous leaders, and the epoch numbers divide the stream into log entries generated by different leaders. **Transactions with a smaller epoch number must be serialized before those with a larger epoch number.** This is guaranteed by Rolis’s failover, because the new leader does not process requests until it has replayed all transactions from the previous epochs.

Given the $\langle \text{epoch}, \text{timestamp} \rangle$ pairs, the replay process will duplicate the final state of the execution, making the replicas consistent with the leader, by always prioritizing the state of transactions with larger epochs/timestamps (via our compare-and-swap method).

We now discuss the more complex case when Paxos replication does not finish replicating all transactions in an epoch, leaving an incomplete tail in the log stream. In this case, the claims we made above still hold as long as the following premise is true: Rolis can ensure a consistent state by stopping replay once it hits a missing log entry, and ignoring all entries afterwards; i.e., Rolis can replay a prefix of the original execution. In this case, we claim that **Rolis’s replay correctly preserves a prefix of the leader’s original proposed serializable execution of transactions.** If a transaction is replayed, due to the watermark mechanism discussed in Section 3.4, any transaction that has a smaller timestamp has been or will be replayed. Therefore, the replay always leads to a state that represents a prefix of the original execution. This is safe because (i) a prefix of a serializable execution is still serializable, and (ii) the execution beyond the

prefix can be discarded—and they are in our design—because they are never revealed to clients.

Since a new leader finishes the replay of previous epochs before proposing new transactions, the transactions in the new epoch will always read the latest changes.

4.3 Adding a new follower replica

In typical SMR implementations, adding a new replica usually requires a snapshot (checkpoint) of the system, which is transmitted to the new replica to bootstrap its state. In our case, Silo does not yet support snapshots: generating snapshots without degrading performance in a multi-core system is challenging, and is not supported by most state-of-art multi-core databases we discussed. So in the spirit of supporting more local database choices (beyond Silo) in our design, we adopt an approach from MongoDB [51] of adding new replicas without using snapshots. Briefly, a new replica first chooses a follower replica as its synchronization source, and performs an asynchronous “pull” from the source by scanning the source’s local database, while the source is still working. Then, the new replica will retrieve all logs from the source and replays those logs. After that, the new replica is up to date and can join the replication group. The key reason why this solution works is that the log replay in Rolis is idempotent: repeating the same log entry multiple times does not change the system state.

5 Implementation

Compare-and-swap in replay. The compare-and-swap operation in the replay process (§3.4 needs to be atomic. In principle, this can be done by modifying the underlying data structure of Silo (i.e., Masstree [31]) to support this operation using lock-free instructions or spin-locks. This would achieve better performance but it also requires invasive changes and limits the portability of our to other systems. Instead, we take a slightly slower but more flexible approach: wrapping the compare-and-swap as a Silo transaction. Although this adds overhead to the replay, the replay is still faster than the leader’s execution in our tests, as we explain later (§6, Figure 15).

Heartbeat with empty transactions. A small issue we have not discussed yet is how the system commits the tail of the log if no new log entries show up. This issue happens when the system is idle or in the process of shutting down. Consider the example in Figure 7, and assume we attempt to stop the system at 1.5ms. The system cannot replay logs with timestamp above 70 and will wait forever because no new logs will come in, and the watermark cannot be advanced. To solve this issue, we add a special empty transaction in the heartbeat to every Paxos stream. These empty transactions help the replay finish the transactions hanging at the tail of the log.

6 Evaluation

This section presents our evaluation results, focusing on answering the following questions:

- Can Rolis preserve Silo’s performance (including its scalability) in a multi-core setup?
- How does the performance of Rolis compare to that of the state-of-art systems with advanced NIC support, and to traditional software systems?
- How fast can Rolis recover from failures?

6.1 Experimental setup

To evaluate the multi-core scalability of Rolis, all experiments were run on multiple Azure virtual machines based on model Intel Xeon Platinum 8272CL CPU @ 2.60GHz Processor within the same datacenter. Each machine has 32 (hyperthreaded) CPU cores, 128G RAM, and is interconnected via a 16000 Mbps network. Unless otherwise mentioned, the experiments are conducted with 3 replicas. Each trial is run for 30 seconds, the same as Silo’s original test configuration. Throughput and latency are calculated based on the release committed transactions. When testing the scalability over threads, we use the `cgroups` kernel feature to limit the CPU and RAM resources. In addition, we always need 1 extra CPU core to advance the watermark and perform leader election tasks, in the event of resources contention with the database workers and Paxos streams.

Like prior work [40, 41, 44], our experiments do not use networked clients, instead, we bind a workload generator within the servers to simplify the tests.

