



C5: cloned concurrency control that always keeps up

Jeffrey Helt¹ · Abhinav Sharma² · Daniel J. Abadi³ · Wyatt Lloyd¹ · Jose M. Faleiro⁴

Received: 20 September 2024 / Revised: 22 December 2024 / Accepted: 16 January 2025 / Published online: 12 February 2025
© The Author(s), under exclusive licence to Springer-Verlag GmbH Germany, part of Springer Nature 2025

Abstract

Asynchronously replicated primary-backup databases are commonly deployed to improve availability and offload read-only transactions. To both apply replicated writes from the primary and serve read-only transactions, the backups implement a cloned concurrency control protocol. The protocol ensures read-only transactions always return a snapshot of state that previously existed on the primary. This compels the backup to exactly copy the commit order resulting from the primary's concurrency control. Existing cloned concurrency control protocols guarantee this by limiting the backup's parallelism. As a result, the primary's concurrency control executes some workloads with more parallelism than these protocols. In this paper, we prove that this parallelism gap leads to unbounded replication lag, where writes can take arbitrarily long to replicate to the backup and which has led to catastrophic failures in production systems. We then design C5, the first cloned concurrency protocol to provide bounded replication lag. We implement two versions of C5: Our evaluation in MyRocks, a widely deployed database, demonstrates C5 provides bounded replication lag. Our evaluation in Cicada, a recent in-memory database, demonstrates C5 keeps up with even the fastest of primaries.

Keywords Databases · Concurrency control · Replication

1 Introduction

Asynchronously replicated primary-backup databases are the cornerstones of many applications [4, 26, 29, 38, 68]. In these systems, after the primary executes a transaction, it sends the resultant writes to a set of backups. The backups apply the writes to reconstruct the primary's state and execute read-only transactions against their local state. To simultaneously execute writes and read-only transactions, a backup imple-

ments a cloned concurrency control protocol. In addition to providing availability if the primary fails, these protocols improve the database's performance: throughput is increased by serving reads from many backups and latency is reduced by serving reads from a nearby backup.

To reap these benefits without breaking overlying applications, a cloned concurrency control protocol must guarantee monotonic prefix consistency, where it exposes a progressing sequence of the primary's recent states to read-only transactions. This ensures backups never return values from states that did not exist on the primary, thereby helping maintain application invariants.

But monotonic prefix consistency makes no guarantees about how quickly writes replicate to the backup. In theory, they could be delayed indefinitely. To be reliable, a cloned concurrency control protocol must also guarantee bounded replication lag. Intuitively, a transaction's replication lag is the time between when its writes are first observable by reads on the primary and backup. By guaranteeing bounded replication lag, a cloned concurrency control protocol ensures transactions always appear promptly.

Guaranteeing bounded replication lag is important; significant lag has led to catastrophic failures. For instance, GitLab was unavailable for eighteen hours after a workload change

✉ Jeffrey Helt
jhelt@cs.princeton.edu

Abhinav Sharma
abhinavsharma@meta.com

Daniel J. Abadi
abadi@cs.umd.edu

Wyatt Lloyd
wlloyd@princeton.edu

Jose M. Faleiro
jmfaleiro.work@gmail.com

- 1 Princeton University, Princeton, New Jersey, USA
- 2 Meta Platforms, Menlo Park, California, USA
- 3 University of Maryland, College Park, California, USA
- 4 San Francisco, California, USA

caused such significant lag that replication stopped entirely. In the process of fixing the issue, user data was lost [19, 20]. Similarly, several times in past years, Meta routed all user requests away from a data center because too many of that location's backups had excessive lag.

To guarantee bounded replication lag, a cloned concurrency control protocol must apply writes with as much parallelism as was used by the primary's concurrency control protocol. But guaranteeing monotonic prefix consistency severely constrains the cloned concurrency control protocol, making it difficult to execute with sufficient parallelism. For example, consider two concurrent transactions with both conflicting and non-conflicting writes. If the primary employs two-phase locking [6], the non-conflicting writes can execute in parallel, and the commit order is determined by the lock acquisition order on the first conflicting write. Once their commit order is chosen, however, monotonic prefix consistency mandates that a backup's state reflects it. Thus, the cloned concurrency control protocol must ensure the transactions are serialized correctly, potentially constraining its parallelism.

In the past, slow I/O devices bottlenecked the primary and backup, dominating differences in parallelism. But low-latency persistent storage and large main memories removed this bottleneck, so the primary's concurrency control and the backup's cloned concurrency control protocols are now directly competing.

Existing protocols differ in how much parallelism they leverage while executing writes. On one end of the spectrum are single-threaded protocols [48, 57]. On the other are transaction- [24, 32, 45, 51] and page-granularity [4, 54, 68] protocols. The former execute non-conflicting transactions in parallel; the latter execute writes to different pages in parallel. Such protocols have the potential to keep up with similarly restricted primaries, e.g., single-threaded cloned concurrency control with a single-threaded primary. But no existing protocol can keep up with an unrestricted primary.

In fact, existing protocols cannot keep up with a primary that uses two-phase locking; there are workloads where such a primary always executes with more parallelism. Using these workloads, we prove neither class of protocols guarantees bounded replication lag. In turn, implementations of these protocols are not reliable because changes to the workload or the primary's concurrency control can suddenly lead to unbounded replication lag.

In this paper, we present C5, the first cloned concurrency control protocol to provide bounded replication lag. To always keep up, C5's insight is that the backup's protocol must execute writes at the same granularity as the primary's concurrency control protocol. Thus, its cloned concurrency control has commensurate constraints (C5) with the primary. Because the primary executes writes to non-conflicting rows in parallel, C5 uses a row-granularity protocol.

But row-granularity execution introduces several challenges. First, applying individual row writes to the backup's state can lead to permanent violations of monotonic prefix consistency, where the backup's state ceases to match the primary's. To avoid such violations, C5's scheduler calculates the necessary metadata for its workers to correctly order writes to each row. Second, row-granularity execution does not guarantee monotonic prefix consistency for read-only transactions because transactional atomicity and commit order are not necessarily respected. Imposing additional constraints on workers, however, could reintroduce replication lag. Instead, C5's snapshotter uses three progressing snapshots, ensuring reads observe a consistent state without constraining execution.

We show formally that a row-granularity protocol never imposes more constraints on the backup's execution than a valid concurrency control protocol imposes on the primary's. Thus, C5 can, in theory, always match the primary's parallelism.

In practice, however, row-granularity execution is necessary but not sufficient to provide bounded replication lag. Other bottlenecks, such as a slow scheduler, may get in the way. We thus implement two versions of C5, C5-MyRocks and C5-Cicada, to confirm it always keeps up. C5-MyRocks is backward-compatible and deployed in production at Meta. Making it backward-compatible, however, required some additional constraints to the parallelism in our design. We thus also implemented C5-Cicada, which faithfully implements our design (without additional constraints) and demonstrates C5 can keep up with a cutting-edge concurrency control protocol.

We compare C5-MyRocks and C5-Cicada to a state-of-the-art, transaction-granularity protocol [24]. While it keeps up on some workloads, unbounded replication lag is lurking nearby: Simple optimizations that improve the primary's throughput cause the protocol to lag. In contrast, our implementations always keep up.

In sum, this paper's contributions stem from our key insight that cloned concurrency control must have commensurate constraints with the primary: first, we prove neither a transaction-granularity nor a page-granularity protocol can always keep up with a two-phase locking primary; next, we prove a commensurate-granularity protocol has the potential to keep up with an unrestricted primary; finally, we describe such a protocol, C5, implement two versions of it, and demonstrate both always keep up in practice.

Further, Sect. 8 describes experience from deploying C5 at Meta. Our experience echoes our evaluation: The simple single-threaded cloned concurrency control that was previously deployed could often keep up with the primary, but large replication lag would be exposed by workload changes. The deployment of C5 eradicated these issues and led to noticeably better reliability.

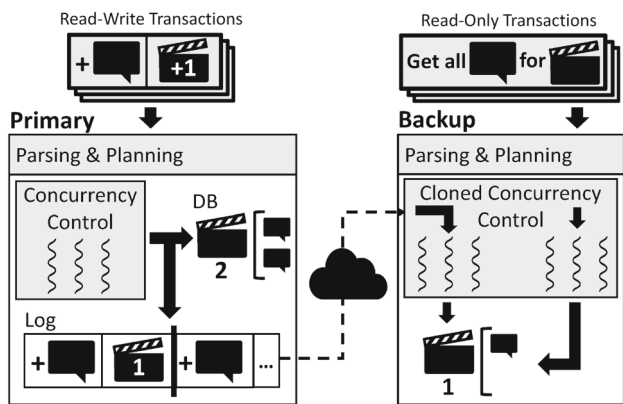


Fig. 1 Transaction processing in primary-backup

2 Background

This section provides the needed background on primary-backup databases, cloned concurrency control protocols, and their guarantees.

2.1 Motivating example

Throughout the paper, we use the following motivating example. Alice, Bob, and Charlie use a social media platform to share and comment on videos. The platform stores its videos and comments in a database. One table stores each video’s name and metadata, including a per-video comment counter; a second table stores each comment’s text and metadata. When a user comments on a video, an application server executes a transaction of two operations: it first inserts a new row in the comment table and then increments the video’s counter.

The platform replicates the primary’s database at a set of backups. The primary implements a concurrency control protocol, and each backup implements a cloned concurrency control protocol.

2.2 Primary-backup replication

Figure 1 shows an overview of the primary and backup’s processing as they execute transactions. For each operation in a read-write transaction, the primary parses it and plans its execution. Each plan, which may include row queries, local computation, and row writes (i.e., inserts, updates, and deletes), is then executed. For instance, to increment a video’s comment counter, the primary reads the counter’s current value from the video’s row in the video table, increments it, and writes the result back to the row. After all operations execute, the transaction commits by writing to the primary’s database and flushing a log of its changes to stable storage.

The primary then sends a copy of its log to the backup. The log reflects a total order of the writes applied by the primary, determined by the primary’s transaction commit order and the order of each transaction’s operations. The log includes, for each transaction, the written rows and metadata to demarcate its writes from those of others [13, 24, 32, 39, 45, 46, 48, 51, 55].

The backup’s cloned concurrency control protocol reads the operations in the log and schedules them for execution by worker threads, bypassing parsing and planning. The workers apply operations to the backup’s copy of the database. The protocol also executes read-only transactions using separate threads.

2.3 Monotonic prefix consistency

While an asynchronously replicated backup’s state inevitably lags, it is ideally otherwise indistinguishable from the primary. Intuitively, the backup should expose a progressing sequence of the primary’s recent states. This intuitive behavior is provided by many existing systems [4, 13, 24, 32, 39, 45, 46, 48, 51, 54, 55, 57, 68, 69]. We refer to this guarantee here as monotonic prefix consistency (MPC).

We define monotonic prefix consistency relative to the primary’s log of transactions. It comprises two guarantees: First, the backup’s state must reflect the changes of a contiguous prefix of transactions. Second, the sequence of states exposed to read-only transactions must reflect prefixes of monotonically increasing length.

In our example application, MPC ensures read-only transactions never see a mismatch between the number of comments on a video and the video’s comment counter; each transaction’s changes appear atomically. Further, MPC ensures comments never seem to disappear. Once a comment becomes visible to a user, all future states exposed by the same backup will include it. Although beyond our scope here, MPC can be guaranteed across multiple backups using sticky sessions [63] or with client-tracked metadata.

Monotonic prefix consistency also maintains implicit application invariants. For instance, suppose Alice first updates her default video permissions to share her future videos only with Bob and then uploads a new video. To make these changes, a transaction first updates her default access control list to only include Bob, and a subsequent transaction adds the new video. An implicit invariant, that Charlie should not see the new video, is expressed by the order of the two transactions. MPC preserves such invariants because states always reflect contiguous prefixes of the log.

2.4 Bounded replication lag

Monotonic prefix consistency specifies a cloned concurrency control protocol’s correctness but does not clarify its perfor-

mance requirements. For instance, if Alice calls Bob after commenting on his video, Bob should ideally see her new comment by the time he receives her call. Given only MPC, the comment may be delayed at the backup for an arbitrarily long time.

We define replication lag as the time between when a transaction's changes are included in the state returned by the primary and backup (For the purposes of this paper, we assume the log is always delivered promptly to the backup). More precisely, we say a transaction T is included in the state returned by the primary or backup once either its writes or later writes are returned to reads. To include a transaction in the returned state requires the backup's protocol to do one of the following: (1) it can eagerly apply the transaction's changes to its copy of the database, making them visible to future reads without additional processing beyond that required to execute the read at the primary [4, 24, 32, 39, 45, 48, 51, 54, 55, 57, 68]; or (2) it can defer part of the execution of the transaction's changes until a corresponding read arrives [69]. For each T , we then define $f_p(T)$ and $f_b(T)$ as the real time when the primary and backup respectively include T in their state. For eager protocols, $f_b(T)$ is the first time at which an arriving read would see T . For lazy protocols, $f_b(T)$ is the first time at which an arriving read would see T , plus the additional time required to finish any deferred execution.

