# Weakly Durable High-Performance Transactions

Yun-Sheng Chang[*]          Yu-Fang Chen       Hsiang-Shang Ko
*MIT CSAIL*       *Institute of Information Science, Academia Sinica, Taiwan*

## Abstract

Existing disk-based database systems largely fall into two categories—they either provide very high performance but few guarantees, or expose the transaction abstraction satisfying the full ACID guarantees at the cost of lower performance. In this paper, we present an alternative that achieves the best of both worlds, namely good performance and transactional properties. Our key observation is that, because of the frequent use of synchronization primitives, systems with strong durability can hardly utilize the extremely high parallelism granted by modern storage devices. Thus, we explore the notion of weakly durable transactions, and discuss how to safely relax durability without compromising other transactional properties. We present AciKV, a transactional system whose design is centered around weak durability. AciKV exposes to users the normal transactional interface, but what sets it apart from others is a new "persist" primitive that decouples durability from commit. AciKV is a middle ground between systems that perform fast atomic operations, and ones that support transactions; this middle ground is useful as it provides similar performance to the former, while prevents isolation and consistency anomalies like the latter. Our evaluation using the YCSB benchmark shows that AciKV, under workloads that involve write requests, exhibits more than two orders of magnitude higher throughput than existing strongly durable systems.

## 1   Introduction

Existing disk-based database systems often sit on the two ends of the performance-guarantee spectrum. At one end of the spectrum, we have systems that can perform operations in an extremely efficient manner, but often come with very few guarantees, for instance, the atomicity of a single operation. At the other end of the spectrum, we have transactional database systems that guarantee the appealing atomicity, consistency, isolation, and durability (ACID) properties, but often have significantly lower performance compared with ones that have no support for transactions. This state of affairs forces users to pick between good performance and useful abstraction.

We argue that this unsatisfactory situation can be primarily attributed to the pursuit of *strong durability*, which induces high synchronization overhead. To validate this argument, we ran an experiment on some popular trans-
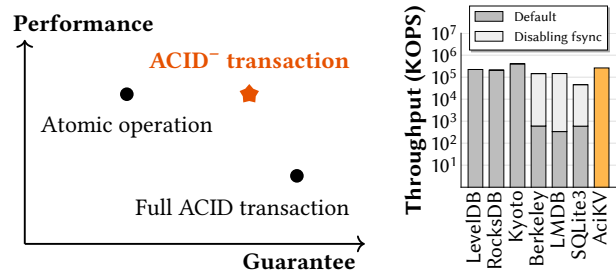


**Figure 1**: Evaluation results of three non-transactional systems (LevelDB, RocksDB, Kyoto Cabinet), three transactional systems (Berkeley DB, LMDB, SQLite3), and this work (AciKV) under a 50R-50W workload. See §4.1 for the detailed experimental setup.

actional systems in two modes: a default strongly durable mode, and a weakly durable mode that disables the issuing of fsync[1] to the underlying file system. As shown in Figure 1, all systems in the weakly durable mode performs at least two orders of magnitude faster than when they are in the default strongly durable mode. A similar experiment conducted recently at the level of file systems and of disks [58] also aligns with our result. This analysis motivates us to explore the notion of *weakly durable transactions*, which aims to perform comparably to systems capable of carrying out fast operations, while providing the useful transaction abstraction.

Another motivation for this work is hardware trend. Modern storage devices, such as flash-based solid-state drives (SSDs) and emerging non-volatile memory (NVM) technology, are often structured in a highly parallel way. While they promise to provide tremendous throughput, strongly durable systems can hardly utilize the granted parallelism because of the frequent use of synchronization primitives.

A common technique that can also alleviate the synchronization overhead is *group commit*, which groups multiple transactions into one batch, and delays issuing a synchronization primitive (e.g., an fsync) until the end of each batch. The problem with group commit, however, is an inherent trade-off between throughput and latency—a large batch increases throughput but suffers from higher latency, and a small batch has a lower latency but also a lower throughput. In §4.2, we experimentally evaluate the two techniques, weak durability and group commit, and discuss their differences in detail.

---

[*]Work done while the author was affiliated with Academia Sinica.

[1]Throughout this paper, we use fsync to denote file system synchronization primitives including fdatasync and sync.

At this point, one might ask: Isn't disabling `fsync` good enough? The answer is no, as simply disabling `fsync` can easily lead to *internal inconsistency*. This means that the consistency invariants (e.g., an allocated data block should always be referenced by some metadata) enforced internally by the system may no longer hold true after a full system crash (e.g., a power failure). A mild consequence is resource leaks, and a more severe one is unusable corrupted databases [60]. The reason internal inconsistency may arise is that file system operations may be reordered on a crash [4, 41], and without `fsync`, the database system cannot enforce certain write orderings for crash safety. Thus, the LMDB developer faithfully warns its users about the risk of enabling such option [15]:

> "MDB_NOSYNC ... This optimization means a system crash can corrupt the database or lose the last transactions if buffers are not yet flushed to disk. ... However, if the filesystem preserves write order and the MDB_WRITEMAP flag is not used, transactions exhibit ACI (atomicity, consistency, isolation) properties and only lose D (durability)."

SQLite3 also gives their users a similar warning [53]. But there is a more convoluted issue. In the context of database transactions, consistency (i.e., C in ACID) is often specified with a set of *integrity constraints*. For instance, one might require the account balance to be non-negative, or the user identifier to be unique. Without any restriction on weak durability, it is possible for a crash to create *external inconsistency*; that is, a database user may observe a violation of the integrity constraints after a crash. In §2.1, we give an example of external inconsistency induced by weak durability, and elaborate more on its cause.

Naturally, our exploration of weakly durable transactions starts with a theoretical analysis on how to *safely* relax durability to avoid inconsistency; §2.2 shows that a mild adaptation of the standard serializability theory [3] is sufficient for this purpose. The refined weak durability, along with other transactional properties, gives rise to what we call the ACID⁻ properties. Table 1 shows an interface of weakly durable transactional systems. The crux of weak durability

| Category | Primitive |
|---|---|
| Data | *get(k), getrange(k1, k2), put(k, v), delete(k)* |
| Transaction | *begin, commit, abort* |
| Durability | *persist* |

**Table 1**: An interface of weakly durable transactional systems.

is manifested in the new *persist* primitive, which decouples durability from *commit*: committed transactions are not guaranteed to survive across crashes; only when they encounter a *persist* do they become persistent. However, users can rest assured that no inconsistency will arise because of crashes.

As a demonstration, we built AciKV, which implements the interface in Table 1 and guarantees ACID⁻. While the main benefit of weak durability is a significant reduction in synchronization overhead, we believe the *permissiveness* of weak durability also implies a simpler and more efficient system design. Thus, we built AciKV from scratch using weak durability as its "first-class" property (rather than building it to ensure strong durability, and trying to relax durability later on). In §3.1, we introduce the recovery mechanism of AciKV, including an efficient client-server synchronization protocol and a shadow paging technique [33], which together allow AciKV to easily create a consistent snapshot in a crash-safe manner. Compared with the more commonly used ARIES-style [36] write-ahead logging (WAL) that often has to depend on other layers such as buffers, or to assume a certain locking protocol, AciKV's recovery mechanism is much simpler, in that it *fully decouples crash safety from other layers*. We also discuss why the conventional wisdom that shadow paging is not suitable for concurrent writes [24] may not be the case for weak durability. In §3.2, we describe AciKV's latch-free two-level index design that fits well with the *batch-processing* nature of weak durability. The index consists of an in-memory concurrent skip list and an on-disk B+-tree. The skip list is employed for absorbing insertions, so that the B+-tree structure is guaranteed to remain the same until the next *persist*, at which point the two data structures are merged. This means that no index latches are required. Moreover, unlike other multi-level index structures such as the LSM-tree [39], AciKV does not need an extra mechanism to protect data in the skip list, as the data are changes made after the last *persist*.

Using the YCSB benchmark suite [17], we compare AciKV with six popular disk-based database systems, three transactional and three non-transactional. Evaluation results show that, under workloads that involve writes, the throughput of AciKV can be two or even three orders of magnitude higher than that of the transactional systems, and is on par with that of the non-transactional ones. Further, to understand the benefit of building AciKV around weak durability, we also compare AciKV with the transactional systems running in their weakly durable mode. Results show that AciKV is the only system that consistently scales across a wide range of workloads.

**Limitations of AciKV.**

- AciKV is a *storage engine*, whose main responsibilities are concurrency control and crash recovery. Thus, its interface supports only the generic byte-array data type and a record-at-a-time data access method, and has no mechanisms for specifying integrity constraints. System designers can build their own data model (e.g., SQL) on top of AciKV.

- AciKV lacks features commonly seen in database systems, such as compression and replication. We believe that

AciKV's design and properties are compatible with these features, so adding them to AciKV should be feasible.

**Contributions of this work.**

- We introduce the notion of weakly durable transactions, which aims to provide the useful transaction abstraction with enough guarantees, namely ACID$^-$, while maintaining good performance.

- We present AciKV, a transactional key-value store whose design is centered around weakly durable transactions.

- We evaluate AciKV by comparing it with six popular disk-based database systems using the YCSB benchmark.

## 2 From ACID to ACID$^-$

Weakly durable transactional systems guarantee ACID$^-$, whose semantics follows that of standard ACID except for Durability, which is changed to guarantee that only those committed transactions made persistent by the *persist* primitive (shown in Table 1) will survive across crashes. In this section, we first review "ad-hoc" weak durability seen in existing transactional systems (§2.1); in particular, we look at a simple example of inconsistency that may arise, showing that this kind of weak durability is potentially unsafe. Then, we present our theoretical analysis on how to safely relax durability to maintain consistency (§2.2).