Our experiments run TPC-C [18] and YCSB++ transactional benchmark. TPC-C is a common benchmark for OLTP workloads. YCSB++ is a simple transactional benchmark which is the benchmark derived from YCSB workload F, which consists of 50% Reads and 50% RMW [6, 47]. The set of records for each transaction in both benchmarks are selected uniformly from the entire database.

Figure 9 shows the percentage and number of read/write operations of each transaction type in TPC-C and YCSB++. For TPC-C we follow the official ratio for the five transaction types: NewOrder (NEW), Payment (PAY), OrderStatus (ORDER), StockLevel (STOCK), and Delivery (DLVR). We use TPC-C to demonstrate Rolis’s performance when dealing with complicated transactions. For YCSB++ we have two transaction types: Read-Only (READ) and Read-Modify-Write (RMW). The total data space for YCSB++ is 1 million keys (or per partition). We perform 4 updates per RMW operation and 4 reads per READ operation, and items are accessed at random. We use YCSB++ to push the limit of Rolis with a high transactional throughput workload. The batch sizes we use for TPC-C and YCSB++ are 1,000 and 10,000 respectively.

| TPC-C | NEW | PAY | DLVR | ORDER | STOCK | YCSB++ | READ | RMW |
|------------|--------|---------|------|---------|-------|------------|------|-----|
| Percent | 45% | 43% | 4% | 4% | 4% | Percent | 50% | 50% |
| Get/Scan | avg+23 | avg+3.6 | 50 | avg+3.6 | 3 | Get/Scan | 4 | 4 |
| Insert/Put | avg+24 | avg+4 | 180 | 0 | 0 | Insert/Put | 0 | 4 |

Figure 9. Workload of TPC-C and YCSB++ benchmarks. avg+ stands for an estimated number. Also, we treat scan and get as one read operation, and insert and put as one write operation.

6.2 Performance and scalability

In the first experiment, we evaluate the performance and scalability of Rolis with an increasing number of available CPU cores on benchmark TPC-C and YCSB++. We also compare Rolis’s results to the throughput of Silo on a single machine, i.e., without any replication. Since our Rolis is built on the top of Silo, the throughput of Silo is the upper bound of our implementation. In Figure 10 and Figure 11, we show that the throughput and per-core throughput of Rolis on both the TPC-C and YCSB++ benchmarks are good and scale well as we introduce more threads.

TPC-C benchmark: In Figure 10, the throughput of the Rolis at 32 cores is 1.24M transactions per second, which is 63.7% of Silo. The overhead is mainly caused by transaction (de)serialization and memory copies in the transmission. In our implementation, Rolis needs to serialize/de-serialize 771.8 bytes of data per transaction on average. Figure 11 shows per-core throughput is higher at the first 15 cores, and since then it becomes stable gradually which is the same as Silo. This decreasing tendency is mainly caused by several factors from Silo itself, including increased database size and sharing of resources such as the L3 cache.

YCSB++ benchmark: To better study the scalability of Rolis in different workloads, we run an experiment on YCSB++ benchmark, which is much simpler than TPC-C benchmark. In this simple but super-high throughput scenario, Rolis still can scale well, the throughput using 32 cores can be up to 10.3M transactions per second, which is 8.25× compared to TPC-C experiment. Rolis can retain 77.3% of the throughput of Silo on YCSB++ benchmark which is higher than TPC-C benchmark, mainly because YCSB++ has a smaller write-set.

6.3 Comparison with software implementations

In this experiment, we set the number of replicas of each partition to 3. We compare Rolis with the two-phase locking (2PL) with Paxos-based replication implementation in Janus [32] and with Calvin [40], which is deterministic concurrency control and replication algorithm implemented in STAR [30] on YCSB++ benchmark reported in Figure 12.

2PL. 2PL is the most widely used pessimistic concurrent control protocol. We use the implementation in Janus, a partitioned distributed data store. Each transaction performs 4

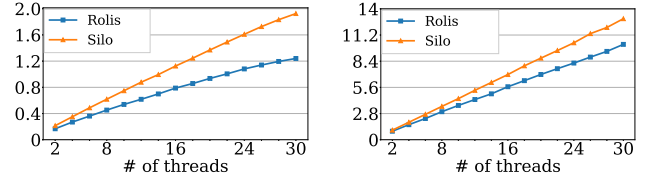


Figure 10. Throughput (million transactions per second in y-axis) over worker threads on TPC-C and YCSB++ benchmark.