A cloned concurrency control protocol guarantees bounded replication lag if there exists some finite time L such that for all workloads W and for all transactions T in W , $f_b(T) - f_p(T) \leq L$ (Transactions and workloads are defined more precisely in Sect. 3.1). In practice, guaranteeing bounded lag ensures Bob never waits long to see Alice's comment.

2.5 Hardware trends and practical considerations

Section 2.3 described the need for monotonic prefix consistency. One straightforward mechanism to guarantee MPC is to process log records serially. Indeed, several widely used database systems, including MySQL [48] and PostgreSQL [57], employ serial log application. But serial log replay can leave a backup prone to unbounded replication lag.

Unbounded replication lag occurs because the primary can process transactions in parallel using multiple cores, but backups are restricted to effectively using a single core to process transactions. This parallelism gap is increasingly problematic due to the emergence of low-latency storage devices and the affordability of large main-memory capacities.

With high-latency storage devices and a paucity of main memory, primaries were bottlenecked by the high latency of I/O. This fundamentally gave backups an advantage over the primary, since the primary generally does not know in

advance what it needs to access until shortly before it makes the request. Meanwhile, backups have the complete list of pages that it will need to access in the log records. Backups can use this a priori knowledge to read ahead into the log and prefetch the appropriate database pages from persistent storage before applying log records [40, 42], thereby ensuring the database pages are already available in memory before they need to be accessed. Meanwhile, the primary needs to pay context switching or spinning costs waiting for the required pages to arrive. Thus, the backup's ability to take advantage of read-ahead mechanisms mitigates the impact of the parallelism gap because backups could effectively hide the cost of I/O on the critical path of log application.

With large main memories and decreasing I/O latency, read ahead provides diminishing returns. This is because primaries can maintain a larger fraction of their working set in memory. Further, even when data must be fetched from stable storage, the cost of a main-memory miss is significantly less burdensome due to the low latency of I/O. Therefore, read ahead no longer suffices in hiding the parallelism gap between primaries and backups.

As these hardware trends continue, it is increasingly important to address the primary-backup parallelism gap. One approach used by a variety of database systems, both old [32] and new [24], is to analyze transactions in the log for conflicts, allowing backups to execute non-conflicting transactions in parallel. Unfortunately, this paper shows that scheduling log records at the granularity of non-conflicting transactions is insufficient to guarantee at least as much parallelism as the primary.

Intuitively, this is because concurrency control protocols schedule transactions based on dynamic detection of conflicts. Under two-phase locking, for example, the primary can dynamically acquire locks on records as it executes a particular transaction; locks on a particular database object are only acquired when the transaction accesses the object. As a consequence, locks on objects accessed at the end of a transaction are held for a much shorter duration than those acquired on objects accessed at the beginning of a transaction. Applications routinely exploit this behavior to access highly contended records at the end of a transaction, minimizing the duration for which contended locks are held. Indeed, recent research proposes techniques to reorder a transaction's constituent statements such that accesses to contended records are performed as late as possible [72]. Transaction-granularity scheduling effectively undoes these application-level optimizations by acquiring all of a transaction's locks at the beginning of its execution.

Furthermore, the past decade has seen a proliferation of increasingly sophisticated concurrency control protocols for modern database systems [33, 49, 67, 70]. This makes addressing the primary-backup parallelism gap more challenging. Advances in concurrency control can increase

parallelism on primaries, but backups do not benefit because most concurrency control protocols do not guarantee equivalence to a predefined serial order. Indeed, applications often employ weak isolation levels [3], such as read committed, which permit even more flexible schedules than serializability. Any solution to the primary-backup parallelism gap must work despite these advances and be resilient to advances in future protocols. What is needed, therefore, is a log application mechanism that is provably guaranteed to exploit at least as much parallelism as the primary, regardless of concurrency control protocol and isolation level the primary uses.

Finally, parallel log application must also take into account that queries can only observe transactionally contiguous log prefixes (Sect. 2.3). Executing non-conflicting log records in parallel naturally creates holes in the applied log. In order to prevent unbounded query staleness, parallel log application must therefore ensure that the contiguous prefix of the log observable to queries is constantly advancing.

3 Unbounded lag in existing protocols

Guaranteeing bounded replication lag is challenging. To satisfy monotonic prefix consistency, the backup’s cloned concurrency control protocol must ensure the backup’s state converges to the primary’s. To accomplish this, existing protocols serialize conflicting writes [4, 24, 32, 39, 45, 48, 51, 54, 55, 57, 68, 69].

Serialization limits the backup’s parallelism. But to always guarantee bounded replication lag, the backup’s protocol must be able to match the parallelism used by an unrestricted primary’s concurrency control protocol on every workload. Otherwise lag can grow arbitrarily large in some cases.

Transaction- and page-granularity cloned concurrency control protocols are the current best approaches. The former assume logical logs, and the latter assume physical redo logs [4, 54, 68]. In transaction-granularity protocols, writes conflict if they modify the same row, and the protocol serializes transactions with conflicting writes [24, 32, 45, 51]. Page-granularity protocols serialize writes to each page [4, 54, 68]. Both, however, fail to guarantee bounded replication lag because for some workloads, a primary executes with more parallelism.

We show how a transaction-granularity protocol can lag by returning to our motivating example. Suppose Alice, Bob, and Charlie simultaneously comment on the same video. Figure 2 shows a primary and backup’s executions of the six resultant operations. The primary uses two-phase locking [6] and stored procedures (i.e., no parsing and planning). The backup implements a transaction-granularity protocol [24, 32, 45, 51].

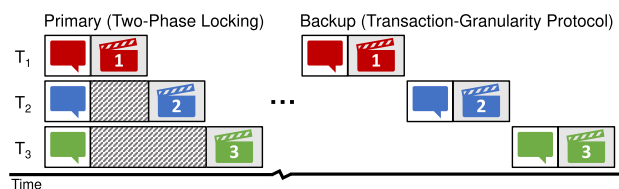


Fig. 2 Primary (2PL) and backup (transaction-granularity protocol) executions when three users comment on the same video. Diagonal lines depict waiting for a lock

On the primary, three threads insert rows in the comments table in parallel, but updates to the video’s comment counter are serialized by a row lock. On the backup, however, the transaction-granularity protocol serially executes all of the operations. Even if the backup uses different workers to execute each transaction (as shown), their execution is not faster than that of one worker. Thus, a fundamental gap exists between the parallelism available to the primary and backup, which can cause arbitrarily long lag.

Page-granularity protocols have similar issues. Concurrency control protocols using logical locking may allow concurrent transactions to update distinct rows residing on the same physical page [6, 53]. Such concurrency control can cause arbitrarily long replication lag when paired with page-granularity cloned concurrency control because writes that execute in parallel on the primary are serialized on the backup. Thus, a fundamental gap again exists.

In the rest of this section, we prove this problem is general to all transaction- and page-granularity protocols, and thus, no such protocol guarantees bounded replication lag.

3.1 Transaction-granularity protocols

System model A database \mathcal{D} stores sets \mathcal{K} of keys and \mathcal{V} of values. The database’s state is a mapping from \mathcal{K} to \mathcal{V} . A transaction T is an ordered set of operations (reads and writes) on individual keys. For simplicity, we assume each value is uniquely identifiable, so two identical transactions are, too. We define $\mathcal{R}(T)$ and $\mathcal{W}(T)$ as the sets of keys read and written by the operations in T . We define transaction arrival times at the primary and backup as $a_p(T)$ and $a_b(T)$, respectively. Finally, we define the real time at which a transaction is included in the primary and backup’s state as $f_p(T)$ and $f_b(T)$, respectively (as in Sect. 2.4).

We assume a primary-backup system where both have m cores. In isolation, the primary’s cores each execute an operation in $e > 0$ time units. The primary uses 2PL [6], so an operation may wait for a lock if there is a concurrent operation on the same key. If there are multiple conflicting operations, assume they are granted the lock in the order requested. To account for both eager and lazy cloned concurrency control protocols, assume the backup’s cores execute each operation

in $0 < d \leq e$ time units. We assume $d \leq e$ because backups often avoid some processing.

When the primary finishes executing transaction T 's operations, it records T 's writes in its log. $T_1 < T_2$ denotes T_1 precedes T_2 in the log. The log is then sent to the backup. For simplicity, we assume this occurs instantaneously.

A workload $W \in \mathcal{W}$ is a tuple $(\mathcal{T}, A_{\mathcal{T}})$ where \mathcal{T} is a set of transactions and $A_{\mathcal{T}}$ is a function from \mathbb{R} to finite sets of transactions $T \in \mathcal{T}$ representing the transaction arrival process at the primary. \mathcal{W} is the set of all definable workloads.

Transaction T 's replication lag is given by $f_b(T) - f_p(T)$. A cloned concurrency control protocol has finite replication lag for workload $W = (\mathcal{T}, A_{\mathcal{T}})$ if there exists some finite L such that for all $T \in \mathcal{T}$, $f_b(T) - f_p(T) \leq L$, and it guarantees bounded replication lag if it has finite replication lag for all $W \in \mathcal{W}$.

Definitions and assumptions. A transaction-granularity cloned concurrency control protocol guarantees that for all pairs of transactions T_1 and T_2 , if $\mathfrak{W}(T_1) \cap \mathfrak{W}(T_2) \neq \emptyset$ and $T_1 < T_2$, then all of T_1 's writes execute before any of T_2 's (This definition matches existing implementations [24, 32, 45, 51])

The proof below requires $m > \lceil \frac{e}{d} \rceil$, but this assumption is reasonable in practice. Server CPUs commonly contain at least 64 physical cores [28]. Thus, the assumption is not satisfied only if the backup executes operations more than 63 times faster than the primary. Stored procedures and sophisticated concurrency control protocols [33, 49, 67, 70] make such an advantage unlikely.

Theorem 1 *If $m > \lceil \frac{e}{d} \rceil$, then a primary-backup system using two-phase locking on its primary and a transaction-granularity cloned concurrency control protocol on its backup cannot guarantee bounded replication lag.*

Proof Assume we have a primary-backup system as described above, and assume to contradict that it guarantees bounded replication lag. Then there exists some L such that for all $W \in \mathcal{W}$, the system executes all transactions in W with replication lag $\leq L$.

We now construct a workload $W \in \mathcal{W}$ that includes at least one transaction with replication lag greater than L . Each transaction comprises $m \geq n > \lceil \frac{e}{d} \rceil$ writes, and there are $\lceil \frac{L}{nd-e} \rceil$ such transactions. Because $nd > e$, the number of transactions is well-defined. The first $n - 1$ writes of each transaction modify unique keys, and the last updates key k_0 . Define $A_{\mathcal{T}}$ such that a new transaction arrives at the primary every e time units, starting at time 0.

Because the primary uses 2PL, it executes the first $n - 1$ writes of each transaction in parallel but serializes their final updates to k_0 . For convenience, we index the transactions in the order they appear in the primary's log.

For the first set of m transactions, $f_p(T_0) = ne, \dots$, and $f_p(T_{m-1}) = (n + m - 1)e$. Because $m \geq n$, the core that executed T_0 is free when T_m arrives. Thus, T_m finishes e time units after T_{m-1} . In general, we see $f_p(T_i) = (n + i)e$.

The backup uses a transaction-granularity protocol, so it serially executes all writes in the workload. Thus, the backup finishes executing T_0 at $n(e + d)$. By construction, $nd > e$, so $f_b(T_0) > f_p(T_1)$. Thus, the backup immediately starts executing T_1 after T_0 . The same is true for all subsequent transactions. In general, we see $f_b(T_i) = ne + (i + 1)nd$.

Thus, in general, $f_b(T_i) - f_p(T_i) = ne + (i + 1)nd - (n + i)e = i(nd - e) + nd$. For the final transaction T in the workload, $i = \lceil \frac{L}{nd-e} \rceil$, and thus $f_b(T) - f_p(T) = \lceil \frac{L}{nd-e} \rceil (nd - e) + nd \geq \frac{L}{nd-e} (nd - e) + nd > \frac{L}{nd-e} (nd - e)$. Equivalently, $f_b(T) - f_p(T) > L$, a contradiction.