### 2.1 Ad-hoc weak durability

ACID has been the gold standard of transactional systems for decades. We briefly recap its definition here. Atomicity says that the effect of each transaction should appear as a whole to other transactions. Consistency ensures that each transaction always observes a consistent state satisfying all the given integrity constraints; also, on a crash, the recovered state should be consistent. Isolation creates the illusion that each transaction is the sole owner of the database, even though multiple transactions may be executed concurrently. Durability guarantees that the effect of committed transactions will always survive across crashes. By promising these properties, transactions can often greatly relieve the programming burden of users.

But supporting the strongest form of ACID is expensive. Many transactional systems therefore provide options to relax some of the guarantees, so that users can "pay" only for the guarantees they actually need. For instance, many systems can run transactions at a lower isolation level. In this work, however, we will discuss only systems with atomicity and serializability, the highest isolation level (and focus on the interplay between consistency and durability). Some transactional systems also have options to relax durability, often achieved by not issuing slow synchronization primitives such as fsync. As shown in Figure 1, these options have the potential to significantly improve the performance. To distinguish this kind of weak durability from the one we

propose in this work, we refer to the former as *ad-hoc weak durability*.

The downside of ad-hoc weak durability, however, is the risk of internal and external inconsistency when encountering a crash. As we have covered internal inconsistency in §1, here we focus on the external one, that is, violations of integrity constraints. The reason ad-hoc weak durability is subject to inconsistency is that it does not specify *which transactions are allowed to be discarded on a crash.* Below we illustrate the issue with an example. Suppose we have a database consisting of two integer objects $x$ and $y$, and an integrity constraint $x < y$. Consider the following serial history:

$$r_1(x, 0)r_1(y, 1)w_1(y, 2)c_1 r_2(x, 0)r_2(y, 2)w_2(x, 1)c_2$$

In the history, transaction T1 first observes $(x, y) = (0, 1)$, updates $y$ to 2 without violating the constraint $x < y$, and commits; transaction T2 then observes $(x, y) = (0, 2)$, so it is also safe for T2 to update $x$ to 1. Assume a crash occurs after T2 commits. While both transactions preserve the constraint, inconsistency can still arise in the post-crash state when we naively allow the weakly durable system to discard *any* transaction on a crash—if the system chose to discard the effect of T1, the resulting state would be $(x, y) = (1, 1)$, violating $x < y$.

The previous example suggests that we need a refined version of weak durability to safely relax durability. Before diving into our theoretical analysis, we begin with some intuition. Note that there are four possible scenarios after the crash: (i) only T1, (ii) only T2, (iii) both T1 and T2, (iv) none of them survives across the crash, but only the second case can give rise to inconsistency. The root cause of the problem is that *T2 depends on T1's effect*, or more concretely, T2 observes T1's update of $y$ from 1 to 2. Naturally, we should refine our definition of weak durability to incorporate the *dependencies* between transactions.

Note that we by no means claim that ad-hoc weak durability *will* create inconsistency, and indeed, a particular system may be free from inconsistency due to certain implementation artifacts. The point being made is that we should justify *why* a weakly durable system is not subject to inconsistency.

### 2.2 Consistency of weakly durable systems

To make sure theoretically that consistency will not be lost if we adopt weak durability, we analyze the standard crash consistency argument for database systems with (strong) durability, serializability, and consistency preservation by all transactions [3: Section 6.2]. We will see that, with some mild adaptation, essentially the same argument also works for weakly durable systems.

The high-level argument is fairly straightforward. Let modifications to a database be represented as a history $H$ of read, write, commit, or abort operations performed by transactions. When the database system crashes, (strong) durability guarantees that the recovered database retains ex-

actly the effects of those operations performed by *committed* transactions in $H$. This retained part of $H$ is called the *committed projection* of $H$ and denoted by $C(H)$. Serializability then guarantees that the recovered database is the same as the result of executing all the transactions in $C(H)$ in some serial order. Finally, since each of the transactions preserves consistency, the recovered database must be consistent.

The durability guarantee is more involved than it seems, as it may not even be theoretically possible to retain exactly the effects of the operations in $C(H)$ due to the *dependencies* among the operations in $H$ (which form a partial order, so in general $H$ may be an acyclic directed graph rather than a sequence). Here we do not need to discuss specific forms of dependency (e.g., read-after-write dependency), and only need to assume that if an operation $op$ depends on another operation $op'$, then the effect of $op'$ must happen before $op$ can take effect. (For example, $op$ and $op'$ may be "reading an object $x$" and "writing to $x$", and the effect of $op$, namely getting a particular value $v$ of $x$, cannot happen if $op'$ does not take effect because $v$ has to be written by $op'$ first.) Dependencies may prevent the database system from retaining exactly the effects of the operations in $C(H)$, because an operation in $C(H)$ may depend on some other operation not in $C(H)$ and could not take effect. Therefore, to make durability at least theoretically possible (before we can move on to think about how to implement durability), we should ensure that $C(H)$ is a *prefix* (or a "downward closed" subset) of $H$, meaning that for every operation in $C(H)$, all the operations it depends on are in $C(H)$ as well. This can be ensured by committing transactions with some care, and one way is to enforce the following *prefix preservation* property (which is comparable to "recoverability" [3: Section 2.4]): whenever an operation of a transaction $T$ depends on an operation of another transaction $T'$, if $T$ is committed, then $T'$ must be committed before $T$ is. This implies that $C(H)$ is a prefix of $H$: every operation $op$ in $C(H)$ belongs to a committed transaction, so for any operation $op'$ in $H$ on which $op$ depends, the transaction containing $op'$ must be committed, and hence $op'$ is in $C(H)$ as well.

For a database system with weak (instead of strong) durability, we can reuse essentially the same crash consistency argument, except that $H$ can now include *persist* operations, and what the system retains after a crash is the effect of the *persistently committed projection* $PC(H)$, which consists of the operations of the transactions committed before any *persist* operation in $H$. Like in the case of strong durability, we should ensure that $PC(H)$ is a prefix of $H$ so that weak durability is at least theoretically possible. For this, the same prefix preservation property is sufficient: given an operation $op$ of a persistently committed transaction $T$, any operation $op'$ on which $op$ depends must belong to a transaction $T'$ which, by prefix preservation, is committed before $T$ is, and hence also persistently committed (due to the same *persist* operation that makes $T$ persistently committed).
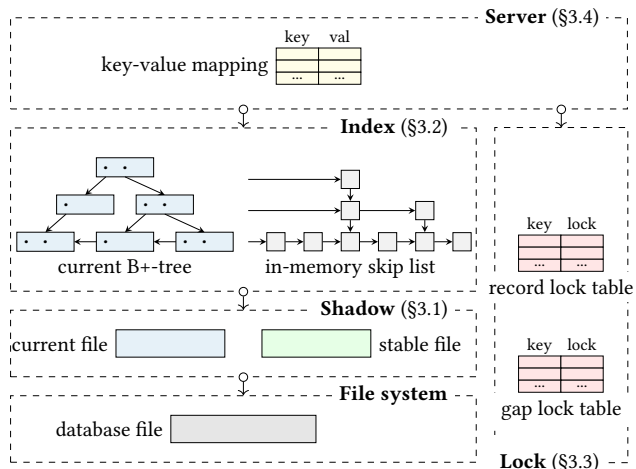


**Figure 2**: Overview of the AciKV design. Each layer implements its own state on top of its underlying layers. ⇕: use.

The above analysis shows that, practically, it is possible to reuse some existing scheduling mechanism with prefix preservation and serializability to implement weakly durable systems that guarantee crash consistency. Indeed this is what we do for our system AciKV.

## 3 AciKV Design

In this section, we present the design of AciKV. Figure 2 shows its main layers. Starting from a database file, AciKV builds on top of it a *shadow paging* mechanism (§3.1) to ensure crash safety. More concretely, the shadow guarantees to revert the database to a consistent state after a system crash. Similarly to all database systems, AciKV uses an indexing mechanism (§3.2) to speed up point and range queries. Its index has two parts: an on-disk B+-tree and an in-memory skip list. On receiving a *persist*, AciKV will merge the skip list into the B+-tree and persistently create a consistent snapshot with the shadow for durability. To ensure serializability of database transactions, AciKV adopts the strong strict two-phase locking (SS2PL) protocol (§3.3) on records and key intervals. On top of the indexing mechanism and SS2PL, AciKV provides standard operations (in addition to *persist*) of key-value stores for its users (§3.4). For simplicity, we omit layers (e.g., the buffer layer) that are standard and whose absence does not obscure the exposition of AciKV's behavior.

### 3.1 Shadow paging

Dealing with crashes is hard, as one has to thoroughly consider every possible post-crash states and carefully order certain writes to prevent crash-safety bugs. To ensure crash safety, AciKV relies on the *shadow paging* (or simply shadow) technique [9, 24, 27, 33, 47, 48]. We choose to implement a modern and efficient one presented in [9]. Below we briefly describe its design; for more details (in particular its correctness), we refer readers to the original paper.

4

At the core of the shadow is an in-memory page table, mapping logical addresses (i.e., the ones upper layers observe) to physical addresses (i.e., the ones being used to index the underlying database file). Modifications are made in a out-of-place manner: On receiving a *write* request, the shadow assigns a free physical page to hold the data, and updates the corresponding entry in the page table to point to the assigned page. Notice that due to the out-of-place update, the old data remain intact on the database file; this is crucial as the recovery procedure may need the old data to revert the database to a consistent state.