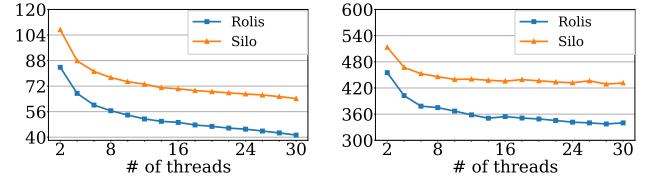


Figure 11. Per-core throughput (thousand transactions per second in y-axis) over worker threads on TPC-C and YCSB++ benchmark.

read-write access or 4 read access by incrementing 4 random chosen keys. Each partition takes up 1 CPU core resource for transactions execution. To minimize the extra cost of inter-process communication, we make each transaction only access a single partition. This gives the partitioned test targets extra advantage over Rolis.

Calvin. Calvin uses a central sequencer to determine the order of batched transactions which are sent to all replicas to execute deterministically later. We use the Calvin implementation in STAR [30] instead of the original one because the STAR-based implementation can have a multi-threaded lock manager instead of a single-threaded lock manager to have a better CPU utilization for a fairer comparison. We pre-populate each partition with 1 million key-value pairs before starting the experiments. Similar to 2PL, all keys in one transaction are performed on the same partition for best performance. Meanwhile, we generously add extra 4 CPU core resources for multiple lock managers and sequencer threads for each experiment run. In this implementation, transactions are generated in-place on the servers similarly to Rolis.

In Figure 12, we can observe that throughput of 2PL and Calvin is scaled correspondingly as we expected in Section 6.2 due to our perfect partition setting. However, 2PL only can achieve 137 K throughput with 28 partitions. The reasons are twofold: (1) 2PL implementation is a client-server architecture (2) Rolis’s OCC-based implementation has more

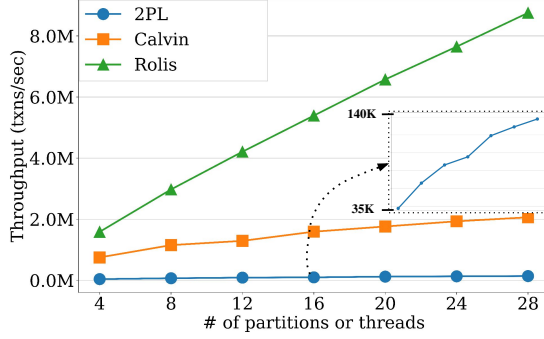


Figure 12. Comparisons with traditional software implementations: throughput on YCSB++ benchmark.

advantages over 2PL implementation due to the low contention setting. Rolis can have a much higher throughput than Calvin and 2PL. For Calvin, it needs a central sequencer to determine the order for a batch of transactions before they start execution which is expensive compared to Rolis. Similar to Calvin, 2PL also needs extensive coordination among replicas and holds all locks before transactions execution.

6.4 Comparison with kernel-bypass systems

We find it hard to compare to hardware-optimized systems in a real evaluation. Many of the state-of-art systems are not open-sourced (i.e., FaRM [11]) or depend on special platform features (i.e., DrTM+R [5] depends on Intel’s Restricted Transactional Memory). Meerkat [39] is the only system we found that runs with kernel bypass (DPDK) enabled NICs to compare with. Meerkat targets scaling multi-core performance in replication. It has an advanced fast quorum-based replication and transaction protocol. The replication and transaction are mixed so it faces the same problem we pointed out that the long latency of replication may compromise performance. Meerkat uses DPDK to reduce the message latency to resolve this issue.

When we try to deploy Meerkat we find that it is a good example of why a pure software solution like Rolis would have an advantage in maintenance and portability. In our experience, Meerkat does not run out-of-box on Azure because of the RPC Meerkat uses, eRPC, relies on a specific version of NIC driver. Azure happens to not support it. A newer version eRPC works on Azure without that driver, but the API changes in it make it incompatible with Meerkat. In short, to bring Meerkat into running, we spent several extra weeks investigating the issues and made code-level changes to upgrade the eRPC/DPDK library in Meerkat.

For Meerkat tests, we pre-load the entire database with 1 million data items per CPU core to keep the contention level constant as we increase the numbers of cores.

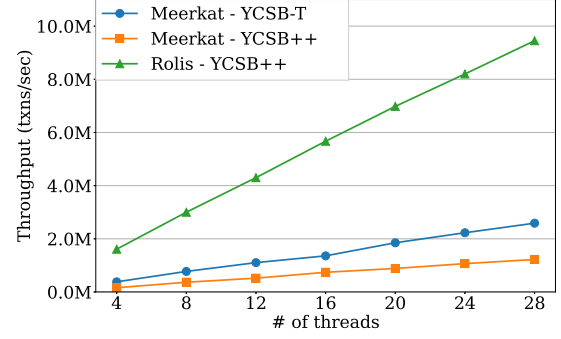


Figure 13. Comparisons with Meerkat: throughput on YCSB-T and YCSB++ benchmark.