The result above shows that if the primary has sufficient cores, then a transaction-granularity protocol cannot guarantee bounded replication lag. To simplify our formalism, the proof assumes the primary uses 2PL and serializable isolation [5, 56]. We note three important extensions: First, the theorem applies if the primary uses weaker isolation [3] because it can only accelerate the primary.

Second, a similar result can be derived for some optimistic protocols [6, 36]. For example, a similar execution to the one in Fig. 2 is possible with multi-version timestamp ordering (MVTSO) [6]. Using MVTSO, the three transactions still insert comments in parallel. If they then read the comment counter, write its new value, and perform validation serially, in timestamp order, all three transactions will commit, and a fundamental gap will again exist. We leave the generalization of our formal framework to optimistic concurrency control to future work.

Third, because the above proof assumes $0 < d \leq e$, it also applies if there is primary-specific processing, such as parsing and planning. This additional processing can be accounted for by increasing e , and the theorem holds as long as $m > \lceil \frac{e}{d} \rceil$. If it does not, then the primary-backup system may be able to guarantee bounded replication lag but only because bottlenecks on the primary make it easy for an inefficient cloned concurrency control protocol to keep up, for example, if bottlenecks in logging, persistence, or log transfer restrict the primary's parallelism.

Despite these cases, however, solving replication lag remains urgent: First, there are many cases where these are not bottlenecks. For parsing and planning, many deployments use stored procedures. For persistence, mechanisms such as early lock release [9, 18, 30] and epoch-based group commit [8, 67] help decouple transaction throughput from I/O latency. Second, we expect advances in research and technology, such as non-volatile memory [27], to eventually remove these bottlenecks.

3.2 Page-granularity protocols

Databases historically assumed their data resided on disk. Thus, when persisting writes, they often locked data pages while flushing changes to disk [53, 57]. Page-granularity cloned concurrency control protocols [4, 54, 68] leverage this fact to match their primary’s granularity: Since writes to the same page are serialized on the primary, the backup can keep up with it despite also serializing writes to the same page.

The proof below, however, shows page-granularity protocols cannot keep up with an unrestricted primary. In practice, we believe they are more likely to keep up with an unrestricted primary than a transaction-granularity protocol because they impose fewer constraints on parallelism. For instance, writes may be spread across different pages, which is the best case for a page-granularity protocol. Yet our result shows they do not fundamentally solve the replication lag problem. Optimizations to the primary and changes to workloads can lead to unbounded lag.

3.3 Proof

We use the same system model as in Sect. 3.1.

Definitions and assumptions. We first define a function $p : O \times \mathbb{N}$ where O is the set of all possible operations. $p(o)$ is then the index of the page updated by operation o . A page-granularity protocol guarantees that for all pairs of operations o_1 and o_2 , if $p(o_1) = p(o_2)$ and $o_1 < o_2$, then the backup executes o_1 before o_2 (i.e., $o_1 < o_2$). For simplicity, we assume each operation updates one page. A logical operation that updates multiple pages can be modeled by splitting it into multiple physical operations.

Like the proof for transaction-granularity protocols in Sect. 3.1, the proof below constructs a workload where the primary executes with more parallelism than the backup. This is possible in some concurrency control protocols where concurrent transactions may update rows residing on the same page [53]. Because page-granularity protocols are efficient, however, the degree to which the primary’s parallelism can exceed the backup’s is more limited than when the backup uses a transaction-granularity protocol. The proof thus requires an extra assumption to ensure the primary executes with sufficient parallelism.

Let s be the index of the page that is the target of the most operations and $|S|$ be the number of operations that target this page. Formally, $s = \arg \max_i |\{o \mid p(o) = i\}|$ and $S = \{o \mid p(o) = s\}$. Recall we defined $e > 0$ and $d \geq 0$ as the number of time units required for the primary and backup’s cores to execute an operation in isolation. To ensure the primary can execute enough operations in parallel targeting the same page, the proof below requires $|S| > \lceil \frac{e}{d} \rceil$.

Like the comparable assumption in Sect. 3.1, this assumption is reasonable in practice. Page-granularity protocols cause problems for backups when there are multiple updates to the same page that are parallelized at the primary but serialized at the backup. The processing hardware limits the maximum parallelism possible at the primary to be at the granularity of a cache line: it is possible for different cache lines to be updated concurrently by different threads or processes. With a typical cache line size of 64B—more than enough to store a row with two integer columns—and a typical page size of 4 KiB, 64 rows can be stored on the same page, each on a different cache line. Thus, again, the assumptions will not be satisfied only if the backup can execute operations more than 63 times faster than the primary.

Theorem 2 *If $m > \lceil \frac{e}{d} \rceil$ and $|S| > \lceil \frac{e}{d} \rceil$, then a primary-backup system that uses two-phase locking on its primary and a page-granularity cloned concurrency control protocol on its backup cannot guarantee bounded replication lag.*

Proof Assume we have a primary-backup system as described above, and assume to contradict that it guarantees bounded replication lag. Then there exists some L such that for all $W \in \mathcal{W}$, the system executes all transactions in W with replication lag $\leq L$.

We now construct a workload $W \in \mathcal{W}$ that includes at least one transaction with replication lag greater than L . Let $n = \min(m, |S|)$. By assumption, $n > \lceil \frac{e}{d} \rceil$. The workload comprises a set of $\lceil \frac{nL}{nd-e} \rceil$ transactions of one write each. By assumption, $nd > e$, so the number of transactions is well-defined. Define the transaction arrival process $A_{\mathcal{T}}$ such that batches of n transactions arrive at the primary every e time units, starting at 0. Within each batch, each transaction updates a unique key such that $o \in S$. By assumption, $|S| \geq n$.

Transactions within each batch do not conflict, and by assumption, $m \geq n$. Thus, the primary, using 2PL, executes all transactions in a batch in parallel. For convenience, we assign indices to the transactions in the order they appear in the primary’s log.

In the first batch of transactions, $f_p(T_0) = \dots = f_p(T_{n-1}) = e$. At this moment, the primary’s cores are all idle and the next batch arrives. In general, $f_p(T_i) = \lfloor \frac{i+n}{n} \rfloor e$.

The backup uses a page-granularity protocol, so it serially executes all writes in the workload. At time e , n transactions arrive at the backup. It finishes executing T_0 at $e + d, \dots$, and T_{n-1} at $e + nd$. Since $nd > e$, T_n arrives at the backup before it finishes executing T_{n-1} , so the backup immediately starts executing T_n after T_{n-1} . This is also true of subsequent transactions. In general, $f_b(T_i) = e + (i + 1)d$.

We conclude by showing that for the final transaction T in the workload, $f_b(T) - f_p(T) > L$, and thus, the assumed system does not guarantee bounded replication lag. First, we

derive a lower bound on replication lag:

$$\begin{aligned}
 e + (i + 1)d - \left\lfloor \frac{i + n}{n} \right\rfloor e &\geq e + (i + 1)d - \left(\frac{i + n}{n} \right) e \\
 &= d + i \left(\frac{nd - e}{n} \right) \\
 &> i \left(\frac{nd - e}{n} \right)
 \end{aligned}$$

Letting $i = \left\lceil \frac{nL}{nd - e} \right\rceil$, we can further simplify the right-hand side: $\left\lceil \frac{nL}{nd - e} \right\rceil \left(\frac{nd - e}{n} \right) \geq \frac{nL}{nd - e} \left(\frac{nd - e}{n} \right) = L$

The result above shows that if the primary has sufficient cores and pages are sufficiently large, then a page-granularity protocol cannot guarantee bounded replication lag. By the same reasoning as in Sect. 3.1, this result implies three extensions: First, the result applies if the primary uses weaker isolation [3]; second, a similar result can be derived for some optimistic protocols [6, 36]; and third, the result applies to the (perhaps more realistic) cases where there is primary-specific processing, such as parsing and planning.

4 C5 design

C5 achieves two competing goals: it ensures bounded replication lag and guarantees monotonic prefix consistency for read-only transactions. To accomplish this, C5 comprises three components: a scheduler, a set of workers, and a snapshotter. The scheduler and workers ensure bounded replication lag by executing writes at a sufficiently fine granularity, and the snapshotter guarantees read-only transactions only see changes that are valid under monotonic prefix consistency. Together they implement C5’s row-granularity cloned concurrency control protocol.

As shown in Sect. 3, transaction- and page-granularity protocols fail to provide bounded replication lag because they cannot always execute with the same parallelism as the primary. The primary executes writes to different rows in parallel, so to provide bounded lag, C5’s workers execute writes at row granularity.¹

Unconstrained row-granularity execution, however, can lead to permanent violations of monotonic prefix consistency because conflicting writes may execute in the wrong order. For instance, suppose two transactions T and U each update rows x and y . If different workers execute the resultant writes (denoted $w_T[x]$, $w_T[y]$, $w_U[x]$, and $w_U[y]$), $w_T[x]$

¹ Some concurrency control protocols allow two transactions to update a row’s cells in parallel [25]. For ease of exposition, we assume they cannot, but rows are not fundamental to our design—C5 could be adapted for finer granularities.

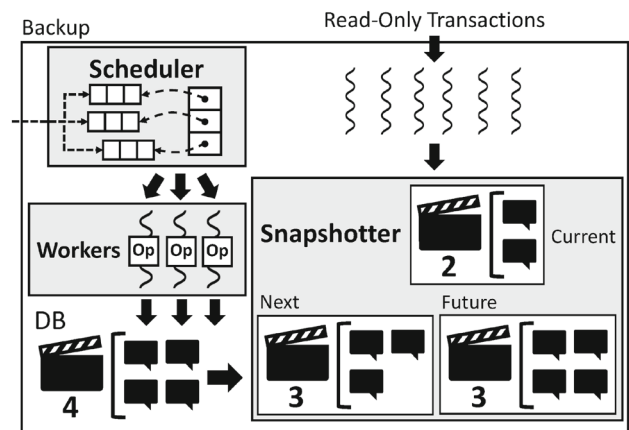


Fig. 3 C5’s scheduler, workers, and snapshotter

may finish before $w_U[x]$ and $w_U[y]$ before $w_T[y]$. If there are no further writes to these rows, the backup will forever reflect $w_U[x]$ and $w_T[y]$, violating transactional atomicity and thus MPC. C5’s scheduler helps avoid permanent consistency violations by constraining the workers’ execution. These constraints ensure writes to each row are applied in the same order as on the primary. Thus, each row reflects monotonically increasing prefixes of the log.

But per-row monotonicity is not sufficient for guaranteeing global monotonic prefix consistency. In the example above, a write to a third row z from a third transaction V may be scheduled after T and U but applied first. If this occurs, then a read-only transaction of rows x , y , and z would violate monotonic prefix consistency. Instead, C5’s snapshotter uses a set of three progressing database snapshots to allow uninterrupted execution of non-conflicting writes while guaranteeing MPC.

Figure 3 shows C5’s design. The scheduler orders writes and schedules them for execution by the workers. The snapshotter exposes a monotonic-prefix consistent view of the database to read-only transactions, which are executed by a separate set of threads. We now describe C5’s components in turn.

4.1 Row-granularity scheduling and execution

As described in Sect. 2.2, the backup continuously receives a log of operations from the primary, including the rows written by each operation and metadata to delimit transactions. To guarantee bounded replication lag, C5’s workers must execute individual row writes while obeying the constraints specified by the scheduler. To avoid permanent consistency violations, the scheduler logically constructs a FIFO queue for each row whose order reflects the order of the row’s writes in the log.

As the scheduler processes writes, it assigns each a sequence number, which reflects the write’s position in the

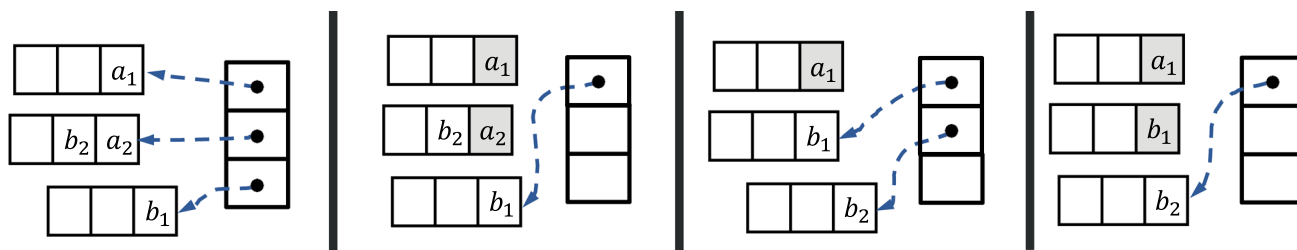


Fig. 4 Left to right shows the scheduler’s queues as two workers execute four writes. Per-row queues are shown horizontally and are ordered right to left. The scheduler queue is shown vertically and is ordered top

to bottom. Grey writes are being executed. This figure shows the design of C5 while Figs. 6 and 11 shows how we modify or refine this design in our MyRocks and Cicada implementations respectively

log. The scheduler then enqueues the write in the appropriate FIFO queue.