On receiving a *flush* request, the page table will be stored on the database file for durability, in the form of table differences (i.e., the deltas), and occasionally when running out of space to accommodate the deltas, in the form of a full table image. We call the full table image plus the deltas the *stable table*. To ensure write ordering and durability, `fsync`s are issued as needed. The garbage collector is also designed in a way that will not remove data pointed to by the stable table.
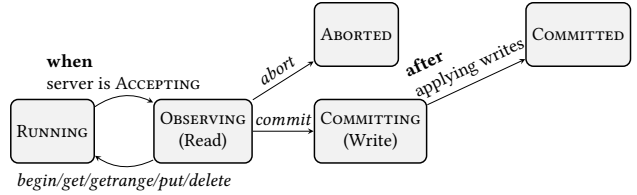
On recovery, the full table image is used as the base image, and the deltas are applied sequentially to rebuild the stable table. In addition, since all data pointed to by the stable table are still there, the procedure can safely recover the file state to the time when the last *flush* command succeeds.

The shadow has a simple specification, consisting of two files *current* and *stable*. The upper layer can *read* and *write* only the current one, and crash-atomically store a *snapshot* of the current one into the stable one with the *flush* command. On crashes, there is no warranty about the contents of the current file, but the stable file is guaranteed to remain intact so that during *crash recovery*, the shadow can safely bring the snapshot back from the stable to the current one. On top of the shadow, system designers can implement arbitrary data structures; for instance, AciKV implements a B+-tree.
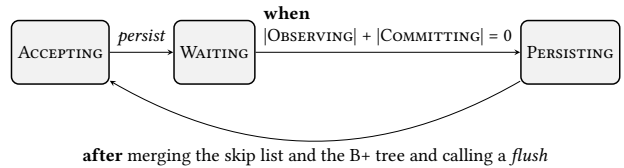
With the shadow at our disposal, ensuring crash safety is effectively reduced to *creating a consistent snapshot*. This means that we can update the stable storage as if crashes are absent, which frees us from reasoning about the intricate crash behavior and thus greatly simplify the design of upper layers.

**Creating consistent snapshot.** The semantics of *persist* gives "consistent snapshot" a precise definition: It is the state where committed transactions have applied all their writes, in the order governed by SS2PL (§3.3) to ensure serializability, and uncommitted transactions have not applied any of its writes. Next we describe how AciKV creates a consistent snapshot efficiently.

As each transaction can run for an arbitrarily long time, it would be unreasonable to wait for all active transactions to commit or abort before creating a snapshot. Thus, AciKV uses an efficient protocol to orchestrate the transactions and the database server, as depicted in the form of state transition diagrams in Figure 3.



(a) Client transition diagram.



**after** merging the skip list and the B+ tree and calling a *flush*

(b) Server transition diagram.

**Figure 3**: Synchronizing clients and the database server.

We first explain the transition diagram of a client, as shown in Figure 3(a). Initially, all clients are in the RUNNING state. A client can perform database operations only in the OBSERVING state. So in order to perform an operation, the client has to first *enter* the server by transitioning its state from RUNNING to OBSERVING. This is possible only when the server is in the ACCEPTING state. An *abort* operation changes the client state from OBSERVING to ABORTED. A *commit* operation changes the client state from OBSERVING to COMMITTING. To prevent uncommitted transactions from updating the server, each planned modification to the server has to first go to the per-transaction local write set, and is applied to the server during the client's COMMITTING phase. In the COMMITTING state, the client updates the server with its local write set and move to COMMITTED only when all the entries in the local write set are processed. We say that a client *leaves* the server when it moves to RUNNING, ABORTED, or COMMITTED.

Next we explain the transition diagram of the server, as shown in Figure 3(b). Initially, the server is in the ACCEPTING state. On receiving a *persist*, the server first moves to the WAITING state and blocks all clients from entering. When the server is in the WAITING state, some clients might still in OBSERVING or COMMITTING states, that is, still in the server. The server will wait until all clients have left and then move to the PERSISTING state to create a snapshot.

The protocol ensures the property that when the server is PERSISTING, no client is OBSERVING or COMMITTING. Notice that at the moment when the server just moves to the PERSISTING state, all clients are either RUNNING, ABORTED, or COMMITTED. The latter two states have no chance to transition to the OBSERVING and COMMITTING states. When the server is in the PERSISTING state, the client cannot move from RUNNING to either OBSERVING or COMMITTING. Hence the property is guaranteed. This property further ensures that only committed transactions apply their writes to the

server. For uncommitted ones, their writes only stay in their local write set.

One important assumption of the protocol is that checking transition guard conditions (e.g., checking if the state is ACCEPTING) and transitioning to the next state (e.g., to the OBSERVING state) are done *atomically*. Without this assumption, we can easily create an example to break the protocol's property. Assume we have only one client, and we always have one server. Initially, the client is in the RUNNING state and the server is in the ACCEPTING state. The client wants to enter the server so it checks the server's state, which is now ACCEPTING. Before the client moves to OBSERVING, the server has moved to the WAITING state, checks the guard condition "no client is OBSERVING or COMMITTING", and immediately moves to PERSISTING. The client moves to OBSERVING only after the server moves to PERSISTING. Now the property "when the server is PERSISTING, no client is in OBSERVING or COMMITTING" is broken.

Thus, we develop a mechanism to ensure checking transition guards and transitioning to the next states are done atomically. The implementation of the synchronization protocol is shown in Figure 4.

```
atomic_uint n_accessing;          int server_enter(void)
bool accepting;                   {
mutex_t mutex;                        n_accessing++;
                                      mfence();
void server_persist(void)             if (!accepting) {
{                                         n_accessing--;
    mutex_lock(&mutex);                   return 1;
    accepting = false;                }
    mfence();                         return 0;
    while (n_accessing != 0)      }
        ;
    do_persist();                 void server_leave(void)
    mfence();                     {
    accepting = true;                 mfence();
    mutex_unlock(&mutex);             n_accessing--;
}                                 }
```

**Figure 4**: Client-server protocol implementation.

When a client requests to enter the server, it increases n_accessing before it actually enters the OBSERVING state. Then it checks whether the server is in ACCEPTING by reading the variable accepting, and enters OBSERVING if accepting is indeed true. The transition from OBSERVING to COMMITTING does not change the value of n_accessing, as it represents the total number of clients in both states. A client decreases n_accessing by one when leaving the server. Similarly, when the server attempts to enter the PERSISTING state, it first sets accepting to false, which denotes that it is already in the WAITING state. Then it waits until n_accessing becomes zero and then moves to PERSISTING in order to do the actual task of *persist*. After finishing the task, the server sets accepting to true to indicate that it has transitioned back to ACCEPTING. We use memory fences to ensure the reads and writes to n_accessing and accepting follow the program order. The mechanism guarantees that (1) the server is in ACCEPTING at the time when client transitions to OB-

SERVING, and (2) |OBSERVING| + |COMMITTING| = 0 at the time when the server transitions to PERSISTING.

**Discussion.** It is known that ARIES-style WAL often performs better than shadow paging in systems where concurrent write transactions are supported [24, 36]. The reason is that with shadow paging, when a transaction is ready to write back the dirty pages at commit time, it must ensure that other transactions have not partially modified those pages, otherwise inconsistency may arise on a crash. Consequently, transactions have to frequently synchronize with each other. With weak durability, however, the need to urgently synchronize with other transactions no longer exists: it is only when users ask for durability that synchronization is required.

On the other hand, the main benefit of our specific shadow technique is *simplicity*, in that it has a simple specification, and does not depend on layers other than the file system. This means that we can separately test or even formally verify [9] the implementation against a simple specification. In contrast, WAL often depends on other layers like buffers, or assumes a certain locking protocol, making the already difficult task of testing or verifying crash safety even harder.

## 3.2 Latch-free two-level index structure

Weak durability offers a great opportunity for *batch processing*—operations within two consecutive *persist*s naturally form a batch. Moreover, weak durability frees us from worrying about the durability of operations generated in the last batch. AciKV embodies these ideas in its index design by employing a latch-free two-level index structure that scales well across a wide range of workloads.

As shown in Figure 2, the index consists of an in-memory latch-free concurrent skip list [22, 44] and an on-disk B+-tree. The interesting part of the two-level design is that it only supports merging a batch of records, rather than inserting a single record, into the B+-tree. Single-record insertions will go to the skip list. As insertions of records are absorbed by the skip list, and deletions, like many other systems [2, 45], are implemented as a special form of updates that set the records to a tombstone value, the B+-tree structure is guaranteed to remain the same within the same batch. This *structural invariance property* is appealing for two primary reasons, as described below.