DPDK on Azure we use offers a faster user-space packet processing framework by leveraging the advantages of high-performance NIC with FPGA. Our experiments show the latency between two VMs in the same cluster is ~ 30 us in the 99.9th percentile. Figure 13 shows Meerkat scales to 32 threads and 3.02M transactions per second on the YCSB-T benchmark which can match numbers reported in the Meerkat paper. Meerkat can achieve 1.35M transactions per second up to 32 cores on the YCSB++ benchmark. In numbers, Rolis can achieve up to 6 \times throughput of Meerkat but note that it is not exactly an apple-to-apple comparison. A major reason is that Meerkat uses networked clients while Rolis’s clients are embedded with the same process as the server. But still, the high performance of Rolis gives it much room for tolerating the extra overhead of adding actual networked clients. In theory, adding client-server communication will cause $\sim 1/3$ overhead for a 3-replica setup, as that the leader needs to communicate with the client in addition to the replicas (including itself).

6.5 Failure recovery

We conduct an experiment on replica failures with 3 replicas deployed in the same data center. In each run, we use fully loaded 4, 8, 16 threads respectively to test failure recovery under different workloads. In the experiment, we kill the leader replica after the system runs for 10 seconds and observe the system’s recovery. The recovery time consists of three parts: timeout through heartbeats, leader election, and replaying transactions in the old epoch. Recovery time is largely determined by the timeout through heartbeats. We set a relatively high timeout, 1 second, to avoid false positive detection and the following expensive recovery process.

After the old leader is killed, a new leader will be elected among alive follower replicas and continue processing requests. As shown in Figure 14, we track the system’s throughput in a 100 ms interval and observe a drastic drop at 10 seconds because the leader replica is killed at this moment, and then the election Paxos component on the follower replica

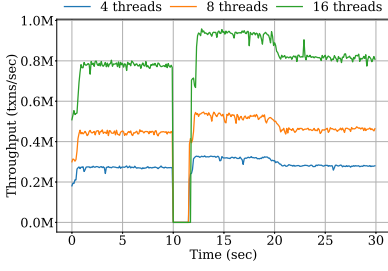


Figure 14. Failover test with timeout set to 1s (TPC-C)

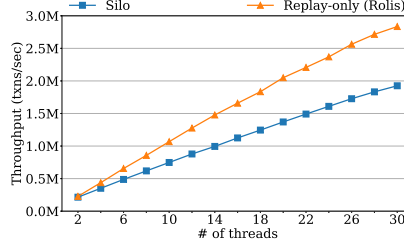


Figure 15. Compare the cost of replay-only in the Rolis with Silo (TPC-C)

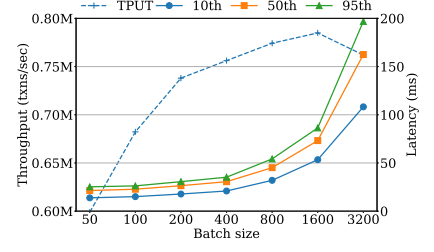


Figure 16. Latency and throughput with different batch size (16 threads, TPC-C)

loses heartbeat from the leader and starts a new round of the election to elect a new leader replica from alive follower replicas after a timeout, which blocks the system approximately for 1.5–2 seconds. Watermark cannot advance during this downtime.

After a new leader is elected, the throughput of the system quickly climbed to a peak that is higher than before as the system is trying to commit transactions queued up during the crash. After that, at around 20s, the system returns to a stable level of throughput, which is still slightly better than before the crash happens. This is because the system running at 2 replicas costs less in network communication than running at 3 replicas.

6.6 Silo vs replay-only

In this experiment, we conduct an evaluation of Silo versus replay-only in which we evaluate the throughput of replaying transactions on follower replicas with watermark control and Paxos disabled. This helps us understand the performance of the replay module, and especially whether it would be a bottleneck with the hacks described in Section 5. In this experiment, we pre-generate transaction logs from an independent Silo run, and then load these logs into replay threads’ memory.

Figure 15 shows that the throughput of replay-only at 32 cores is 2.88M transactions per second, which is 47.2% better than the original Silo’s execution. This improvement is mainly because the replay only processes write-set and ignores read-set, which saves much the cost in the workload. The throughput of replay-only scales very well because there are no complicated synchronizations and every key-value update can work independently and in parallel. This result concludes that the replay is not the bottleneck in the system.