A write is safe to execute when it reaches the head of its FIFO queue and the prior head has finished executing. This assumes the backup receives the log of each row’s writes and the scheduler processes them in order (The log shipping subsystems in many commercial databases satisfy this assumption [48, 53, 57]). Given this, the scheduler is assured that when it processes a write, all conflicting writes that precede it in the log are either already in the queue or executing.

To keep replication lag small, the scheduler ensures workers execute safe writes promptly. To do so, the scheduler uses a FIFO queue to order the queues described above. To avoid ambiguity, we refer to a scheduler queue and per-row queues. Thus, a worker chooses the next write for execution by first removing the per-row queue at the head of the scheduler queue and then executing the write at its head. When the worker finishes executing the write, the per-row queue is reinserted into the scheduler queue.

To demonstrate how C5’s scheduler and workers operate, we return to our motivating example when Alice and Bob concurrently comment on the same video. Assume Alice’s transaction *A* commits first. It performs two operations: operation *a*₁ inserts one comment row, and *a*₂ increments the video’s comment counter. Bob’s transaction *B* performs comparable operations *b*₁ and *b*₂.

Figure 4 shows how two workers execute the four writes. The first panel shows the initial data structures after processing all operations. *a*₁ and *a*₂ then begin executing (second panel). (*b*₁ remains queued because there are only two workers.) *a*₂ finishes before *a*₁, so its corresponding queue is reinserted at the tail of the scheduler queue (third panel). Finally, *b*₁ starts executing (fourth panel). The workers continue until they execute all of the writes.

We now prove that row-granularity execution never imposes more constraints on the backup than any concurrency control protocol imposes on the primary.

4.1.1 Row-granularity execution can keep up

Both the primary’s concurrency control and the backup’s cloned concurrency control protocols can be viewed as functions from a set of logs to a set of sets of execution schedules. Given a log, the primary’s threads and the backup’s workers execute its writes according to one of the schedules in its image. We say a schedule is valid if the schedule of writes produces an equivalent database state as serially executing the writes in the log. In the remainder of this section, we only consider the set of valid protocols, those whose images contain only sets of valid schedules, and denote it as \mathfrak{S} . Note that a primary’s concurrency control protocol is always in \mathfrak{S} because the primary’s durability guarantees that its log, when serially executed, reproduces its state.

Let $w_T[x]$ denote a write to row x by transaction T . As before, $T \prec U$ denotes T precedes U in the log, and in a slight abuse of notation, let $w_T[x] \prec w_U[y]$ denote write $w_T[x]$ precedes write $w_U[y]$ in the log. Similarly, $w_T[x] < w_U[y]$ denotes $w_T[x]$ precedes $w_U[y]$ in an execution schedule. A row-granularity protocol guarantees that for all logs and all pairs of writes $w_T[x]$ and $w_U[y]$, if $x = y$ and $T \prec U$, then $w_T[x] < w_U[y]$ in all of its schedules.

Theorem 3 *Let $R \in \mathfrak{S}$ be a row-granularity protocol. Given a log, there does not exist a valid protocol $P \in \mathfrak{S}$ that imposes fewer constraints on its corresponding set of execution schedules than R .*

Proof (Sketch) Given a log and two transactions T and U such that $T \prec U$, R imposes one constraint on the possible executions of their writes: if $w_T[x]$ conflicts with $w_U[x]$, then $w_T[x] < w_U[x]$.

Assume to contradict there is a valid protocol P that does not impose the above constraint. Then in one of the resulting schedules, $w_U[x] < w_T[x]$. A serial execution of the log, however, always executes $w_T[x]$ before $w_U[x]$ because $T \prec U$ and thus $w_T[x] \prec w_U[x]$. As a result, this execution is

not equivalent to the serial execution of the log, contradicting that P is valid.

The proof shows that all valid concurrency control protocols, regardless of isolation level, must impose, at a minimum, the constraints imposed by a row-granularity protocol. If a backup employs such a protocol, then regardless of how much parallelism is exploited by the primary's concurrency control during its execution, an execution with an equal degree of parallelism is available to the backup. Thus, the proof shows row-granularity execution never hampers the backup's ability to keep up.

The proof, however, elides many practical details. As a result, while it shows that a backup using row-granularity execution can keep up in theory, it does not guarantee that a specific design or implementation will actually execute according to the necessary schedule to keep up in practice.

Regarding the design, a row-granularity scheduler, for instance, may impose more constraints than are strictly necessary for row-granularity execution, and these additional constraints may prevent the backup from keeping up in some cases. Regarding the implementation, a poor one, such as one that improperly uses concurrency mechanisms, may cause the scheduler to bottleneck the backup. Similarly, components outside the scope of a cloned concurrency control protocol may prevent it from keeping up. For instance, row-granularity execution may reduce cache locality when executing writes and could in theory make the backup's workers slower than the primary's threads.

To avoid overly complicating the formalism above, we thus do not claim to prove that a specific design or implementation guarantees bounded replication lag (We leave these theoretical investigations to future work). Instead, our experimental evaluations in Sects. 6 and 7.3 verify that the principles learned here translate to bounded lag in practice.

4.2 Snapshotter and read-only transactions

C5's snapshotter uses database snapshots to guarantee monotonic prefix consistency without blocking workers. To do so, it requires tight control over the creation and updating of the snapshots (as described below). The required operations are not backward-compatible with the storage engines in some commercial databases [14, 21, 47, 58] but can be implemented efficiently in many modern databases where workers can assign timestamps to their writes [10, 31, 35, 36]. We elaborate on how this difference affects the snapshotters of C5-MyRocks and C5-Cicada in Sects. 5.2 and 7.2.

The snapshotter creates new snapshots from the database. Logically, a snapshot is a sequence of writes and is initially empty. Workers apply writes (i.e., inserts, updates, and deletes) to a snapshot. Two snapshots \mathcal{S}_1 and \mathcal{S}_2 can be merged to produce a third \mathcal{S}_3 that reflects the writes applied

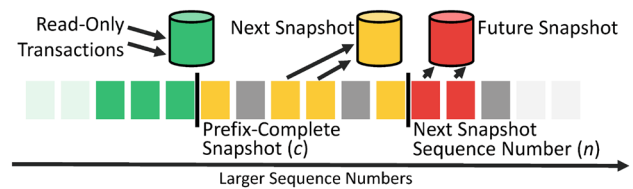


Fig. 5 C5's snapshotter. Writes in color (grey) have finished (not finished) executing

to both, with all writes in \mathcal{S}_1 ordered before those in \mathcal{S}_2 . Finally, the latest version of a row's value can be read from a snapshot.

The snapshotter uses the operations above to maintain three database snapshots, logically representing the current, next, and future. The current snapshot is initially empty, always prefix-complete, and serves read-only transactions. The next and future snapshots are initially empty. Workers only modify the next and future snapshots.

Figure 5 illustrates how the snapshotter incorporates a write into the current snapshot while maintaining monotonic prefix consistency. C5's snapshotter uses two sequence numbers to delimit the three snapshots. The current snapshot includes all writes up to sequence number c . All writes with sequence numbers between c and n (inclusive) update the next snapshot; all writes with sequence numbers greater than n update the future snapshot.

When all writes with sequence numbers between c and n finish executing, the current and next snapshots together form a new, prefix-complete snapshot. The snapshotter then merges them, and the result replaces the current. At the same time, it performs four additional operations: c is updated to reflect the new current snapshot; n is advanced; the next snapshot is replaced with the future snapshot; and a new future snapshot is created (We elaborate on how these steps are implemented in Sects. 5.2 and 7.2).

To satisfy monotonic prefix consistency, the snapshotter always aligns n with a transaction boundary. Thus, the next snapshot always reflects a set of complete transactions before being merged.

Because they execute against different snapshots, workers and read-only transactions execute in parallel. But to guarantee bounded replication lag, workers must be given higher scheduling priority than read-only transactions threads. To avoid starvation, we assume they execute on separate cores (beyond the m assumed in Sect. 4.1.1), or if they execute on the same cores, there are enough spare CPU cycles to process all read-only transactions.

5 C5-MyRocks implementation

C5-MyRocks was developed to solve replication lag at Meta, so backward compatibility and ease of deployment were

primary concerns. To remain backward-compatible with MyRocks (a fork of MySQL that uses RocksDB as its storage engine [13, 14, 48]), C5-MyRocks imposes some additional constraints on its execution, beyond those discussed in Sect. 4. In this section, we describe the implementation, highlighting how it leverages MyRocks's existing features [13, 48] and differs from our design.

5.1 Scheduling and execution

The scheduler leverages MyRocks's row-based-logging subsystem [13, 48], so C5-MyRocks does not require any changes to the primary or its log. For each operation, the log includes the set of written keys and corresponding values.

MyRocks's logging subsystem uses the one-thread-per-transaction execution model: the same worker executes all of a transaction's writes. To keep its changeset small (630 lines of C++ code), C5-MyRocks thus also enforces this constraint, since changing the system's execution model is known to be a complex and invasive operation [22].

To adapt our design to follow this model we need to ensure that once a worker executes one write in a transaction, it will then execute all later writes. We accomplish this by having the scheduler add metadata linking the writes in a transaction and modifying the execution model, so a worker now executes a transaction's writes in order.

More precisely, the scheduler processes a transaction in the log in three steps: First, it builds a linked list of its writes. Second, it adds each write to their per-row queues. Third, it adds an entry for only the transaction's first write in the scheduler queue.

Workers are now either unassigned or assigned to a transaction. An unassigned worker dequeues a write from the head of the scheduler queue, which will be the first write in a transaction. This logically assigns the worker to that transaction. First, it waits until the write reaches the head of its per-row queue, i.e., it is safe to execute, and then executes it. Then it follows the pointer to the next write in the transaction and repeats this process. If there is no next write, i.e., the transaction has been fully applied, the worker returns to being logically unassigned, so it will dequeue its next write from the scheduler queue.

Figure 6 illustrates C5-MyRocks's scheduling and execution for the two transactions described in the example in Sect. 4.1. The one-thread-per-transaction execution model results in several differences. First, b_1 is executed second in contrast with Fig. 4, which executed a_2 second. Further, when the second worker finishes executing b_1 , it will always execute b_2 next, after waiting until a_2 is executed and dequeued.

C5-MyRocks's one-thread-per-transaction execution model and having workers pick up transactions in commit order greatly simplified the implementation, but as demonstrated above, its implementation is more constrained than

our design. These constraints decrease the parallelism available in two ways compared to our design. First, workers must follow pointers to the next write in a transaction where they may need to wait for that write to be at the head of its per-row queue. In contrast, in C5's design, the worker would instead immediately execute another write from the scheduler queue.

This potentially decreases the amount of parallelism provided by each worker because a worker may block even though there is another write available for execution. The total amount of parallelism in the workload is unchanged, however, and additional workers can be used to execute the available writes in such cases.

Second, workers need to execute the writes of each transaction in order, e.g., a_1 and then a_2 in Fig. 6. In contrast, in C5's design, a transaction's writes can be executed in parallel, e.g., a_1 and a_2 in parallel in Fig. 4. This decreases the amount of total parallelism available to the backup. This decrease does not lead to unbounded replication lag, however, when the primary also uses the one-thread-per-transaction model because the primary's parallelism is limited in exactly the same way as the backup's. We expect the backup to only use the one-thread-per-transaction model when the primary does, so this component of our backward-compatible design has as much parallelism as its corresponding primary and thus can keep up.

5.2 Snapshotter and read-only transactions

The storage engines in some widely deployed databases [14, 21, 47, 58], including MyRocks, cannot easily implement the operations described in Sect. 4.2. Unfortunately, they are also complex, comprising tens to hundreds of thousands of lines of code [14, 21]. To again keep C5-MyRocks's changeset small, we opted to implement its snapshotter without requiring changes to RocksDB.

In MyRocks, snapshots are read-only and can only be taken of the database's current state. Neither workers nor the snapshotter have fine-grained control over which writes are included in a snapshot (e.g., by taking a snapshot as of some specified timestamp or version number). As a result, the snapshotter must impose some additional constraints on the workers to ensure the entire database is prefix-consistent when taking a new snapshot.

Instead of three snapshots, C5-MyRocks's snapshotter logically maintains two: It uses a current snapshot, which is always prefix consistent and used to serve read-only transactions. The next snapshot, however, is replaced by the database. c still tracks the writes included in the current snapshot.