First, the index does not require latches to protect the tree structure; along with the fact that the implementation of a concurrent skip list is simpler and more scalable than that of a concurrent B+-tree [44], this combination leads to better scalability. Second, records are guaranteed to remain at the same location within the same batch, which fits well with our use of the local write set (§3.1). With this guarantee, a transaction can now simply store the location of records they wish to update in its write set, and at commit time, apply the writes directly to the stored location if the batch has not changed since the time the transaction retrieves the record.
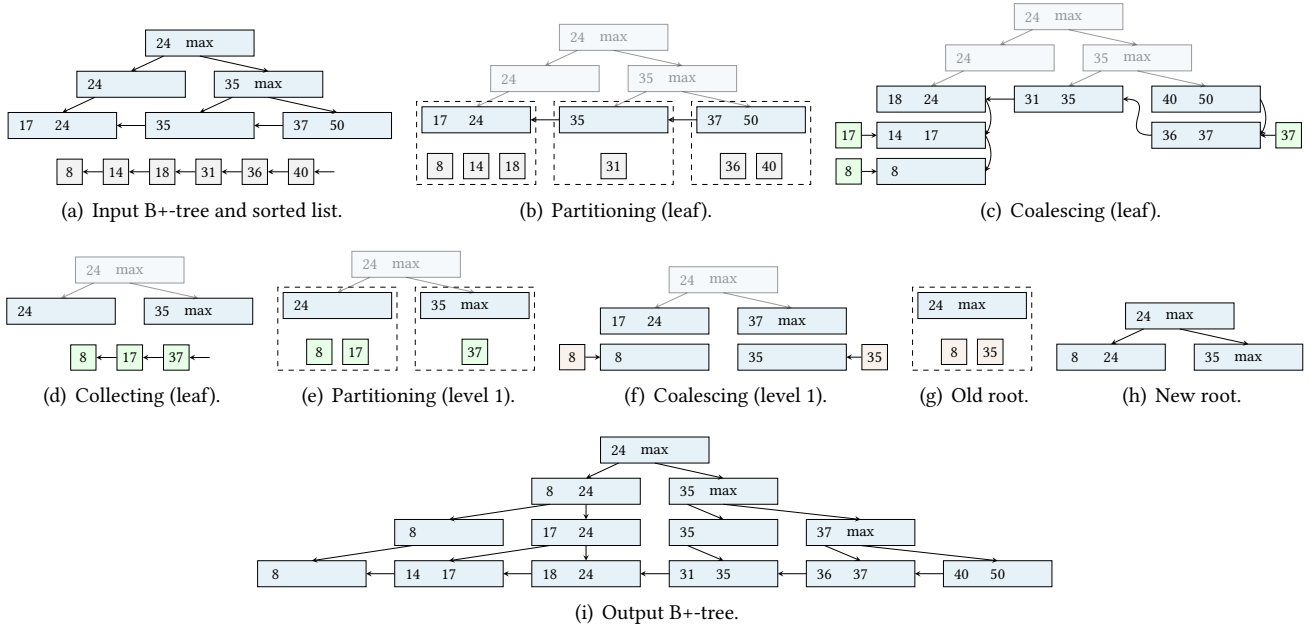
**Figure 5**: Merging the in-memory skip list into the on-disk B+-tree (each number shown in the figure represents a record, but we only show its key and omit its value). (a) The merging procedure takes a B+-tree and a sorted list as input. (b) Each record is assigned to the tree node whose range covers the key of the record; for instance, as the middle leaf node should contain keys within (24, 35], the record with key 31 is assigned to this node. (c) After the assignment, records are coalesced into each tree node. For nodes that still have space (e.g., the middle one), simple insertions suffice; otherwise, a merge sort is performed on the records to be inserted and records in the tree node. During the merge sort, new tree nodes are allocated and internal pointers (colored in green) pointing to the new nodes are generated. (d) The internal pointers are then collected to form the input of level 1. (e and f) The same procedure is then applied to level 1. (g) We have reached the root, but the root cannot accommodate the internal pointers (colored in brown), (h) so a new root is created. (i) shows the output B+-tree.

This removes the need to relocate records at commit time in most cases.

**Merging skip list into B+-tree.** On receiving a *persist*, AciKV merges the in-memory skip list into the on-disk B+-tree for durability. The merging algorithm targets one level at a time, starting from the leaf level and going upward to the root. At each level, the merging proceeds in three phases: *partitioning*, *coalescing*, and *collecting*, as illustrated in Figure 5 with an example. First, the input list is partitioned into multiple sublists, with each sublist containing keys to be coalesced into the same B+-tree node. Then, for each sublist, a dedicated worker thread inserts each record in the sublist into the tree node when there is enough space; otherwise, it performs a merge sort on the tree node and the sublist. When a new pointer is generated due to a split, it does not immediately propagate to the upper levels. Rather, the new pointer is collected in an output list, and linked with other output lists to produce the input for the next level.

A similar (and more detailed) merging algorithm is also described in [50], along with some optimizations such as balancing the inserting load and utilizing SIMD instructions to speed up sorting.

**Discussion.** The multi-level design of AciKV's index shares a similar structure to log-structured merge-trees (LSM-

trees) [39], and also appears in some research prototypes [54] and commercial systems [7]. But there are some notable differences in the index design of AciKV and of prior systems: First, AciKV only inserts new records into the skip list, and directly modifies records that are already present in the skip list or in the B+-tree. In contrast, other systems direct both insertions and modifications to their write-efficient data structure. As a result, the merging load of AciKV is smaller than that of prior systems when the workload involves a few updates. Another difference is that AciKV does not need a protective mechanism to ensure that the records in the skip list are not lost upon a crash, as they are inserted *after* the last *persist*. The permissiveness of weak durability is the key to enable this design choice.

### 3.3 Strong strict two-phase locking (SS2PL)

To ensure serializability, AciKV follows the standard SS2PL locking protocol to acquire shared and exclusive locks on records and on key gaps [35]. For completeness, we briefly describe the main idea here.

Before reading and writing a record, each transaction must obtain the shared and, respectively, exclusive record-lock on the key of the record. Additionally, a range query also has to make sure that the key range it concerns will not be affected by future insertions. This is done with a special form of

locks called *gap-locks*, which serve as "physical surrogates for logical properties" [26]. This means that by obtaining a gap-lock on key $k$ (a physical entity), the transaction essentially owns the range between $k$ and the key immediately preceding $k$ (a logical range). For instance, if a transaction wants to retrieve records whose key lies between 3 and 6, and currently the database has records with key 1, 4, and 8, then the transaction would have to acquire shared gap-locks on 4 and 8. Similarly, a transaction has to acquire an exclusive gap-lock for the range into which it wishes to insert the record. If the key exceeds the largest key currently existing in the database, then the transaction would have to acquire a gap-lock on a special sentinel key, which, in effect, locks the entire rightmost range. To avoid deadlocks, AciKV adopts the no-wait policy [59]; that is, when a transaction fails to acquire a lock, it simply aborts.

**Prefix preservation of SS2PL.** As mentioned in §2.2, to ensure crash consistency, we can use a scheduling mechanism that possesses the prefix preservation property (in addition to serializability). SS2PL indeed has this property, as it requires every transaction to hold its locks (both shared and exclusive) until the transaction terminates; consequently, any operation $op$ that depends on another operation $op'$ of a different transaction can only take effect after the transaction executing $op'$ commits.

### 3.4 Top-level operations

To put the pieces together, we go through how AciKV handles each primitive.

***get(k).*** The transaction first acquires a shared record-lock on $k$, searches for $k$ in the local write set and then the index, and returns $v$ if $(k, v)$ is found.

***getrange(k1, k2).*** The transaction first acquires a shared gap-lock on the smallest key that is greater than or equal to $k2$ in the index. Then it retrieves all records whose key belongs to the range $[k1, k2]$ from the local write set and the index, acquires a shared gap-lock and a shared record-lock on each record, and returns the retrieved records.

***put(k, v).*** The transaction first checks whether $k$ is present in the local write set; if so, it directly updates the corresponding entry and returns. If not, it acquires an exclusive record-lock on $k$, and search for key $k$ in the index. If $k$ is present in the index, it stores $(k, v)$ together with the location (which can be a node address, tagged as List, or a B+-tree record location, tagged as Tree) in its local write set. Otherwise, it acquires an exclusive gap-lock on the smallest key greater than or equal to $k$, allocate a local write set entry to store $(k, v)$, and tags the location of this entry as None.

***delete(k).*** The transaction acquires an exclusive record-lock on $k$, and searchs for $k$ in the local write set and the index. If $k$ is found, it effectively performs a *put* but the update value is set to a tombstone value; in AciKV, a tombstone value is set to any value whose length is zero. Otherwise, the transaction does nothing.

***begin.*** The transaction simply records the current epoch.

***commit.*** The transaction compares the current epoch with the one recorded in *begin*. If they match, the transaction directly applies the write set to the skip list or the B+-tree based on the location information stored in the local write set. If they do not match, meaning that a *persist* has happened in the midst of *begin* and *commit*, the locations are thus invalid as *persist* has merged the skip list into the B+-tree. In this case, the transaction has to search for key $k$ in the B+-tree for each write set entry which previously resides in the index (i.e., an entry whose location is tagged as List or Tree) to find out their new locations. With the new locations, the transaction can now apply the write set to the index. Finally, it releases all the acquired locks and pinned buffers, and resets the local write set.

***abort.*** The transaction simply releases all the acquired locks and pinned buffers, and resets the local write set.

***persist.*** The transaction merges the skip list into the B+-tree, writes back all the dirty buffers to the shadow paging layer, calls a *flush* to crash-atomically update the stable B+-tree, and finally advances the current epoch.

## 4 Evaluation

We ran experiments to answer the following questions:

- What are the trade-offs imposed by weak durability and by group commit? (§4.2)
- How does AciKV perform compared with existing disk-based database systems? (§4.3)
- What are the benefits of building AciKV around weak durability? (§4.4)
- What is the recovery time (§4.5) and memory overhead (§4.6) of AciKV?