6.7 Batch size factor analysis

Latency and throughput on batching is always a trade-off on systems: a larger batch increases the level of parallelism but introduces longer latency. To understand the impact of batching and batch size pickup strategy, we conduct an

experiment varying the batch size on the TPC-C benchmark where we use 16 database worker threads.

Figure 16 shows how changing the batch size would affect Rolis’s throughput and latency. Compare to a relatively small batch size (50), batching with the size 1600 can increase the system throughput by 31.2%. The throughput grows fast as we increase the batch size up to 400, then it keeps increasing at a slower speed until the batch size is 1600. After that, increasing batch size has an opposite impact on the throughput.

We also display the latency at the 10th percentile, 50th percentile, and 95th percentile, respectively as we increase the batch size. Latency of a transaction in this experiment is defined as the duration between the time the client issues a transaction and the time this transaction is released to the client. Figure 16 shows that the latency is increasing as the batch size grows as we expected because higher batch size implies that Rolis should take a longer time to commit a log entry in a Paxos stream and a longer time to allow a transaction to go beyond the advanced watermark to be ready for the replay. Rolis achieves a median latency of 108.2 ms and a 95th percentile latency of 196.8 ms with a high batch size of 3200.

Based on these results, we choose the batch size 1000 for the TPC-C benchmark and 10,000 for the YCSB++ benchmark in other experiments for its good throughput performance and relatively low latency.

7 Related work

We review related works in the literature on transactional, replicated, multi-core systems from the following aspects.

Single-machine multi-core transactional systems. Recent works on optimizing multi-core transactional systems performances primarily focus on multi-core scalability in throughput. A large body of them (Silo [41], Doppel [33], TicToc [48], IC3 [43], Hekaton [10], CoroBase [15], ERMIA [22], Cicada [26]) focus on optimizing concurrency control to build multi-core friendly databases. Zen [27] and

FOEDUS [23] further discuss optimizations targeting non-volatile memory. All these works are complementary to Rolis as they can potentially replace the local database of Rolis.

While these systems scale well with multi-core systems, they cannot be easily extended to a replicated environment because multi-core execution is inherently non-deterministic, and keeping replicated systems consistent requires each replica reflects the same order of execution which is a significant challenge. Rolis uses multiple Paxos streams and the watermark tracking mechanism to address this challenge.

Traditionally, single-machine databases consider improving reliability by supporting checkpoints to disks and recovery from them after rebooting from a crash. This provides a weaker fault-tolerance guarantee than a replication-based solution like ours in that it does not deal with the network asynchrony and it usually takes a much longer time to recover. For example, a well-optimized checkpoint approach, SiloR [50], needs several minutes to recover a Silo instance.

Replicating transactional systems. There is a line of transactional systems [7, 13, 28–30, 32, 40, 49, 52] that are optimized to achieve performance and provide availability through the replication [46]. Eris [25] and Harmonia [52] both exploit programmable switches by moving concurrency control to switches or detecting read-write conflicts in the network to improve distributed transactions. Some use more sophisticated replication protocols, such as inconsistent replication [49] and asymmetric replication [30] protocols to provide fault tolerance. Deterministic databases [28, 40] are able to efficiently run transactions across different replicas without coordination overhead. Compare to these systems, Rolis presents a much higher throughput. The major reason is that Rolis is designed to duplicate a higher baseline—the multi-core single-machine database—to a replicated setting, and it does so by using an execute-then-replay model and masking away the high network by intensive pipelining.

A recent popular approach [5, 11, 20, 39] to improve multi-core scalability is using kernel bypass network. Kernel bypass abstractions (e.g., DPDK, RDMA) are able to provide both low latency and high throughput. Similar to kernel bypass, NetChain [19] can also significantly reduce the latency of Paxos nodes coordination by caching key-values stores in switches. DrTM+R [5], FaRM [11], and FaSST [20] use kernel bypass to achieve scalability of distributed transactions with primary-backup replication [9, 35] by log shipping. However, they all mix transactions execution and replication protocol together to commit transactions where the primary has to receive all ACKs from backups before committing transactions. Unlike those works, Rolis does not rely on the low-latency provided by advanced hardware for high-throughput. Instead, it decouples transactions execution and replication protocol carefully to minimize the impact of replication as much as possible.