To merge the current and next snapshots, the snapshotter performs the following: First, it chooses n , the sequence number of the last write to be included in the next snapshot. To ensure the merged snapshot is prefix-consistent, choosing

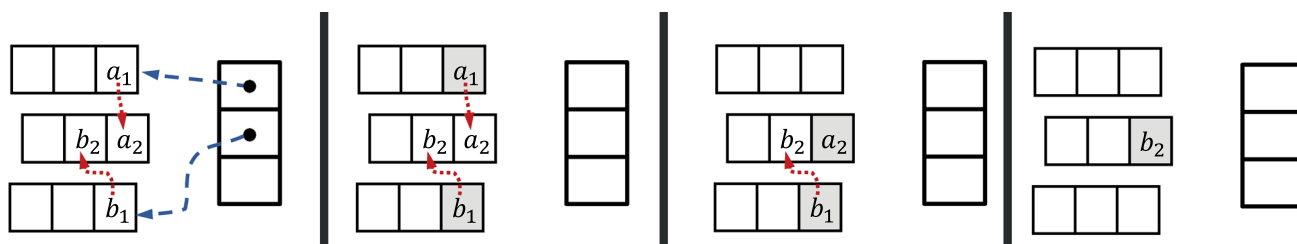


Fig. 6 Scheduling and execution in C5-MyRocks that follows MyRocks one-thread-per-transaction execution model. Blue arrows denote pointers from the scheduler queue to each transaction’s first

write that are dequeued by an unassigned worker. Red arrows are pointers linking each transaction’s writes, which are followed by a worker once it is assigned to a transaction. Grey writes are being executed

n also blocks workers from executing writes with sequence numbers greater than n until after the snapshot is taken (If the storage engine supports transactions, as in RocksDB [14], the workers only need delay committing these writes). Second, after all writes with sequence numbers between c and n (inclusive) execute, the snapshotter takes a new snapshot of the database and replaces the current one. c and n then advance as described previously. Advancing n also allows blocked workers to proceed with their writes.

If taking a snapshot is computationally expensive, the blocking above may lead to spikes in replication lag. To combat this, our implementation allows database administrators to tune the approximate snapshot frequency, I , in milliseconds. Because the storage engine may not be prefix consistent exactly every I milliseconds, the snapshotter advances n by an estimate of the number of writes workers will execute in the next I milliseconds.

By tuning I , administrators can ensure replication lag returns to satisfactory levels between snapshots provided lag can decrease between snapshots. This implies the backup must execute each write marginally faster than the primary. But we found this assumption reasonable in practice, and our evaluation further supports this.

6 C5-MyRocks evaluation

Our evaluation explores the following questions:

1. Does C5-MyRocks help engineers avoid potential disasters caused by optimizations of realistic workloads? (Sect. 6.1)
2. Does C5-MyRocks always keep up with the primary? (Sect. 6.2)
3. Does C5-MyRocks guarantee MPC for read-only transactions without causing unbounded replication lag? (Sect. 6.3)

Experimental setup. All experiments ran on the CloudLab Wisconsin platform [11] with three servers located in one

data center: one for load generation, the primary, and the backup. Round-trip times between machines were less than 100 μ S. Each machine had two 2.20 GHz Intel Xeon processors with ten cores each, hyper-threading disabled, 192 GB of RAM, a 10 Gb NIC, and a 480 GB SSD. Our code and experiment scripts are available online [2].

For each experiment, a fixed number of closed-loop clients executed read-write transactions at the primary. The log of writes was then sent to the backup and executed by the cloned concurrency control protocol’s workers. The number of clients and workers were set to maximize the primary and backup’s throughput, respectively, and there were always fewer workers than primary threads.

For experiments including both read-write and read-only transactions (i.e., Sect. 6.3), an additional set of closed-loop clients sent read-only transactions to the backup. The backup’s workers and read-only threads were pinned to separate cores, and their throughput is shown separately.

All results were from 120-second trials. We omit all data from the first and last 15 s of each trial to avoid experimental artifacts. Unless otherwise specified, we ran each experiment five times and report the median result.

To stress the cloned concurrency control, the primary used read committed isolation [5]; further, the log and MyRocks’s state on the primary and backup were kept in memory [14]. For all implementations, we used read-free replication and disabled MyRocks’s 2PL on the backup [13, 14] since the scheduler already prevents conflicting writes from executing concurrently. In all experiments, memory bandwidth and the network were not bottlenecks.

Workloads. We use three workloads. The first is TPC-C [66], an OLTP benchmark simulating an order-entry application. All experiments use one warehouse, so the database initially contains about 300,000 rows. The other two, insert-only and adversarial, are synthetic. In each, the database contains one table, initially with one row, with two integer columns: a primary key and a value.

Each transaction in the insert-only workload comprises a variable number of unique inserts. Each transaction in the

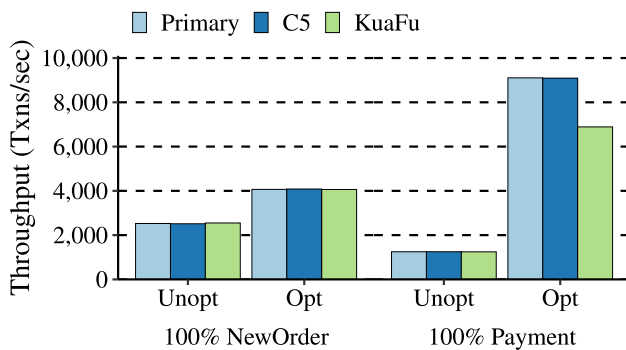


Fig. 7 Throughput on 100% NewOrder and 100% Payment before and after optimization

adversarial workload comprises a variable number of unique inserts and one update. The updates set the same row's value to a random integer, so all transactions conflict.

Baselines. KuaFu [24], a state-of-the-art, transaction-granularity cloned concurrency control protocol, is our baseline. KuaFu's protocol is nearly identical to MySQL 8's write-set-based parallel replication [45] and is strictly better than the database-granularity and epoch-based protocols used in earlier versions of MySQL [46, 48] and its variants [39]. For fair comparison, we re-implemented KuaFu in MyRocks.

6.1 C5-MyRocks prevents potential disasters

Because software and hardware improvements have accelerated the primary's processing, primary-backup systems using transaction-granularity protocols are brittle. Simple changes to a workload may cause unbounded replication lag. To demonstrate this problem, we use TPC-C [66]. While KuaFu [24] keeps up on the standard benchmark workload, simple optimizations and non-standard transaction mixes cause unbounded replication lag with the same protocol [24]. We discuss the optimizations and our results in turn.

We optimize two of TPC-C's transactions: the NewOrder and Payment transactions [66]. In both cases, we defer higher-contention operations as much as possible while preserving application semantics (Similar optimizations were observed in prior work [72]). In the NewOrder transaction, the highest contention write is the increment of the district's next order ID. In the Payment transaction, it is the update to the warehouse's balance [66]. Deferring these writes allows more parallelism on the primary.

Figure 7 shows the primary and backup's throughput while executing read-write transactions for a 100% NewOrder and a Payment workload, each before and after optimization (TPC-C's read-only transactions were not used in these experiments). For the NewOrder workload, the optimization

increases the primary's throughput from 2527 to 4067 transactions per second. For the Payment workload, the primary's throughput increases by over 700% from 1249 to 9105 transactions per second.

KuaFu keeps up with the optimized NewOrder workload. Data dependencies between operations limit how late the district row write can be deferred within each transaction and in turn, limits the primary's parallelism. But KuaFu cannot keep up on the optimized Payment workload; its throughput peaks at 6,889 transactions per second. Conversely, C5-MyRocks keeps up.

If the optimization to the Payment transaction were made to a production workload, KuaFu would cause significant replication lag, with 2216 transactions queuing at the backup every second. This rate is larger than the one that induced replication lag of nearly 2 h in production at Meta (discussed further in Sect. 8). On the other hand, C5-MyRocks thwarts such a disaster.

6.2 C5-MyRocks always keeps up

To validate that C5-MyRocks always keeps up, we measured the primary and C5-MyRocks's throughput using the insert-only and adversarial workloads. The two are on opposite ends of the contention spectrum: no transactions conflict in the former, while all do in the latter. We also present KuaFu's results.

Because all transactions are non-conflicting, the insert-only workload stresses the primary's concurrency control and backup's cloned concurrency control. Here, MyRocks's throughput is about 40,500 transactions per second. C5-MyRocks keeps up with the primary, indicating that its scheduling mechanisms have sufficiently low overhead. As expected, KuaFu also keeps up. With both protocols, incoming writes can be executed immediately.

To verify C5-MyRocks's scheduler is not a bottleneck, we ran the same experiment offline. We loaded the primary with inserts as above but delayed replication. Once all writes finished and the resultant log was transferred, we enabled C5-MyRocks's scheduler and workers. We used sufficient workers, so the scheduler was the bottleneck. C5-MyRocks's scheduler processed 95,683 transactions per second, more than double MyRocks's throughput.

Figure 8 shows each implementation's performance on the adversarial workload. We plot the backup's throughput relative to the primary's as we vary the number of non-conflicting inserts per transaction from 1 to 64. Every transaction updates the same row. Despite the high contention, the primary, using 2PL [6], executes the non-conflicting inserts that precede the conflicting update in parallel. Because all transactions conflict, KuaFu serializes them. Thus, the primary's advantage over KuaFu increases with the number of inserts. KuaFu's throughput drops from 70% to just 38% of the primary's. On

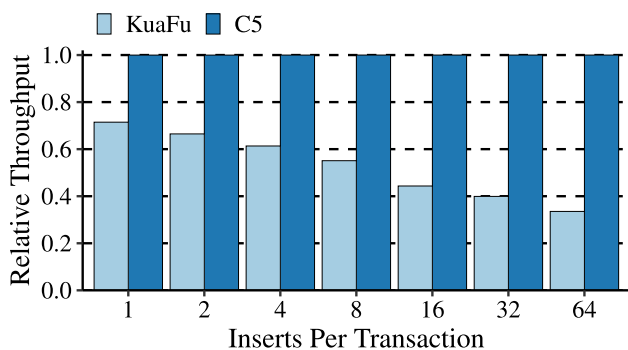


Fig. 8 Backup’s throughput relative to primary’s on adversarial workload as inserts per transaction increases

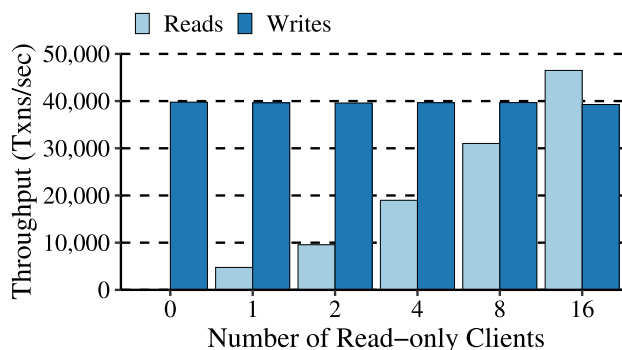


Fig. 10 Backup’s read-only and read-write transaction throughput as the read-only load increases

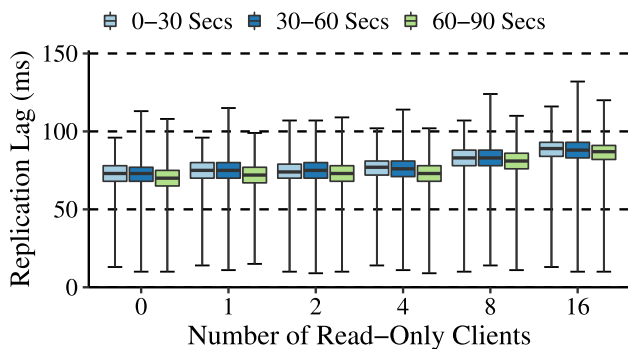


Fig. 9 Replication lag of read-write transactions with increasing number of read-only clients

the other hand, C5-MyRocks executes the non-conflicting inserts in parallel so always keeps up.

6.3 C5 serves reads with bounded lag

C5-MyRocks’s implementation blocks writes from committing while taking a snapshot. But it also exposes a parameter to tune the frequency of snapshots. With periodic snapshots, C5-MyRocks serves read-only transactions in parallel with writes, and thus, steady-state replication lag remains bounded despite additional load from read-only clients. To validate these claims, we measure replication lag as we increase the read-only load on C5-MyRocks.

Figure 9 plots the distribution of replication lag for read-write transactions over three 30-second periods with the insert-only workload. The whiskers show the minimum and maximum, and the boxes show the quartiles. We show one trial; the results from others were similar. For each read-write transaction, we measure replication lag as the difference between when it commits on the primary and when it is included in the current snapshot. Snapshots were taken every 10ms. Each read-only transaction executes a random point query on the table’s primary key; queries could select a nonexistent key. We vary the number of clients from 0 to 16.