### 4.1 Experimental setup

All experiments were done on a host machine with a 6-core 3.2 GHz Intel i7-8700 CPU and 16 GB of DRAM. All database files were stored on an ext4 file system hosted by a Samsung 980 Pro 500-GB NVMe SSD (except for §4.2 where we also use an HDD). We use the YCSB benchmark [17] to evaluate the performance of AciKV. The database contains 20M records, each of which consists of a 16B key and a 100B value. Details of each workload are described below:

**Read-or-write.** Each transaction reads or writes a record whose key is chosen randomly and uniformly from the entire database. The read ratio $r$ controls the percentage of read operations.
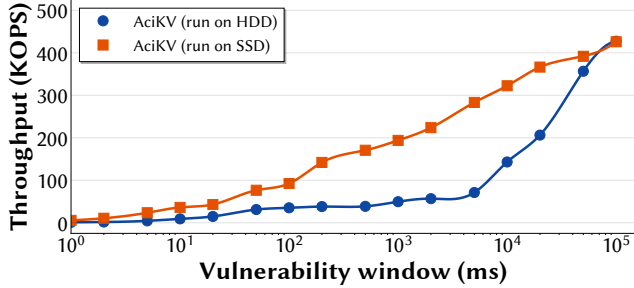
**Figure 6**: Weak durability: trade-off between throughput and vulnerability window.



**Figure 7**: Group commit: trade-off between throughput and latency.

**Insertion.** Each transaction inserts a record whose key is chosen randomly and uniformly from [0, 20M). For this workload, the database is initially empty.

**Range query.** Each transaction randomly and uniformly selects a key $k$ from the entire database, and reads $n$ records starting from $k$, where $n$ is chosen randomly and uniformly from [1, 100].

**Read-modify-write.** Each transaction reads a record whose key is chosen randomly and uniformly from the entire database, and then updates the value of the record.

## 4.2  Weak durability versus group commit

Weak durability and group commit share a similar goal: alleviating the high synchronization overhead due to the strong durability requirement. One interesting fact about AciKV's interface is that it can easily support group commit by delaying the notification that a *commit* has returned until the next *persist*. This way, transactions between two consecutive *persist*s naturally form a batch. This gives us a great opportunity to compare weak durability with group commit on a similar experimental setting, to better understand the trade-off imposed by each technique.

We begin with the configuration of weak durability. We employed a dedicated thread to invoke *persist* for every time interval $k$, which we refer to as the *vulnerability window* to reflect that only transactions committing within the most recent $k$-interval might disappear in the face of a crash. We used the write-only workload, and tuned the value of $k$ from 1ms to 100s. To understand AciKV's performance on slow storage media, we also used an HDD.

Figure 6 shows the throughput with respect to each $k$. As expected, a larger vulnerability window gives rise to a higher throughput. The results clearly indicate that durability requirements have a strong implication for performance. This experiment also demonstrates AciKV's flexibility compared with strong durability—users can pick a suitable value of $k$ for their application by replaying a similar experiment, perhaps with their own workloads. For slower media like HDDs, the throughput also increases at a slower rate. This means that to achieve the same target performance, faster media can have a smaller vulnerability window.
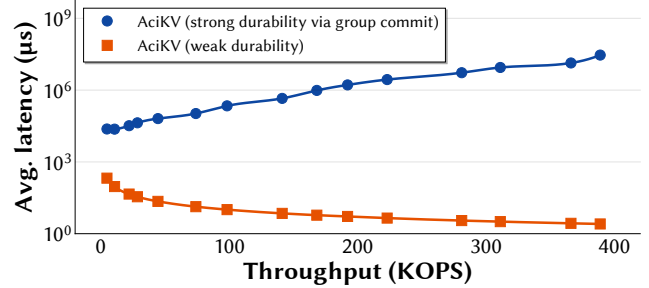
Next we study the trade-off imposed by group commit. Similarly to the prior setting, we tuned the value of $k$ from 1ms to 100s, but this time we reported the average latency against the throughput. We count the latency of weak durability (the orange line) as the time interval between a *begin-commit* pair, and the latency of group commit (the blue line) as the previous latency plus the waiting time, defined as the time interval between the *commit* and the next *persist* it encounters.

Figure 7 shows the results, indicating a clear trend: In the group commit case, the latency increases as the throughput increases, whereas in the weak durability case, the latency decreases as the throughput increases. This shows that group commit forces users to choose between a high throughput or a low latency. The latency difference between the two techniques is more obvious when aiming for a high throughput. For instance, in our experiment, for a target throughput of 300 KOPS, the latency of group commit is five orders of magnitude higher than that of weak durability. The merit of group commit is, of course, strong durability, so it has no vulnerability window.

In summary, with weak durability, users should choose an acceptable vulnerability window for a desired performance. With group commit, users are required to find a balance between throughput and latency. For all subsequent experiments, we set the vulnerability window of AciKV to 5s.

## 4.3  Comparing AciKV with existing systems

The main goal of AciKV is to guarantee transactional properties while maintaining a similar performance to non-transactional systems. Thus, we compare AciKV with six popular disk-based database systems, including three non-transactional systems[2]: LevelDB, RocksDB, and Kyoto Cabinet, and three transactional systems: Berkeley DB, LMDB, and SQLite3. Table 2 summarizes some key design choices of these systems:

We would like to emphasize that every system comes with its own design choice and tuning parameters, and without a total understanding of them, it is hard to make a completely fair comparison. We therefore ask the reader to ignore rela-

---

[2]RocksDB and Kyoto Cabinet also provide some transactional supports, but our experiment did not use them.
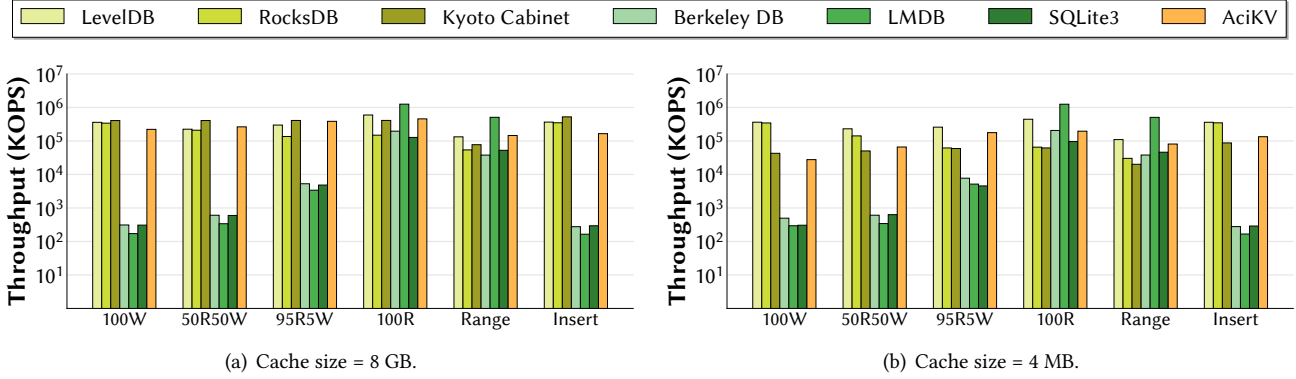
**Figure 8**: YCSB results. For LevelDB and RocksDB, we configured their write buffer size to 120 MB. For SQLite3, we set its journal mode to WAL. For all systems, we configured them in read-only mode for read-only workload.

| System | Interface | Index | Recovery |
|---|---|---|---|
| LevelDB | Non-txn KV | LSM-tree | WAL |
| RocksDB | Non-txn KV | LSM-tree | WAL |
| Kyoto Cabinet | Non-txn KV | B+-tree | WAL + SP |
| Berkeley DB | Txn KV | B+-tree | WAL |
| LMDB | Txn KV | B+-tree | CoW |
| SQLite3 | SQL | B+-tree | WAL |

**Table 2**: Design choices of database systems that are used in our experiment. WAL: write-ahead logging; SP: shadow paging; CoW: copy-on-write.

tively small performance difference, and focus more on the order of magnitude.

Figure 8(a) shows the throughput of each system under different workloads. The first four bar groups are the results of running the read-or-write workload with read ratio set to 0, 0.5, 0.95, and 1, respectively. When only writes are issued, the throughput of AciKV is 662x–1219x higher than that of the three transactional systems, and is within half of the non-transactional ones. The performance difference between AciKV and other transactional systems decreases as the read ratio increases: with the read ratio set to 0.5 and 0.95, the difference drops to 282x–763x and 52x–111x, respectively. When only reads are issued, durability is no longer a factor that can affect performance, and thus the results are roughly at the same order. This is also true for the range-query workload (the 5th bar groups). For the insertion workload (the 6th bar groups), AciKV is 454x–789x faster than the transactional systems, and at worst 4x slower than the non-transactional ones. The difference between the performance of AciKV's updates (209 KOPS) and insertions (140 KOPS) can be mainly attributed to the merging of the skip list and the B+-tree.

The previous experiment uses a large cache (8 GB) that can accommodate the entire database, so disk I/Os are mostly writes for durability. To understand how a smaller cache can affect the performance, we set the cache size to 4 MB (except for LMDB, which does not maintain its own cache, and instead relies on the page cache of file systems) and repeated the previous experiment.

Figure 8(b) shows the results. Our analysis is divided into three parts based on the workload types: (i) For read-only workloads (the 4th and 5th bar groups), a smaller cache degrades the performance of all systems by 25% to 85% (AciKV decreases by 58%). (ii) For write-only workloads (1st and 6th bar groups), all systems except for AciKV and Kyoto Cabinet are not affected by a smaller cache; however, we believe the reason for this is quite different for each system. For LevelDB and RocksDB, we suspect the reason lies in the append-only nature of their LSM-tree design, which does not require reading a page before writing a record, and hence no cache misses are triggered; for Berkeley DB and SQLite3, the reason is that the synchronization overhead due to strong durability is the primary bottleneck affecting the performance, so a few reads due to cache misses are almost negligible. AciKV and Kyoto Cabinet has lower performance with a smaller cache, which is about 10% of the performance when the cache can accommodate the entire database. The reason is that AciKV and Kyoto Cabinet need to, before writing a record, read the page containing the record, and thus can be affected by cache misses. But even with a smaller cache, AciKV still performs 90x–166x and 355x–671x better than the other transactional systems under update and insertion workloads, respectively. (iii) For read-write mix workloads (2nd and 3rd bar groups), the performance degradation due to cache misses is in between read-only and write-only workloads.