Deterministic execution and replay. Previous works [8, 14, 21, 38] have proposed replicating the scheduling information on the leader to all replicas to make the parallel execution (replay) on followers deterministic. Rex [14] is a multi-core friendly Paxos-based replication system. It allows concurrent execution on the leader replica freely while recording the non-deterministic decisions in causal ordered traces. Then follower replicas follow agreed-upon traces by making the same non-deterministic choices in a concurrent replay to reach the same consistent state as the leader. However, as Rex pointed out, there is a great challenge in replicating a transactional database using this way because multi-core transactions require too much locking, which will cause substantial overhead in tracking the synchronization orders between threads. Eve [21] lets replicas execute commands in parallel speculatively, then replicas verify whether they need to rollback and re-execute the commands in case of inconsistencies with high cross-core coordination overhead. Rolis addresses this challenge by choosing not to track the orders between locks, but to track the serialization order at the transaction level at the certification phase of a transaction. This avoids unnecessary coordination overhead in the transaction execution and in the replay.

Scalable-Replay [38] proposed a primary-backup replay-based scheme for replicating a multicore database of Ernia [22]. After executing a transaction, the primary node sends transactions to all backup nodes for replay using a customized multi-version engine on the backups so that transactions can read and write different versions of rows concurrently. This work has comparable throughput to Rolis. The major advantages of Rolis are two-folded. First, Rolis is a consensus-based approach that deals with network asynchrony and hence provides a stronger fault-tolerance level than a traditional primary-backup approach that assumes synchrony or requires a failure detector. Second, Rolis has a much shorter failover time: Scalable-Replay needs several minutes to do a data conversion on the backups before serving new requests because of different storage structures on the primary and backups, while Rolis only needs 1.5–2 seconds for failover as we show in Section 6.5.

8 Conclusion

We have presented Rolis, a new software design to build a speedy and fault-tolerant in-memory transactional system. Rolis adopts an execute-replicate-replay model, and masks the high cost of replication via intensive pipelining. The system is carefully designed so that the execution, replication, and replay are scalable and have nearly zero coordination overhead across cores. Our evaluation shows that Rolis can scale well under different workloads on commodity machines, and achieve high throughput comparable to state-of-art kernel-bypass systems, and orders of magnitude better than conventional software systems.