Replication lag remains bounded in all cases and across all time periods. With 16 read-only clients, the median and maximum lag are about 88 ms and 135 ms, respectively. With zero clients, the median lag is about 73 ms. The latter is lower because it avoids contention on some of MyRocks’s internal data structures between the read-only threads and workers. We expect this contention can be removed with further optimizations.

Figure 10 plots the backup’s throughput of read-only and read-write transactions for the same experiment. C5-MyRocks’s throughput always matches the primary’s (Differences in the primary’s throughput from prior experiments is due to variations in MyRocks’s performance across trials). Further, C5-MyRocks isolates workers from read-only transactions. With steady write throughput, read-only throughput increases from 4755 transactions per second with 1 client to 46,500 transactions per second with 16 clients.

7 C5-Cicada

Unlike C5-MyRocks, our implementation of C5 in Cicada [36], an in-memory multi-version database, is faithful to our design. In this section, we first provide background on Cicada. We then describe the implementation of C5-Cicada’s scheduler, workers, and snapshotter. We conclude with our evaluation, which demonstrates C5 can keep up with a modern concurrency control protocol. Our code and experiment scripts are available online [1].

7.1 Cicada background

Cicada’s multi-version storage engine supports reads, inserts, updates, and deletes. The storage engine is implemented as an array indexed by an internal row ID (Externally meaningful keys are mapped to row IDs through indices). Array entries are linked lists of row versions in descending timestamp order. In addition to a pointer to the next oldest version, a

row version contains data, a status, and read and write timestamps. The latter three are used by Cicada’s concurrency control protocol.

Cicada uses a variant of multi-version timestamp ordering [6, 36]. Each client thread maintains a local clock (Cicada does not support networked clients). The local clocks are loosely synchronized and individually return increasing values.

A client uses its clock to assign a unique timestamp to each transaction. As the transaction executes, each write creates a new row version, and the transaction’s timestamp becomes the version’s write timestamp. Further, reading a version updates its read timestamp to the max of the transaction’s timestamp and the version’s current read timestamp. Cicada uses the timestamps to check if a transaction can commit under serializability [56]. Ordering transactions by their timestamps yields a valid serial schedule.

Logging and replication. Because Cicada does not support networking, logging, persistence, or replication [36], we emulate primary-backup replication on one server.

To this end, we implement a minimal prototype logger to allow replay of the primary’s writes on the backup. The primary only writes logs to memory. After execution and validation but before committing, each client thread logs its changes to a per-thread log. The per-thread logs are coalesced into a single, totally ordered log before the backup’s scheduler, workers, and snapshotter start.

The log is divided into fixed-size segments, each backed by a 2 MiB huge page [37]. Each segment’s header indicates the number of log records it contains. For simplicity, the logger ensures transactions never span segment boundaries.

A client creates a log record for each write in a transaction. Each record contains the following: a table ID, a row ID, the write’s timestamp, and a full copy of the row version. Further, the log contains two metadata fields to be used by the scheduler, as described below: First, each segment header contains a Boolean `preprocessed` flag. Second, each record contains a 64-bit `prev_timestamp` field.

Each log record is split into a header and data, with the former containing everything except the row version. Headers are written from the beginning of the segment, and data are written from the end. Co-locating all headers in a segment was critical for reducing the amount of data the scheduler (discussed further below) needed to process and prevented its throughput from being limited by memory bandwidth, especially in the face of concurrent workers.

7.2 Implementation

Scheduling and execution. For simplicity, the description below assumes we are replicating one table. The implemen-

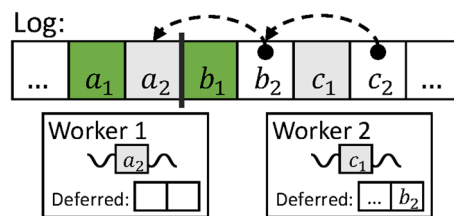


Fig. 11 Scheduling and execution for C5-Cicada that is optimized to enable keeping up with Cicada’s fast primary. The log is shown growing from left to right with green writes already executed and grey writes currently being executed. The thick black line shows how the log is segmented between worker 1 on the left and worker 2 on the right

tation supports multiple tables using a queue for each table and row ID pair.

Because dynamically allocating per-row FIFO queues would prevent the single-threaded scheduler from matching Cicada’s high throughput, it instead (logically) embeds the per-row FIFOs in the log. More specifically, it sets each log record’s `prev_timestamp` to the timestamp of the write to the same row that precedes it. To do so, it maintains a map of row IDs to write timestamps, initially zero. Then to process a record, the scheduler first reads the value from the map using the log record’s row ID, updates the record’s `prev_timestamp`, and finally updates the map’s value with the record’s write timestamp. After the scheduler processes all of a segment’s records, it sets its `preprocessed` flag to true.

Workers are assigned to log segments in a round robin fashion. Once the scheduler processes a segment, the assigned worker starts executing its writes, one for each log record, in three steps: First, it uses the record’s `prev_timestamp` to see if the write is safe to execute. If `prev_timestamp` is equal to the write timestamp of the row version at the head of the storage engine’s version list, then the write should be executed next. Otherwise it is deferred. Second, if the write is safe, the worker allocates a new row version and copies in the necessary data from the log. Third, the new version is installed at the head of row’s version list. Each worker maintains a local FIFO of deferred writes and periodically (after each segment) re-checks them to see if they are now safe to execute.

Figure 11 illustrates the scheduler and workers’ data structures as they execute three transactions of the type described in Sect. 4.1. The bold black line delimits the two log segments, and the scheduler has finished embedding the per-row queues in both. Workers one and two have started processing the left and right segments, respectively, and the former already executed a_1 and the latter b_1 . Worker two has deferred executing b_2 since worker one is still executing a_2 . Instead, worker two executes c_1 , which is already safe to execute. As shown, C5-Cicada thus maintains a distributed, approx-

imate version of the scheduler queue described in Sect. 4.1 comprising the log and each worker's deferred queue.

We refined the general design's global scheduler queue into this approximate, distributed queue to ensure our implementation could keep up with the fast Cicada primary. Logically embedding the per-row queues allows the scheduler to keep up by avoiding the need to allocate per-row queues, while the use of log segmentation and per-worker deferred queues allows execution to keep up by avoiding memory contention from multiple workers accessing a global data structure.

Snapshotter. Cicada's storage engine can efficiently implement the operations described in Sect. 4.2 because workers can explicitly assign timestamps to their writes. This allows workers to write to specific snapshots. Further, read-only transactions can execute against the current snapshot simply by using the sequence number c as a timestamp. Their reads will then reflect any previously executed writes with lesser timestamps. The storage engine thus logically contains the current, next, and future snapshots.

This simplifies C5-Cicada's snapshotter. To merge the current and next snapshots, it simply advances c to n , and replacing the next snapshot with the future snapshot and creating a new future snapshot occur implicitly when c and n advance.

Before advancing c to n , however, the snapshotter must guarantee all writes with timestamps less than or equal to n finish executing. To do this, it cooperates with the workers.

Worker i maintains a local variable c_i as one less than the timestamp of the write it most recently executed. Because log records are ordered by timestamp and each worker processes segments in log order, the worker can guarantee it will never execute another write with a timestamp less than or equal to c_i (assuming no writes are deferred). Each c_i is thus an upper bound on c . When a worker defers a write, it also defers updating c_i until the write executes.

The snapshotter's implementation is simple: in a separate thread, it periodically calculates a new n as the minimum across all c_i and then advances c to n . Since each c_i has only one reader and one writer, no coordination or atomic instructions are needed. With $\times 86$'s total store order, the snapshotter may read a stale copy of a worker's c_i , but this does not violate correctness.

7.3 Evaluation

Our evaluation of C5-Cicada explores whether our design can keep up with advanced concurrency control protocols.

Experimental setup. Experiments ran on the same machines as described in Sect. 6. For the primary, each thread is pinned

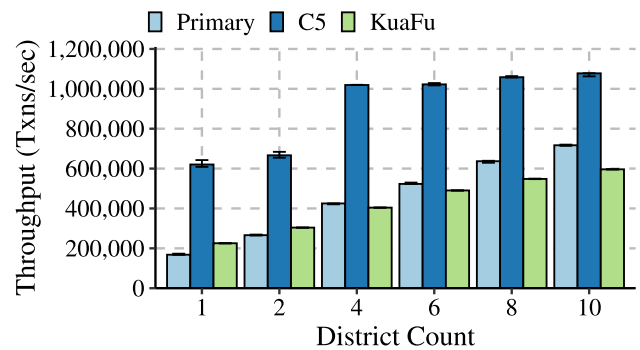


Fig. 12 Primary and backup's throughput on 50%-50% NewOrder-Payment with increasing districts

to its own physical core, and for the backup, the same is true of the scheduler, workers, and snapshotter.

As mentioned in Sect. 7.2, our C5-Cicada implementation emulates primary-backup replication. As a result, unlike in our prior evaluation, the primary and backup ran consecutively on the same machine. We consider a cloned concurrency control protocol as keeping up if its throughput matches or exceeds the primary's.

Because logging slows Cicada, we compare against Cicada's performance without logging, which is an upper bound on that with logging enabled. We again use the optimal number of primary threads and backup workers, with the latter never exceeding the former. We ran experiments five times and report the median. Error bars show the minimum and maximum.

The workloads are the same as in Sect. 6, and we implement KuaFu in Cicada for fair comparison. In fact, we implement two versions of KuaFu, one optimized for low contention and the other for high contention (The latter matches the published pseudocode [24]). We report whichever achieves better performance.

C5-Cicada prevents potential disasters Cicada's concurrency control is much better than MyRocks's at handling contention. Thus, we expect significantly more potential for replication lag.

Our results support this: C5-Cicada is necessary to prevent replication lag with the standard 50–50% NewOrder-Payment workload after applying similar optimizations as in Sect. 6. Cicada achieves 716,950 transactions per second while KuaFu manages only 596,310 transactions per second, lagging by about 17%. C5-Cicada easily keeps up, committing 1,062,533 transactions per second. The results on the unoptimized workload are similar. The primary's throughput is lower, but KuaFu still lags by about 13,000 transactions per second (3%).

We also explore how the primary and two backups behave under varying levels of contention. Figure 12 compares the

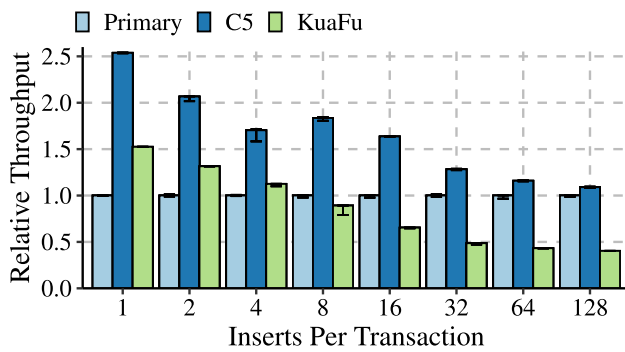


Fig. 13 Primary and backup's throughput on adversarial workload as inserts per transaction increases

throughputs of Cicada, C5-Cicada, and KuaFu on a 50–50% workload as we vary the number of districts from 10 (the standard setting) to 1. As contention increases (i.e., districts decrease), KuaFu lags until 4 districts, but below that, the additional contention harms Cicada's throughput more than KuaFu's by causing significantly higher abort rates (up to about 75%). Thus, with fewer districts, KuaFu keeps up.

To help confirm KuaFu lags due to the constraints imposed on its execution, we re-ran the experiment above but disabled its scheduler's calculation of transaction-granularity constraints. For each number of districts, we compared KuaFu's throughput while using the same number of workers as shown in Fig. 12, and in all cases, KuaFu kept up. For example, with 10 districts and 6 workers, KuaFu's throughput nearly doubles from 596,310 to 1,101,491 transactions per second when its execution is unconstrained. This far exceeds the primary's throughput of 716,950 transactions per second.

These results highlight the complexity of predicting when existing cloned concurrency control protocols will be sufficient to avoid replication lag. As shown in Fig. 12, C5-Cicada always keeps up and thus removes the potential for disasters.

C5-Cicada always keeps up We again validate C5-Cicada using the insert-only and adversarial workloads. On insert-only, Cicada achieves its best performance with transactions of 16 inserts each, amortizing its per-transaction overhead. Here, Cicada inserts about 87M rows/s with 20 threads. KuaFu (96M rows/s with 12 workers) and C5-Cicada (99M rows/s with 10 workers) both keep up. In each case, the scheduler provides sufficient performance.