In summary, under workloads that involve writes, AciKV provides transactional guarantees while achieving similar performance to systems without transactional support. With a smaller cache, the throughput of AciKV drops because of cache misses, but is still much higher than that of transactional systems with a strong durability guarantee.

## 4.4 Multicore scalability

We built AciKV around weak durability. To understand the benefit of doing so, we compare AciKV to other transactional systems that consider strong durability as their "first-class" property, and then provide ad-hoc options to relax durability. As discussed in §1 and in §2.1, enabling these options is
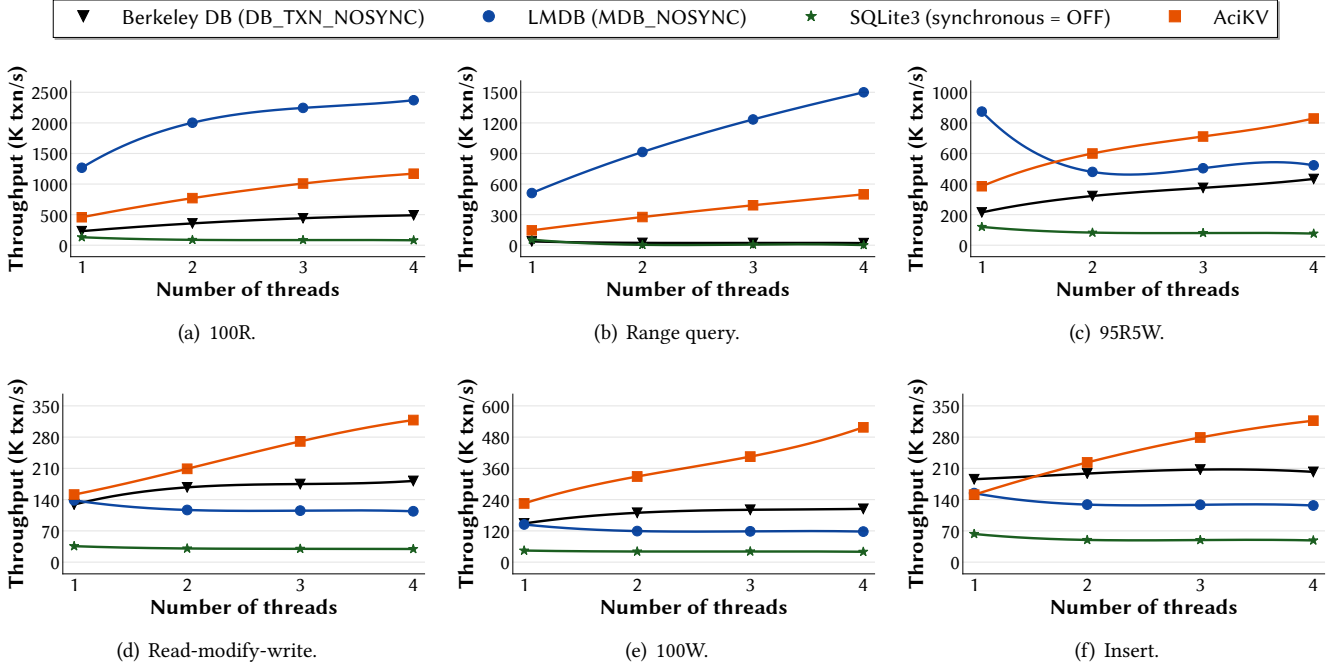
**Figure 9**: Multi-thread YCSB results.

subject to internal and external inconsistency in the event of a crash. This experiment focuses on the performance aspect, especially *multicore scalability*.

In this experiment, we chose the transactional systems Berkeley DB, LMDB, and SQLite3, and opened their own weakly durable option. We used the read-or-write, insertion, range query, and read-modify-write workloads to compare AciKV with the above systems. We changed the number of threads from 1 to 4, and reported the throughput.

Again, we divide our analysis based on workload types. Figure 9(a) and Figure 9(b) show the results of read-only workloads. Both AciKV and LMDB can scale well under these workloads. In particular, LMDB shows great read performance, partly because it is intentionally designed to be a read-optimized database system; for instance, its multi-version concurrency control mechanism frees its readers from acquiring any lock during execution. The throughput of LMDB, however, drops when the workloads involve some writes, as shown in Figure 9(c) and Figure 9(d). For write-only workloads, as shown in Figure 9(e) and Figure 9(f), AciKV, Berkeley DB, and LMDB have similar performance with a single thread, but AciKV performs better than the others when there are multiple threads concurrently writing the database. The primary reason why AciKV can scale under these workloads lies in its two-level latch-free design—by absorbing insertions with a concurrent latch-free skip list, transactions do not need to acquire any index latch when traversing the B+-tree.

In summary, AciKV is the only system that consistently scales across a wide range of workloads. Building a system
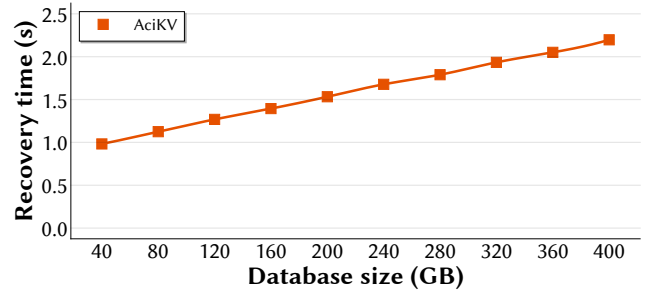


**Figure 10**: Recovery time with respect to different database sizes.

around weak durability not only prevents inconsistency, but also has the potential to discover a better design.

## 4.5 Recovery time

AciKV relies on shadow paging to ensure crash safety (§3.1). Unlike logging-based techniques whose recovery time usually depends on *when* a crash occurs (e.g., the recovery time for crashes appearing at the end of the log is likely to be longer than that at the start of the log), the shadow paging technique AciKV uses has a nearly constant recovery time regardless of when a crash occurs. Instead, the recovery time depends only on the size of the database. This means that it is easier to measure AciKV's recovery time.

Figure 10 shows the recovery time of AciKV with respect to different database sizes. The recovery time increases linearly with the database size. For a 400GB database, the recovery time is less than 3 seconds, which we believe is acceptable in most cases.

## 4.6 Memory overhead

There are two main sources of memory overhead in AciKV: one is the in-memory skip list (§3.2) and the other is the page table used in shadow (§3.1).

We first discuss how we determine the size of the skip list, which depends on two factors: the maximum insertion throughput and the size of each record. For instance, we allocated 10M skip list entries (each of which can accommodate one record) in the previous experiments; since the skip list is reset every time a *persist* is issued, and our persist interval is set to 5s, the maximum insertion throughput is thus 2M insertions per second. Given each record is about 120B, the memory overhead of the skip list is about 1.2 GB.

For the page table used in shadow, the memory overhead is roughly 1/1000 of the database size, as each 4B table entry maps a 4KB logical page to a physical one (which is also 4KB). There are two ways to reduce this overhead: we can use a larger page granularity, or we can store the page table on the database file and implement a caching mechanism for the page table.

## 5 Related work

**High-performance transactional systems.** Some systems are able to provide full ACID guarantees while maintaining a good performance [18, 20, 21, 55–57, 61]. While many of these are main-memory database systems, to sustain a high performance while ensuring durability, they still need an efficient way to cope with failures. For instance, Silo [57, 61] and Hekaton [20] adopt group commit to amortize the disk latency; FaRM [21] relies on NVM; H-Store [55], Spanner [18], and CockroachDB [56] failover to a replica upon system failures.

This work approaches the goal of retaining both the transaction abstraction and a good performance from a different angle, that is, by relaxing the durability requirement. Unlike prior systems that are often deployed as a cloud service, AciKV can be easily hosted on commodity hardware.

**Alleviating synchronization overhead.** Lower down the storage stack, file systems are also aware of the performance implication of durability requirements. OptFS [14], BarrierFS [58], Featherstitch [23], and Horae [32] propose to add lightweight ordering primitives to the file system interface, so that applications can use fsync only when they actually require durability. NoFS [13] proposes a novel technique to fully remove the need to enforce ordering constraints on writes, but it cannot implement atomic operations. Xsyncfs [38] adopts a user-centric approach: a write is guaranteed to be durable when an external observer sees the write. This approach creates for users the illusion that writes are synchronous and keeps the efficiency of asynchronous writes.

**ACID at lower layers.** ACID is not a unique concept in database systems; some file systems and storage devices also provide full or partial support to ACID. Examples include TxFS [28], Valor [52], Mime [10], the Logical Disk [19], Stasis [49], TxFlash [43], Beyond Block IO [40], Mars [16], LightTx [34], X-FTL [29], and Isotope [51].

While these systems have been shown to improve the design of their upper layers (in terms of simplicity and efficiency), one downside of this approach is that it often requires substantial changes to the existing software stack. AciKV shares a similar goal to these systems, that is, providing a useful abstraction to relieve the programming burden of users, but targets a higher abstraction level.