References

- [1] M. K. Aguilera, N. Ben-David, R. Guerraoui, V. J. Marathe, A. Xytkis, and I. Zablatchi. Microsecond consensus for microsecond applications. In *Proceedings of USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, Nov. 2020.
- [2] B. Arun, S. Peluso, R. Palmieri, G. Losa, and B. Ravindran. Speeding up consensus by chasing fast decisions. In *2017 47th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*. IEEE, June 2017.
- [3] M. Burke, A. Cheng, and W. Lloyd. Gryff: Unifying consensus and shared registers. In *Proceedings of USENIX Conference on Networked Systems Design and Implementation (NSDI)*, Feb. 2020.
- [4] M. Burrows. The chubby lock service for loosely-coupled distributed systems. In *Proceedings of ACM Symposium on Operating Systems Principles (SOSP)*, Nov. 2006.
- [5] Y. Chen, X. Wei, J. Shi, R. Chen, and H. Chen. Fast and general distributed transactions using RDMA and HTM. In *Proceedings of the Eleventh European Conference on Computer Systems*, Apr. 2016.
- [6] B. F. Cooper, A. Silberstein, E. Tam, R. Ramakrishnan, and R. Sears. Benchmarking cloud serving systems with YCSB. In *Proceedings of ACM Symposium on Cloud Computing (SoCC)*, June 2010.
- [7] J. C. Corbett, J. Dean, M. Epstein, A. Fikes, C. Frost, J. J. Furman, S. Ghemawat, A. Gubarev, C. Heiser, P. Hochschild, et al. Spanner: Google's globally distributed database. *ACM Transactions on Computer Systems (TOCS)*, Aug. 2013.
- [8] H. Cui, R. Gu, C. Liu, T. Chen, and J. Yang. Paxos made transparent. In *Proceedings of ACM Symposium on Operating Systems Principles (SOSP)*, 2015.
- [9] B. Cully, G. Lefebvre, D. Meyer, M. Feeley, N. Hutchinson, and A. Warfield. Remus: High availability via asynchronous virtual machine replication. In *Proceedings of USENIX Conference on Networked Systems Design and Implementation (NSDI)*, Apr. 2008.
- [10] C. Diaconu, C. Freedman, E. Ismert, P.-A. Larson, P. Mittal, R. Stonecipher, N. Verma, and M. Zwilling. Hekaton: SQL server's memory-optimized OLTP engine. In *Proceedings of ACM International Conference on Management of Data (SIGMOD)*, June 2013.
- [11] A. Dragojević, D. Narayanan, E. B. Nightingale, M. Renzelmann, A. Shamis, A. Badam, and M. Castro. No compromises: Distributed transactions with consistency, availability, and performance. In *Proceedings of ACM Symposium on Operating Systems Principles (SOSP)*, Oct. 2015.
- [12] V. Enes, C. Baquero, T. F. Rezende, A. Gotsman, M. Perrin, and P. Sutra. State-machine replication for planet-scale systems. In *Proceedings of ACM European Conference on Computer Systems (EuroSys)*, Apr. 2020.
- [13] H. Fan and W. Golab. Gossip-based visibility control for high-performance geo-distributed transactions. *Proceedings of International Conference on Very Large Data Bases (VLDB)*, 2021.
- [14] Z. Guo, C. Hong, M. Yang, D. Zhou, L. Zhou, and L. Zhuang. Rex: replication at the speed of multi-core. In *Proceedings of ACM European Conference on Computer Systems (EuroSys)*, Apr. 2014.
- [15] Y. He, J. Lu, and T. Wang. Corobase: coroutine-oriented main-memory database engine. *Proceedings of International Conference on Very Large Data Bases (VLDB)*, 2020.
- [16] M. P. Herlihy and J. M. Wing. Linearizability: A correctness condition for concurrent objects. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, July 1990.
- [17] C. Hong, D. Zhou, M. Yang, C. Kuo, L. Zhang, and L. Zhou. Kuafu: Closing the parallelism gap in database replication. In *2013 IEEE 29th International Conference on Data Engineering (ICDE)*. IEEE, 2013.
- [18] TPC-C is an On-Line Transaction Processing Benchmark. <http://www.tpc.org/tpcc/>.
- [19] X. Jin, X. Li, H. Zhang, N. Foster, J. Lee, R. Soulé, C. Kim, and I. Stoica. Netchain: Scale-free sub-rtt coordination. In *Proceedings of USENIX Conference on Networked Systems Design and Implementation (NSDI)*, 2018.
- [20] A. Kalia, M. Kaminsky, and D. G. Andersen. FaSST: Fast, Scalable and Simple Distributed Transactions with Two-Sided (RDMA) Datagram RPCs. In *Proceedings of USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, Nov. 2016.
- [21] M. Kapritsos, Y. Wang, V. Quema, A. Clement, L. Alvisi, and M. Dahlin. All about eve: Execute-verify replication for multi-core servers. In *Proceedings of USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2012.
- [22] K. Kim, T. Wang, R. Johnson, and I. Pandis. Ermia: Fast memory-optimized database system for heterogeneous workloads. In *Proceedings of ACM International Conference on Management of Data (SIGMOD)*, June 2016.
- [23] H. Kimura. Foedus: Oltp engine for a thousand cores and nvram. In *Proceedings of ACM International Conference on Management of Data (SIGMOD)*, May 2015.
- [24] L. Lamport et al. Paxos made simple. *ACM Sigact News*, Nov. 2001.
- [25] J. Li, E. Michael, and D. R. Ports. Eris: Coordination-free consistent transactions using in-network concurrency control. In *Proceedings of ACM Symposium on Operating Systems Principles (SOSP)*, Oct. 2017.
- [26] H. Lim, M. Kaminsky, and D. G. Andersen. Cicada: Dependably fast multi-core in-memory transactions. In *Proceedings of ACM International Conference on Management of Data (SIGMOD)*, May 2017.
- [27] G. Liu, L. Chen, and S. Chen. Zen: a high-throughput log-free OLTP engine for non-volatile main memory. *Proceedings of International Conference on Very Large Data Bases (VLDB)*, Jan. 2021.
- [28] Y. Lu, X. Yu, L. Cao, and S. Madden. Aria: a fast and practical deterministic oltp database. *The Proceedings of the VLDB Endowment (PVLDB)*, Aug. 2020.
- [29] Y. Lu, X. Yu, L. Cao, and S. Madden. Epoch-based commit and replication in distributed OLTP databases. *Proceedings of International Conference on Very Large Data Bases (VLDB)*, Jan. 2021.
- [30] Y. Lu, X. Yu, and S. Madden. Star: Scaling transactions through asymmetric replication. In *The Proceedings of the VLDB Endowment (PVLDB)*, July 2019.
- [31] Y. Mao, E. Kohler, and R. T. Morris. Cache craftiness for fast multicore key-value storage. In *Proceedings of ACM European Conference on Computer Systems (EuroSys)*, Apr. 2012.
- [32] S. Mu, L. Nelson, W. Lloyd, and J. Li. Consolidating concurrency control and consensus for commits under conflicts. In *Proceedings of USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, Nov. 2016.
- [33] N. Narula, C. Cutler, E. Kohler, and R. Morris. Phase reconciliation for contended in-memory transactions. In *Proceedings of USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2014.
- [34] D. Ongaro and J. Ousterhout. In search of an understandable consensus algorithm. In *Proceedings of USENIX Conference on Annual Technical Conference (ATC)*, June 2014.
- [35] D. Ongaro, S. M. Rumble, R. Stutsman, J. Ousterhout, and M. Rosenblum. Fast crash recovery in RAMCloud. In *Proceedings of ACM Symposium on Operating Systems Principles (SOSP)*, Nov. 2011.
- [36] C. H. Papadimitriou. The serializability of concurrent database updates. *Journal of the ACM (JACM)*, Oct. 1979.
- [37] M. Poke and T. Hoefler. Dare: High-performance state machine replication on rdma networks. In *Proceedings of ACM Symposium on High-Performance Parallel and Distributed Computing (HPDC)*, June 2015.
- [38] D. Qin, A. D. Brown, and A. Goel. Scalable replay-based replication for fast databases. *Proceedings of the VLDB Endowment*, Sept. 2017.
- [39] A. Szekeres, M. Whittaker, J. Li, N. K. Sharma, A. Krishnamurthy, D. R. Ports, and I. Zhang. Meerkat: multicore-scalable replicated transactions following the zero-coordination principle. In *Proceedings of ACM European Conference on Computer Systems (EuroSys)*, Apr. 2018.