Figure 13 compares each cloned concurrency control protocol's performance to Cicada's on the adversarial workload. We plot the backup's throughput relative to the primary's median as we vary the number of non-conflicting inserts per transaction.

C5-Cicada mirrors the primary and executes the non-conflicting inserts in parallel, so it always keeps up. This advantage is especially evident as the number of inserts per

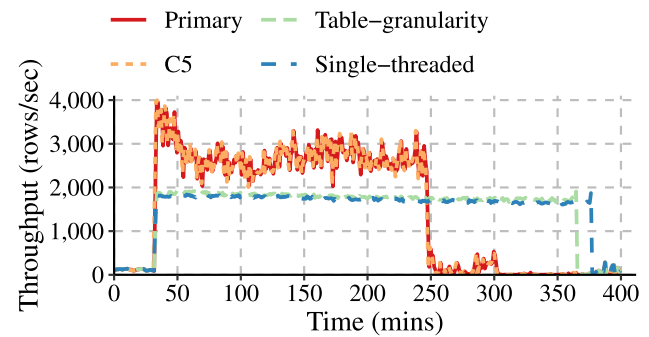


Fig. 14 Lag at Meta used to reach 2h daily and take 2h to recover. With C5, it remains below 3s

transaction increases from 4 to 8. With more parallel work per transaction, C5-Cicada leverages additional workers and its relative throughput actually increases.

On the other hand, the primary's advantage over KuaFu increases with the number of inserts per transaction. With 128 inserts per transaction, KuaFu's throughput is just 40% of the primary's.

8 Deployment experience

MyRocks databases are deployed and used inside the globally distributed data centers at Meta. Each shard of their social graph uses asynchronous primary-backup replication. Further, their internal cloud provides asynchronously replicated, multi-tenant MyRocks instances. Before deploying C5-MyRocks, ensuring short replication lag throughout their infrastructure was a persistent challenge.

A version of C5-MyRocks has been deployed in production since mid-2017, and full deployment finished in early 2019. Deploying C5-MyRocks at Meta has notably improved the reliability of their applications built atop asynchronous primary-backup databases. Since its deployment, anecdotally, the number of complaints about replication lag by on-call engineers drastically decreased.

Figure 14 shows one example of significant lag fixed by C5-MyRocks. It plots throughput over time for one shard. During daily periods of high insert load, the primary's throughput exceeded the backup's, and lag grew to over two hours both with MySQL 5.6's default, single-threaded cloned concurrency control [13, 48] and Meta's earlier table-granularity protocol. After the load spike ends, both protocols took two hours to return lag to zero. C5-MyRocks eradicated the issue, keeping lag below three seconds.

Similarly, after live videos were deployed, popular videos became significant sources of contention. As in the motivating example in Sect. 2.1, all comments on a video caused

writes to a single row. With prior cloned concurrency control protocols, popular videos caused significant replication lag.

Solving replication lag also had secondary benefits: First, deploying C5-MyRocks revealed bottlenecks in downstream systems, such as Wormhole [62], which have since been fixed by Meta's engineers. Second, Meta uses cross-region replication. Multiple copies of the data exist on servers within a primary region, and writes asynchronously replicate to backups in other regions. Short lag reduces the number of times that data must be fetched from other regions to satisfy read-your-writes consistency for clients with recent writes. Further, if the entire primary region fails while all machines in the backup region are lagging, some unreplicated writes may be lost. With C5-MyRocks, the probability of user data loss and the magnitude of such loss if it occurs are both significantly reduced. Finally, C5-MyRocks eliminated a noisy-neighbor problem experienced by applications deployed on Meta's internal cloud. Applications using a multi-tenant MyRocks instance previously would sometimes experience replication lag caused by others sharing the instance.

An unexpected deadlock MySQL has a layered architecture where the replication layer interacts with the storage engine using a well defined API. The replication layer does not have fine grained control over the storage layer. Due to this limitation, Meta engineers discovered a subtle deadlock during deployment. Since C5 is fully capable of scheduling transactions on the storage engines correctly, the storage engine need not hold any row locks. However, some storage engines do not allow disabling row locks completely. For instance, the InnoDB storage engine always holds row locks, including gap locks.

Further, some shards at Meta have a requirement to commit transactions on backups in the same order as they were committed on the primary. This requirement is in addition to monotonic prefix consistency and is used to ensure a backup's log exactly matches its primary's.

We satisfy this requirement by blocking a commit of a transaction until the previous transaction has committed. However, we discovered that sometimes on the storage engines that hold gap locks, this might lead to a deadlock between two independent transactions which are waiting for ordering commits at the replication layer and at the same time contending on a gap lock at the storage engine layer. For example, transactions T_1 and T_2 can be completely independent but end up contending on a gap lock at the storage engine layer.

C5's replication layer has no way of knowing this before scheduling the transactions; it is only aware of the rows that the transactions modify and uses this information to deem them (in)dependent. If T_1 and T_2 are scheduled to run in parallel and T_2 holds a gap lock first, T_1 will end up blocking

behind T_2 . But if commit ordering is enabled, T_2 will not be able to commit before T_1 . Since gap locks are not released until transaction commits, this results in a deadlock.

Oracle MySQL's implementation of parallel replication was also affected by this bug [50] and they came up with a deadlock detection scheme by exposing a callback from the storage engine that tells the upper layers when two transactions conflict on a lock. We borrowed the same solution in C5-MyRocks: We detect these deadlocks using the storage engine callback and rollback and re-execute the transaction that comes later in the commit order.

9 Related work

Deterministic Concurrency Control Deterministic concurrency control protocols [15–17, 60, 65, 71] ensure that database state is a deterministic function of the input log. C5's processing of writes from the primary is inspired by such protocols. These include the up-front resolution of write-write conflicts prior to executing them [15] and its representation of permissible execution schedules of writes [17, 71]. Databases employing deterministic concurrency control, however, do not designate a single replica as a primary and others as backups. They instead employ active replication [64, 65]. C5, on the other hand, is applicable to primary-backup systems where the primary is non-deterministic.

Database Recovery and Replication Database recovery [7, 34, 43, 61, 74] and replication [41, 44, 52, 59, 73] are two problems closely related to cloned concurrency control. In database recovery, changes to the database are logged and stored on stable storage. If a database fails, a recovery protocol creates a new copy of the database. Database replication extends database recovery to reduce recovery time by replicating and applying changes to the backup while the primary executes transactions. If the primary fails, the backup executes a synchronization protocol to bring it into a consistent state before processing new transactions. But database replication (and thus recovery) is simpler than cloned concurrency control because backups do not serve read-only transactions, so the backup only needs to be prefix-consistent before processing new transactions. A cloned concurrency control protocol must always be able to serve read-only transactions from a prefix-consistent state.

Cloned Concurrency Control No existing asynchronous or semi-synchronous cloned concurrency control protocol can guarantee bounded replication lag (Semi-synchronous protocols require a log of a transaction's writes to be persisted at the backup before the transaction commits at the primary). Synchronous protocols [12] trivially guarantee it because the primary and backup coordinate before a transaction commits, but they reduce the primary's performance.

Thus, asynchronous and semi-synchronous are more widely deployed [4, 26, 29, 38, 68].

Transaction- [24, 32, 45, 51] and page-granularity [4, 54, 68] protocols cannot guarantee bounded replication lag. By similar reasoning, coarser granularity protocols, such as those using groups of transactions [39, 46, 55], cannot either.

To the best of our knowledge, Query Fresh [69] is the only existing row-granularity cloned concurrency control protocol, but as discussed below, it does not guarantee bounded lag. It is a semi-synchronous protocol, and to reduce its workers' processing, instantiation of the backup's copy of the database is deferred to read-only transaction threads. These threads load row versions as necessary to return correct results (as defined by monotonic prefix consistency).

On one hand, Query Fresh's lazy instantiation can cause arbitrarily large replication lag by forcing read-only transaction threads to traverse large portions of the log. We provide a detailed description of this case in the appendix of our technical report [23]. On the other, lazy instantiation allows Query Fresh to avoid some of the work that C5 does by applying all writes. An interesting avenue of future work could explore a partially lazy approach.

10 Conclusion

We presented C5, the first cloned concurrency protocol to provide bounded replication lag. C5 comprises three parts: a scheduler, workers, and a snapshotter. C5 is backed by multiple theoretical results showing the necessity of its row-granularity protocol. We also presented two implementations: C5-MyRocks and C5-Cicada. The former is backward-compatible with MyRocks and deployed at Meta, while the latter faithfully implements our design. We demonstrated experimentally they always keep up with their primaries.

Acknowledgements We thank the anonymous reviewers and our shepherd for their helpful comments and feedback. We are also grateful to Princeton's systems group for their comments on earlier versions of this paper. This work was supported by the National Science Foundation under grants CNS-1824130 and IIS-1910613.

References

1. C5-Cicada. <https://github.com/princeton-sns/c5-cicada-exp> (2022)
2. C5-MyRocks. <https://github.com/princeton-sns/c5-myrocks-exp> (2022)
3. Adya, A.: Weak consistency: A generalized theory and optimistic implementations for distributed transactions. Ph.D. thesis, MIT, Cambridge (1999)
4. Antonopoulos, P., Budovski, A., Diaconu, C., Hernandez Saenz, A., Hu, J., Kodavalla, H., Kossmann, D., Lingam, S., Minhas, U.F., Prakash, N., et al.: Socrates: The new SQL server in the cloud. In: Proc. ACM International Conference on Management of Data (SIGMOD), pp. 1743–1756. Association for Computing Machinery, Amsterdam, The Netherlands (2019)
5. Berenson, H., Bernstein, P., Gray, J., Melton, J., O'Neil, E., O'Neil, P.: A critique of ANSI SQL isolation levels. In: Proc. ACM International Conference on Management of Data (SIGMOD), pp. 1–10. Association for Computing Machinery, San Jose (1995)
6. Bernstein, P.A., Hadzilacos, V., Goodman, N.: Concurrency Control and Recovery in Database Systems. Addison-Wesley, Reading (1987)
7. Cha, S.K., Song, C.: P*TIME: Highly scalable OLTP DBMS for managing update-intensive stream workload. In: Proc. International Conference on Very Large Data Bases (VLDB), pp. 1033–1044. VLDB Endowment, Toronto (2004)
8. Chandramouli, B., Prasad, G., Kossmann, D., Levandoski, J., Hunter, J., Barnett, M.: FASTER: A concurrent key-value store with in-place updates. In: Proc. ACM International Conference on Management of Data (SIGMOD), pp. 275–290. Association for Computing Machinery, Houston (2018)
9. DeWitt, D.J., Katz, R.H., Olken, F., Shapiro, L.D., Stonebraker, M.R., Wood, D.A.: Implementation techniques for main memory database systems. In: Proc. ACM International Conference on Management of Data (SIGMOD), pp. 1–8. Association for Computing Machinery, Boston (1984)
10. Diaconu, C., Freedman, C., Ismert, E., Larson, P.A., Mittal, P., Stonecipher, R., Verma, N., Zwilling, M.: Hekaton: SQL server's memory-optimized OLTP engine. In: Proc. ACM International Conference on Management of Data (SIGMOD), pp. 1243–1254. Association for Computing Machinery, New York (2013)
11. Duplyakin, D., Ricci, R., Maricq, A., Wong, G., Duerig, J., Eide, E., Stoller, L., Hibler, M., Johnson, D., Webb, K., Akella, A., Wang, K., Ricart, G., Landweber, L., Elliott, C., Zink, M., Cecchet, E., Kar, S., Mishra, P.: The design and operation of CloudLab. In: Proc. USENIX Annual Technical Conference (ATC), pp. 1–14. USENIX Association, Renton (2019)
12. Elnikety, S., Pedone, F., Zwaenepoel, W.: Generalized snapshot isolation and a prefix-consistent implementation. Tech. Rep. IC/2004/21, School of Computer and Communication Sciences, EPFL, Lausanne (2004)
13. Facebook: MyRocks GitHub Wiki (2019). <https://github.com/facebook/mysql-5.6/wiki>
14. Facebook: RocksDB: A Persistent Key-value Store (2020). <https://rocksdb.org/>
15. Faleiro, J.M., Abadi, D.J.: Rethinking serializable multiversion concurrency control. Proc. Very Large Data Bases Endow. (PVLDB) **8**(11), 1190–1201 (2015)
16. Faleiro, J.M., Abadi, D.J., Hellerstein, J.M.: High performance transactions via early write visibility. Proc. Very Large Data Bases Endow. (PVLDB) **10**(5), 613–624 (2017)
17. Faleiro, J.M., Thomson, A., Abadi, D.J.: Lazy evaluation of transactions in database systems. In: Proc. ACM International Conference on Management of Data (SIGMOD), pp. 15–26. Association for Computing Machinery, Snowbird (2014)
18. Gawlick, D., Kinkade, D.: Varieties of concurrency control in IMS/VS fast path. IEEE Data Eng. Bull. **8**(2), 3–10 (1985)
19. GitLab: GitLab.com Database Incident (2017). <https://about.gitlab.com/2017/02/01/gitlab-dot-com-database-incident/>
20. GitLab: Postmortem of Database Outage of January 31 (2017). <https://about.gitlab.com/2017/02/10/postmortem-of-database-outage-of-january-31/>
21. Google: LevelDB (2020). <https://github.com/google/leveldb>
22. Hellerstein, J.M., Stonebraker, M., Hamilton, J.: Architecture of a database system. Found. Trends Databases **1**(2), 141–259 (2007)
23. Helt, J., Sharma, A., Abadi, D.J., Lloyd, W., Faleiro, J.M.: C5: Cloned concurrency control that always keeps up (2022).