**Weak durability in the field.** Many storage systems are by default weakly durable. For instance, almost none of the major file systems guarantee that a `write` is made durable unless users explicitly issue an `fsync`. While weak durability substantially improves the performance of these systems, it also has some undesirable consequences: notably, (i) complex specification of how the systems should behave on a crash, and (ii) an enormous space of post-crash states, which in turn lead to many crash safety bugs [30, 37, 41, 60]. As a result, recent research has proposed to reduce the number of post-crash states resulting from weak durability with some ordering guarantees [8, 42], to formally specify crash behavior [4, 31, 46], or to formally verify crash safety [1, 5, 6, 9, 11, 12, 25] of their systems.

AciKV's weak durability is much more benign, in that it is complemented with other transactional guarantees. For instance, users can expect that their transactions are atomic with respect to crashes, and that their integrity constraints will not be accidentally violated in the face of a crash.

## 6 Conclusion

We introduce the notion of weakly durable transactions (ACID⁻) that aims to prevent isolation and consistency anomalies, while benefit from the high I/O parallelism granted by modern storage devices. We discuss, from a theoretical point of view, how to achieve ACID⁻, and use ACID⁻ as the new requirement to build AciKV. We conduct extensive experiments to evaluate AciKV; results show that AciKV can outperform existing transactional systems by orders of magnitude. We believe AciKV represents a useful middle ground between systems that provide high throughput but have few guarantees, and systems that expose the useful transaction abstraction but suffer from high synchronization overhead.

# References

[1] S. Amani, A. Hixon, Z. Chen, C. Rizkallah, P. Chubb, L. O'Connor, J. Beeren, Y. Nagashima, J. Lim, T. Sewell, J. Tuong, G. Keller, T. Murray, G. Klein, and G. Heiser. Cogent: Verifying high-assurance file system implementations. In *Proceedings of the Twenty-First International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS '16, pages 175–188, New York, NY, USA, 2016. Association for Computing Machinery.

[2] J. Arulraj, A. Pavlo, and S. R. Dulloor. Let's talk about storage & recovery methods for non-volatile memory database systems. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, SIGMOD '15, pages 707–722, New York, NY, USA, 2015. Association for Computing Machinery.

[3] P. A. Bernstein, V. Hadzilacos, and N. Goodman. *Concurrency Control and Recovery in Database Systems*. Addison-Wesley, 1987.

[4] J. Bornholt, A. Kaufmann, J. Li, A. Krishnamurthy, E. Torlak, and X. Wang. Specifying and checking file system crash-consistency models. In *Proceedings of the Twenty-First International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS '16, pages 83–98, New York, NY, USA, 2016. Association for Computing Machinery.

[5] T. Chajed, J. Tassarotti, M. F. Kaashoek, and N. Zeldovich. Argosy: Verifying layered storage systems with recovery refinement. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI '19, pages 1054–1068, New York, NY, USA, 2019. Association for Computing Machinery.

[6] T. Chajed, J. Tassarotti, M. F. Kaashoek, and N. Zeldovich. Verifying concurrent, crash-safe systems with perennial. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*, SOSP '19, pages 243–258, New York, NY, USA, 2019. Association for Computing Machinery.

[7] F. Chang, J. Dean, S. Ghemawat, W. C. Hsieh, D. A. Wallach, M. Burrows, T. Chandra, A. Fikes, and R. E. Gruber. Bigtable: A distributed storage system for structured data. In *7th USENIX Symposium on Operating Systems Design and Implementation*, OSDI '06, Seattle, WA, 2006. USENIX Association.

[8] Y.-S. Chang and R.-S. Liu. OPTR: Order-preserving translation and recovery design for SSDs with a standard block device interface. In *Proceedings of the 2019*

*USENIX Conference on Usenix Annual Technical Conference*, USENIX ATC '19, pages 1009–1023, USA, 2019. USENIX Association.

[9] Y.-S. Chang, Y. Hsiao, T.-C. Lin, C.-W. Tsao, C.-F. Wu, Y.-H. Chang, H.-S. Ko, and Y.-F. Chen. Determinizing crash behavior with a verified snapshot-consistent flash translation layer. In *14th USENIX Symposium on Operating Systems Design and Implementation*, OSDI '20, pages 81–97. USENIX Association, Nov. 2020.

[10] C. Chao, R. English, D. Jacobson, A. Stepanov, E. Stepanov, J. Wilkes, R. Wagner, and S. I. Mime: A high performance parallel storage device with strong recovery guarantees. Technical Report HPL−CSP−92−9 rev 1, Hewlett-Packard Laboratories, 1992.

[11] H. Chen, D. Ziegler, T. Chajed, A. Chlipala, M. F. Kaashoek, and N. Zeldovich. Using crash hoare logic for certifying the FSCQ file system. In *Proceedings of the 25th Symposium on Operating Systems Principles*, SOSP '15, pages 18–37, New York, NY, USA, 2015. Association for Computing Machinery.

[12] H. Chen, T. Chajed, A. Konradi, S. Wang, A. undefinedleri, A. Chlipala, M. F. Kaashoek, and N. Zeldovich. Verifying a high-performance crash-safe file system using a tree specification. In *Proceedings of the 26th Symposium on Operating Systems Principles*, SOSP '17, pages 270–286, New York, NY, USA, 2017. Association for Computing Machinery.

[13] V. Chidambaram, T. Sharma, A. C. Arpaci-Dusseau, and R. H. Arpaci-Dusseau. Consistency without ordering. In *Proceedings of the 10th USENIX Conference on File and Storage Technologies*, FAST '12, page 9, USA, 2012. USENIX Association.

[14] V. Chidambaram, T. S. Pillai, A. C. Arpaci-Dusseau, and R. H. Arpaci-Dusseau. Optimistic crash consistency. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, SOSP '13, pages 228–243, New York, NY, USA, 2013. Association for Computing Machinery.

[15] H. Chu. Lightning memory-mapped database manager (lmdb). http://www.lmdb.tech/doc/group_ _mdb.html, Apr. 2021.

[16] J. Coburn, T. Bunker, M. Schwarz, R. Gupta, and S. Swanson. From ARIES to MARS: Transaction support for next-generation, solid-state drives. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, SOSP '13, pages 197–212, New York, NY, USA, 2013. Association for Computing Machinery.

[17] B. F. Cooper, A. Silberstein, E. Tam, R. Ramakrishnan, and R. Sears. Benchmarking cloud serving systems with ycsb. In *Proceedings of the 1st ACM Symposium on Cloud Computing*, SoCC '10, pages 143–154, New York, NY, USA, 2010. Association for Computing Machinery.

[18] J. C. Corbett, J. Dean, M. Epstein, A. Fikes, C. Frost, J. Furman, S. Ghemawat, A. Gubarev, C. Heiser, P. Hochschild, W. Hsieh, S. Kanthak, E. Kogan, H. Li, A. Lloyd, S. Melnik, D. Mwaura, D. Nagle, S. Quinlan, R. Rao, L. Rolig, Y. Saito, M. Szymaniak, C. Taylor, R. Wang, and D. Woodford. Spanner: Google's globally-distributed database. In *10th USENIX Symposium on Operating Systems Design and Implementation*, OSDI '12, pages 261–264, Hollywood, CA, Oct. 2012. USENIX Association.

[19] W. de Jonge, M. F. Kaashoek, and W. C. Hsieh. The logical disk: A new approach to improving file systems. In *Proceedings of the Fourteenth ACM Symposium on Operating Systems Principles*, SOSP '93, pages 15–28, New York, NY, USA, 1993. Association for Computing Machinery.

[20] C. Diaconu, C. Freedman, E. Ismert, P.-A. Larson, P. Mittal, R. Stonecipher, N. Verma, and M. Zwilling. Hekaton: Sql server's memory-optimized oltp engine. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, SIGMOD '13, pages 1243–1254, New York, NY, USA, 2013. Association for Computing Machinery.

[21] A. Dragojević, D. Narayanan, E. B. Nightingale, M. Renzelmann, A. Shamis, A. Badam, and M. Castro. No compromises: Distributed transactions with consistency, availability, and performance. In *Proceedings of the 25th Symposium on Operating Systems Principles*, SOSP '15, pages 54–70, New York, NY, USA, 2015. Association for Computing Machinery.

[22] K. Fraser. *Practical lock-freedom*. PhD thesis, University of Cambridge, UK, 2004.

[23] C. Frost, M. Mammarella, E. Kohler, A. de los Reyes, S. Hovsepian, A. Matsuoka, and L. Zhang. Generalized file system dependencies. In *Proceedings of Twenty-First ACM SIGOPS Symposium on Operating Systems Principles*, SOSP '07, pages 307–320, New York, NY, USA, 2007. Association for Computing Machinery.

[24] J. Gray, P. McJones, M. Blasgen, B. Lindsay, R. Lorie, T. Price, F. Putzolu, and I. Traiger. The recovery manager of the system r database manager. *ACM Comput. Surv.*, 13(2):223–242, June 1981.

[25] T. Hance, A. Lattuada, C. Hawblitzel, J. Howell, R. Johnson, and B. Parno. Storage systems are distributed systems (so verify them that way!). In *14th USENIX Symposium on Operating Systems Design and Implementation*, OSDI '20, pages 99–115. USENIX Association, Nov. 2020.

[26] J. M. Hellerstein, M. Stonebraker, and J. Hamilton. Architecture of a database system. *Found. Trends Databases*, 1(2):141–259, Feb. 2007.

[27] D. Hitz, J. Lau, and M. Malcolm. File system design for an NFS file server appliance. In *Proceedings of the USENIX Winter 1994 Technical Conference on USENIX Winter 1994 Technical Conference*, WTEC '94, page 19, USA, 1994. USENIX Association.