2020.

- [40] A. Thomson, T. Diamond, S.-C. Weng, K. Ren, P. Shao, and D. J. Abadi. Calvin: fast distributed transactions for partitioned database systems. In *Proceedings of ACM International Conference on Management of Data (SIGMOD)*, May 2012.
- [41] S. Tu, W. Zheng, E. Kohler, B. Liskov, and S. Madden. Speedy transactions in multicore in-memory databases. In *Proceedings of ACM Symposium on Operating Systems Principles (SOSP)*, Nov. 2013.
- [42] C. Wang, J. Jiang, X. Chen, N. Yi, and H. Cui. APUS: Fast and Scalable Paxos on RDMA. In *Proceedings of ACM Symposium on Cloud Computing (SoCC)*, Sept. 2017.
- [43] Z. Wang, S. Mu, Y. Cui, H. Yi, H. Chen, and J. Li. Scaling multicore databases via constrained parallel execution. In *Proceedings of ACM International Conference on Management of Data (SIGMOD)*, June 2016.
- [44] Z. Wang, H. Qian, J. Li, and H. Chen. Using restricted transactional memory to build a scalable in-memory database. In *Proceedings of ACM European Conference on Computer Systems (EuroSys)*, Apr. 2014.
- [45] X. Wei, J. Shi, Y. Chen, R. Chen, and H. Chen. Fast in-memory transaction processing using rdma and htm. In *Proceedings of the 25th Symposium on Operating Systems Principles*, Oct. 2015.
- [46] M. Wiesmann, F. Pedone, A. Schiper, B. Kemme, and G. Alonso. Understanding replication in databases and distributed systems. In *Proceedings 20th IEEE International Conference on Distributed Computing Systems*. IEEE, 2000.
- [47] YCSB open-source implementation. <https://github.com/brianfrankcooper/YCSB>.
- [48] X. Yu, A. Pavlo, D. Sanchez, and S. Devadas. Tictoc: Time traveling optimistic concurrency control. In *Proceedings of ACM International Conference on Management of Data (SIGMOD)*, June 2016.
- [49] I. Zhang, N. K. Sharma, A. Szekeres, A. Krishnamurthy, and D. R. Ports. Building consistent transactions with inconsistent replication. *ACM Transactions on Computer Systems (TOCS)*, Dec. 2018.
- [50] W. Zheng, S. Tu, E. Kohler, and B. Liskov. Fast databases with fast durability and recovery through multicore parallelism. In *Proceedings of USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, Oct. 2014.
- [51] S. Zhou and S. Mu. Fault-tolerant replication with pull-based consensus in MongoDB. In *Proceedings of USENIX Conference on Networked Systems Design and Implementation (NSDI)*, Apr. 2021.
- [52] H. Zhu, Z. Bai, J. Li, E. Michael, D. Ports, I. Stoica, and X. Jin. Harmonia: Near-linear scalability for replicated storage with in-network conflict detection. *Proceedings of International Conference on Very Large Data Bases (VLDB)*, Nov. 2019.