- <https://doi.org/10.48550/arXiv.2207.02746>. <https://arxiv.org/abs/2207.02746>
24. Hong, C., Zhou, D., Yang, M., Kuo, C., Zhang, L., Zhou, L.: KuaFu: Closing the parallelism gap in database replication. In: Proc. IEEE International Conference on Data Engineering (ICDE), pp. 1186–1195. Institute of Electrical and Electronics Engineers, Brisbane (2013)
 25. Huang, Y., Qian, W., Kohler, E., Liskov, B., Shriram, L.: Opportunities for optimism in contended main-memory multicore transactions. Proc. Very Large Data Bases Endow. (PVLDB) **13**(5), 629–642 (2020)
 26. Instagram: Instagrator Part 2: Scaling Our Infrastructure To Multiple Data Centers (2015). <https://instagram-engineering.com/instagrator-pt-2-scaling-our-infrastructure-to-multiple-data-centers-5745cbad7834>
 27. Intel: Intel Optane Persistent Memory (2020). <https://www.intel.com/content/www/us/en/architecture-and-technology/optane-dc-persistent-memory.html>
 28. Intel: Intel Server Products (2020). <https://www.intel.com/content/www/us/en/products/servers.html>
 29. Jeffrey, M.C., Ying, V.A., Subramanian, S., Lee, H.R., Emer, J., Sanchez, D.: Harmonizing speculative and non-speculative execution in architectures for ordered parallelism. In: Proc. IEEE/ACM International Symposium on Microarchitecture (MICRO), pp. 217–230. Institute of Electrical and Electronics Engineers, Fukuoka City (2018)
 30. Johnson, R., Pandis, I., Stoica, R., Athanassoulis, M., Ailamaki, A.: Aether: a scalable approach to logging. Proc. Very Large Data Bases Endow. (PVLDB) **3**(1–2), 681–692 (2010)
 31. Kim, K., Wang, T., Johnson, R., Pandis, I.: ERMIA: Fast memory-optimized database system for heterogeneous workloads. In: Proc. ACM International Conference on Management of Data (SIGMOD), pp. 1675–1687. Association for Computing Machinery, San Francisco (2016)
 32. King, R.P., Halim, N., Garcia-Molina, H., Polyzois, C.A.: Management of a remote backup copy for disaster recovery. ACM Trans. Database Syst. (TODS) **16**(2), 338–368 (1991)
 33. Larson, P.R., Blanas, S., Diaconu, C., Freedman, C., Patel, J.M., Zwillig, M.: High-performance concurrency control mechanisms for main-memory databases. VLDB Endow. **5**(4), 298–309 (2011)
 34. Lee, J., Kim, K., Cha, S.K.: Differential logging: A commutative and associative logging scheme for highly parallel main memory databases. In: Proc. IEEE International Conference on Data Engineering (ICDE), pp. 173–182. Institute of Electrical and Electronics Engineers, Heidelberg (2001)
 35. Levandoski, J., Lomet, D., Sengupta, S., Stutsman, R., Wang, R.: High performance transactions in Deuteronomy. In: Proc. Conference on Innovative Data Systems Research (CIDR), pp. 1–12. cidrdb.org, Asilomar (2015)
 36. Lim, H., Kaminsky, M., Andersen, D.G.: Cicada: Dependably fast multi-core in-memory transactions. In: Proc. ACM International Conference on Management of Data (SIGMOD), pp. 21–35. Association for Computing Machinery, Chicago (2017)
 37. Linux: libhugetlbfs: preload library to back text, data, malloc() or shared memory with hugepages (2020). <https://linux.die.net/man/7/libhugetlbfs>
 38. Lu, H., Veeraraghavan, K., Ajoux, P., Hunt, J., Song, Y.J., Tobias, W., Kumar, S., Lloyd, W.: Existential consistency: Measuring and understanding consistency at facebook. In: Proc. ACM Symposium on Operating Systems Principles (SOSP), pp. 295–310. Association for Computing Machinery, Monterey (2015)
 39. MariaDB: MariaDB 10 Parallel Replication (2017). <https://mariadb.com/kb/en/parallel-replication/>
 40. Matsunobu, Y.: Making slave pre-fetching work better with SSD. <https://yoshinorimatsunobu.blogspot.com/2011/10/making-slave-pre-fetching-work-better.html/> (2011)
 41. Minhas, U.F., Rajagopalan, S., Cully, B., Abounaga, A., Salem, K., Warfield, A.: RemusDB: transparent high availability for database systems. VLDB J. **22**(1), 29–45 (2013)
 42. Mitzukas, D.: On MySQL replication prefetching. <https://dom.as/2011/12/03/replication-prefetching/> (2011)
 43. Mohan, C., Haderle, D., Lindsay, B., Pirahesh, H., Schwarz, P.: ARIES: a transaction recovery method supporting fine-granularity locking and partial rollbacks using write-ahead logging. ACM Trans. Database Syst. (TODS) **17**(1), 94–162 (1992)
 44. Mohan, C., Treiber, K., Obermarck, R.: Algorithms for the management of remote backup data bases for disaster recovery. In: Proc. IEEE International Conference on Data Engineering (ICDE), pp. 511–518. Institute of Electrical and Electronics Engineers, Vienna (1993)
 45. MySQL: Improving The Parallel Applier With Writese-based Dependency Tracking (2017). <https://mysqlhighavailability.com/improving-the-parallel-applier-with-writese-based-dependency-tracking/>
 46. MySQL: Group Commit of Binary Log (2019). <https://dev.mysql.com/worklog/task/?id=5223>
 47. MySQL: InnoDB (2019). <https://dev.mysql.com/doc/refman/5.6/en/innoDB-storage-engine.html>
 48. MySQL: MySQL 5.6 Reference Manual (2020). <https://dev.mysql.com/doc/refman/5.6/en/>
 49. Narula, N., Cutler, C., Kohler, E., Morris, R.: Phase reconciliation for contended in-memory transactions. In: Proc. USENIX Symposium on Operating Systems Design and Implementation (OSDI), pp. 511–524. USENIX Association, Broomfield (2014)
 50. Nielsen, K.: Bug #74177 –slave-preserve-commit-order causes slave to deadlock and break for some queries. <https://bugs.mysql.com/bug.php?id=74177> (2014)
 51. Oracle: Method of applying changes to a standby database system (2001). Patent No. US6980988B1, Filed Oct. 1, 2001, Issued Dec. 27, 2005
 52. Oracle: Eager replication of uncommitted transactions (2014). Patent No. US9747356B2, Filed Jan. 23, 2014, Issued Aug. 29, 2017
 53. Oracle: Oracle 19 Database Administrator’s Guide (2020). <https://docs.oracle.com/en/database/oracle/oracle-database/19/admin/index.html>
 54. Oracle: Oracle Active Data Guard (2020). <https://www.oracle.com/technetwork/database/availability/dg-adg-technical-overview-wp-5347548.pdf>
 55. Oracle: Understanding Oracle GoldenGate (2020). <https://docs.oracle.com/en/middleware/goldengate/core/19.1/>
 56. Papadimitriou, C.H.: The serializability of concurrent database updates. J. ACM (JACM) **26**(4), 631–653 (1979)
 57. PostgreSQL: PostgreSQL 12.1 Documentation (2020). <https://www.postgresql.org/docs/12/index.html>
 58. PostgreSQL: Snapshot Synchronization Functions (2020). <https://www.postgresql.org/docs/current/functions-admin.html>
 59. Qin, D., Brown, A.D., Goel, A.: Scalable replay-based replication for fast databases. Proc. Very Large Data Bases Endow. (PVLDB) **10**(13), 2025–2036 (2017)
 60. Qin, D., Brown, A.D., Goel, A.: Caracal: Contention management with deterministic concurrency control. In: Proc. ACM Symposium on Operating Systems Principles (SOSP), pp. 180–194. Association for Computing Machinery, Virtual Event (2021)
 61. Schwalb, D., Faust, M., Wust, J., Grund, M., Plattner, H.: Efficient transaction processing for Hyrise in mixed workload environments. In: Proc. International Workshop on In Memory Data Management and Analytics (IMDM), pp. 16–29. Springer, Hangzhou (2014)
 62. Sharma, Y., Ajoux, P., Ang, P., Callies, D., Choudhary, A., Demailly, L., Fersch, T., Guz, L.A., Kotulski, A., Kulkarni, S., Kumar, S., Li, H., Li, J., Makeev, E., Prakasam, K., Renesse, R.V., Roy, S., Seth, P., Song, Y.J., Wester, B., Veeraraghavan, K., Xie,

- P.: Wormhole: Reliable pub-sub to support geo-replicated internet services. In: Proc. USENIX Symposium on Networked Systems Design and Implementation (NSDI), pp. 351–366. USENIX Association, Oakland (2015)
63. Terry, D.B., Demers, A.J., Petersen, K., Spreitzer, M.J., Theimer, M.M., Welch, B.B.: Session guarantees for weakly consistent replicated data. In: Proc. International Conference on Parallel and Distributed Information Systems (PDIS), pp. 140–149. Institute of Electrical and Electronics Engineers, Austin (1994)
64. Thomson, A., Abadi, D.J.: The case for determinism in database systems. Proc. Very Large Data Bases Endow. (PVLDB) **3**(1–2), 70–80 (2010)
65. Thomson, A., Diamond, T., Weng, S.C., Ren, K., Shao, P., Abadi, D.J.: Calvin: Fast distributed transactions for partitioned database systems. In: Proc. ACM International Conference on Management of Data (SIGMOD), pp. 1–12. Association for Computing Machinery, Scottsdale (2012)
66. Council, T.P.C.: TPC Benchmark C revision 5.11 (2010)
67. Tu, S., Zheng, W., Kohler, E., Liskov, B., Madden, S.: Speedy transactions in multicore in-memory databases. In: Proc. ACM Symposium on Operating Systems Principles (SOSP), pp. 18–32. Association for Computing Machinery, Farmington (2013)
68. Verbitski, A., Gupta, A., Saha, D., Corey, J., Gupta, K., Brahmadesam, M., Mittal, R., Krishnamurthy, S., Maurice, S., Kharatishvili, T., Bao, X.: Amazon Aurora: On avoiding distributed consensus for I/Os, commits, and membership changes. In: Proc. ACM International Conference on Management of Data (SIGMOD), pp. 789–796. Association for Computing Machinery, Houston (2018)
69. Wang, T., Johnson, R., Pandis, I.: Query fresh: Log shipping on steroids. Proc. Very Large Data Bases Endow. (PVLDB) **11**(4), 406–419 (2017)
70. Wang, T., Kimura, H.: Mostly-optimistic concurrency control for highly contended dynamic workloads on a thousand cores. Proc. Very Large Data Bases Endow. (PVLDB) **10**(2), 49–60 (2016)
71. Whitney, A., Shasha, D., Apter, S.: High volume transaction processing without concurrency control, two phase commit, SQL or C++. In: International Workshop on High Performance Transaction Systems, pp. 211–217. Springer, Asimolar (1997)
72. Yan, C., Cheung, A.: Leveraging lock contention to improve OLTP application performance. Proc. Very Large Data Bases Endow. (PVLDB) **9**(5), 444–455 (2016)
73. Zamanian, E., Yu, X., Stonebraker, M., Kraska, T.: Rethinking database high availability with RDMA networks. Proc. Very Large Data Bases Endow. (PVLDB) **12**(11), 1637–1650 (2019)
74. Zheng, W., Tu, S., Kohler, E., Liskov, B.: Fast databases with fast durability and recovery through multicore parallelism. In: Proc. USENIX Symposium on Operating Systems Design and Implementation (OSDI), pp. 465–477. USENIX Association, Broomfield (2014)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.