[28] Y. Hu, Z. Zhu, I. Neal, Y. Kwon, T. Cheng, V. Chidambaram, and E. Witchel. Txfs: Leveraging filesystem crash consistency to provide ACID transactions. In *2018 USENIX Annual Technical Conference*, ATC '18, pages 879–891, Boston, MA, July 2018. USENIX Association.

[29] W.-H. Kang, S.-W. Lee, B. Moon, G.-H. Oh, and C. Min. X-FTL: Transactional FTL for SQLite databases. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, SIGMOD '13, pages 97–108, New York, NY, USA, 2013. Association for Computing Machinery.

[30] S. Kim, M. Xu, S. Kashyap, J. Yoon, W. Xu, and T. Kim. Finding semantic bugs in file systems with an extensible fuzzing framework. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*, SOSP '19, pages 147–161, New York, NY, USA, 2019. Association for Computing Machinery.

[31] M. Kokologiannakis, I. Kaysin, A. Raad, and V. Vafeiadis. PerSeVerE: Persistency semantics for verification under ext4. *Proc. ACM Program. Lang.*, 5(POPL), Jan. 2021.

[32] X. Liao, Y. Lu, E. Xu, and J. Shu. Write dependency disentanglement with HORAE. In *14th USENIX Symposium on Operating Systems Design and Implementation*, OSDI '20, pages 549–565. USENIX Association, Nov. 2020.

[33] R. A. Lorie. Physical integrity in a large segmented database. *ACM Trans. Database Syst.*, 2(1):91–104, Mar. 1977.

[34] Y. Lu, J. Shu, J. Guo, S. Li, and O. Mutlu. LightTx: A lightweight transactional design in flash-based SSDs to support flexible transactions. In *Proceedings of the 31st IEEE International Conference on Computer Design*, ICCD '13, pages 115–122. IEEE Computer Society, 2013.

[35] C. Mohan. ARIES/KVL: A key-value locking method for concurrency control of multiaction transactions operating on b-tree indexes. In *Proceedings of the Sixteenth International Conference on Very Large Databases*, VLDB '90, pages 392–405, San Francisco, CA, USA, 1990. Morgan Kaufmann Publishers Inc.

[36] C. Mohan, D. Haderle, B. Lindsay, H. Pirahesh, and P. Schwarz. Aries: A transaction recovery method supporting fine-granularity locking and partial rollbacks using write-ahead logging. *ACM Trans. Database Syst.*, 17(1):94–162, Mar. 1992.

[37] J. Mohan, A. Martinez, S. Ponnapalli, P. Raju, and V. Chidambaram. Finding crash-consistency bugs with bounded black-box crash testing. In *Proceedings of the 13th USENIX Conference on Operating Systems Design and Implementation*, OSDI '18, pages 33–50, USA, 2018. USENIX Association.

[38] E. B. Nightingale, K. Veeraraghavan, P. M. Chen, and J. Flinn. Rethink the sync. In *Proceedings of the 7th Symposium on Operating Systems Design and Implementation*, OSDI '06, pages 1–14, USA, 2006. USENIX Association.

[39] P. O'Neil, E. Cheng, D. Gawlick, and E. O'Neil. The log-structured merge-tree (LSM-tree). *Acta Inf.*, 33(4): 351–385, June 1996.

[40] X. Ouyang, D. Nellans, R. Wipfel, D. Flynn, and D. K. Panda. Beyond block I/O: Rethinking traditional storage primitives. In *2011 IEEE 17th International Symposium on High Performance Computer Architecture*, HPCA '11, pages 301–311, 2011.

[41] T. S. Pillai, V. Chidambaram, R. Alagappan, S. Al-Kiswany, A. C. Arpaci-Dusseau, and R. H. Arpaci-Dusseau. All file systems are not created equal: On the complexity of crafting crash-consistent applications. In *Proceedings of the 11th USENIX Conference on Operating Systems Design and Implementation*, OSDI '14, pages 433–448, USA, 2014. USENIX Association.

[42] T. S. Pillai, R. Alagappan, L. Lu, V. Chidambaram, A. C. Arpaci-Dusseau, and R. H. Arpaci-Dusseau. Application crash consistency and performance with CCFS. *ACM Trans. Storage*, 13(3), Sept. 2017.

[43] V. Prabhakaran, T. L. Rodeheffer, and L. Zhou. Transactional flash. In *Proceedings of the 8th USENIX Conference on Operating Systems Design and Implementation*, OSDI '08, pages 147–160. USENIX Association, 2008.

[44] W. Pugh. Concurrent maintenance of skip lists. Technical report, University of Maryland at College Park, USA, 1990.

[45] P. Raju, R. Kadekodi, V. Chidambaram, and I. Abraham. PebblesDB: Building Key-Value Stores using Fragmented Log-Structured Merge Trees. In *Proceedings of the 26th ACM Symposium on Operating Systems Principles*, SOSP '17, Shanghai, China, October 2017.

[46] T. Ridge, D. Sheets, T. Tuerk, A. Giugliano, A. Madhavapeddy, and P. Sewell. SibylFS: Formal specification and oracle-based testing for POSIX and real-world file systems. In *Proceedings of the 25th Symposium on Operating Systems Principles*, SOSP '15, pages 38–53, New York, NY, USA, 2015. Association for Computing Machinery.

[47] O. Rodeh. B-trees, shadowing, and clones. *ACM Trans. Storage*, 3(4), Feb. 2008.

[48] O. Rodeh, J. Bacik, and C. Mason. Btrfs: The linux b-tree filesystem. *ACM Trans. Storage*, 9(3), Aug. 2013.

[49] R. Sears and E. Brewer. Stasis: Flexible transactional storage. In *Proceedings of the 7th Symposium on Operating Systems Design and Implementation*, OSDI '06, pages 29–44, USA, 2006. USENIX Association.

[50] J. Sewall, J. Chhugani, C. Kim, N. Satish, and P. Dubey. PALM: Parallel architecture-friendly latch-free modifications to b+ trees on many-core processors. *Proc. VLDB Endow.*, 4(11):795–806, Aug. 2011.

[51] J.-Y. Shin, M. Balakrishnan, T. Marian, and H. Weatherspoon. Isotope: Transactional isolation for block storage. In *14th USENIX Conference on File and Storage Technologies*, FAST '16, pages 23–37, Santa Clara, CA, Feb. 2016. USENIX Association.

[52] R. P. Spillane, S. Gaikwad, M. Chinni, E. Zadok, and C. P. Wright. Enabling transactional file access via lightweight kernel extensions. In *7th USENIX Conference on File and Storage Technologies*, FAST '09, San Francisco, CA, Feb. 2009. USENIX Association.

[53] SQLite. How to corrupt an sqlite database file. https://www.sqlite.org/howtocorrupt.html#_failure_to_sync, Apr. 2021.

[54] M. Stonebraker, D. J. Abadi, A. Batkin, X. Chen, M. Cherniack, M. Ferreira, E. Lau, A. Lin, S. Madden, E. O'Neil, P. O'Neil, A. Rasin, N. Tran, and S. Zdonik. C-Store: A column-oriented DBMS. In *Proceedings of the 31st International Conference on Very Large Data Bases*, VLDB '05, pages 553–564. VLDB Endowment, 2005.

[55] M. Stonebraker, S. Madden, D. J. Abadi, S. Harizopoulos, N. Hachem, and P. Helland. The end of an architectural era: (it's time for a complete rewrite). In *Proceedings of the 33rd International Conference on Very Large Data Bases*, VLDB '07, pages 1150–1160. VLDB Endowment, 2007.

[56] R. Taft, I. Sharif, A. Matei, N. VanBenschoten, J. Lewis, T. Grieger, K. Niemi, A. Woods, A. Birzin, R. Poss, P. Bardea, A. Ranade, B. Darnell, B. Gruneir, J. Jaffray, L. Zhang, and P. Mattis. CockroachDB: The resilient geo-distributed SQL database. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, SIGMOD '20, pages 1493–1509, New York, NY, USA, 2020. Association for Computing Machinery.

[57] S. Tu, W. Zheng, E. Kohler, B. Liskov, and S. Madden. Speedy transactions in multicore in-memory databases. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, SOSP '13, pages 18–32, New York, NY, USA, 2013. Association for Computing Machinery.

[58] Y. Won, J. Jung, G. Choi, J. Oh, S. Son, J. Hwang, and S. Cho. Barrier-enabled io stack for flash storage. In *Proceedings of the 16th USENIX Conference on File and Storage Technologies*, FAST '18, pages 211–226, USA, 2018. USENIX Association.

[59] X. Yu, G. Bezerra, A. Pavlo, S. Devadas, and M. Stonebraker. Staring into the abyss: An evaluation of concurrency control with one thousand cores. *Proc. VLDB Endow.*, 8(3):209–220, Nov. 2014.

[60] M. Zheng, J. Tucek, D. Huang, F. Qin, M. Lillibridge, E. S. Yang, B. W. Zhao, and S. Singh. Torturing databases for fun and profit. In *Proceedings of the 11th USENIX Conference on Operating Systems Design and Implementation*, OSDI '14, pages 449–464, USA, 2014. USENIX Association.

[61] W. Zheng, S. Tu, E. Kohler, and B. Liskov. Fast databases with fast durability and recovery through multicore parallelism. In *11th USENIX Symposium on Operating Systems Design and Implementation*, OSDI '14, pages 465–477, Broomfield, CO, Oct. 2014. USENIX Association.