# Protocol-Aware Recovery for Consensus-Based Storage

Ramnatthan Alagappan, Aishwarya Ganesan, Eric Lee[†], Aws Albarghouthi,
Vijay Chidambaram[†], Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau

*University of Wisconsin – Madison*    [†] *University of Texas at Austin*

## Abstract

We introduce *protocol-aware recovery* (PAR), a new approach that exploits protocol-specific knowledge to correctly recover from storage faults in distributed systems. We demonstrate the efficacy of PAR through the design and implementation of corruption-tolerant replication (CTRL), a PAR mechanism specific to replicated state machine (RSM) systems. We experimentally show that the CTRL versions of two systems, LogCabin and ZooKeeper, safely recover from storage faults and provide high availability, while the unmodified versions can lose data or become unavailable. We also show that the CTRL versions have little performance overhead.

## 1 Introduction

Failure recovery using redundancy is central to improved reliability of distributed systems [14, 22, 31, 35, 61, 67]. Distributed systems recover from node crashes and network failures using copies of data and functionality on several nodes [6, 47, 55]. Similarly, bad or corrupted data on one node should be recovered from redundant copies.

In a *static* setting where all nodes always remain reachable and where clients do not actively update data, recovering corrupted data from replicas is straightforward; in such a setting, a node could repair its state by simply fetching the data from any other node.

In reality, however, a distributed system is a *dynamic* environment, constantly in a state of flux. In such settings, orchestrating recovery correctly is surprisingly hard. As a simple example, consider a quorum-based system, in which a piece of data is corrupted on one node. When the node tries to recover its data, some nodes may fail and be unreachable, some nodes may have recently recovered from a failure and so lack the required data or hold a stale version. If enough care is not exercised, the node could "fix" its data from a stale node, overwriting the new data, potentially leading to a data loss.

To correctly recover corrupted data from redundant copies in a distributed system, we propose that a recovery approach should be *protocol-aware*. A *protocol-aware recovery* (PAR) approach is carefully designed based on how the distributed system performs updates to its replicated data, elects the leader, etc. For instance, in the previous example, a PAR mechanism would realize that a faulty node has to query at least $R$ (read quorum) other nodes to safely and quickly recover its data.

In this paper, we apply PAR to replicated state machine (RSM) systems. We focus on RSM systems for two reasons. First, correctly implementing recovery is most challenging for RSM systems because of the strong consistency and durability guarantees they provide [58]; a small misstep in recovery could violate the guarantees. Second, the reliability of RSM systems is crucial: many systems entrust RSM systems with their critical data [45]. For example, Bigtable, GFS, and other systems [7,26] store their metadata on RSM systems such as Chubby [16] or ZooKeeper [4]. Hence, protecting RSM systems from storage faults such as data corruption will improve the reliability of many dependent systems.

We first characterize the different approaches to handling storage faults by developing the *RSM recovery taxonomy*, through experimental and qualitative analysis of practical systems and methods proposed by prior research (§2). Our analyses show that most approaches employed by currently deployed systems do not use any protocol-level knowledge to perform recovery, leading to disastrous outcomes such as data loss and unavailability.

Thus, to improve the resiliency of RSM systems to storage faults, we design a new protocol-aware recovery approach that we call corruption-tolerant replication or CTRL (§3). CTRL constitutes two components: a *local storage layer* and a *distributed recovery protocol*; while the storage layer reliably detects faults, the distributed protocol recovers faulty data from redundant copies. Both the components carefully exploit RSM-specific knowledge to ensure safety (e.g., no data loss) and high availability.

CTRL applies several novel techniques to achieve safety and high availability. For example, a *crash-corruption disentanglement* technique in the storage layer distinguishes corruptions caused by crashes from disk faults; without this technique, safety violations or unavailability could result. Next, a *global-commitment determination* protocol in the distributed recovery separates committed items from uncommitted ones; this separation is critical: while recovering faulty committed items is necessary for safety, discarding uncommitted items quickly is crucial for availability. Finally, a novel *leader-initiated snapshotting* mechanism enables identical snapshots across nodes to greatly simplify recovery.

We implement CTRL in two storage systems that are based on different consensus algorithms: LogCabin [43]

(based on Raft [50]) and ZooKeeper [4] (based on ZAB [39]) (§4). Through experiments, we show that CTRL versions provide safety and high availability in the presence of storage faults, while the original systems remain unsafe or unavailable in many cases; we also show that CTRL induces minimal performance overhead (§5).

## 2 Background and Motivation

We first provide background on storage faults and RSM systems. We then present the taxonomy of different approaches to handling storage faults in RSM systems.

### 2.1 Storage Faults in Distributed Systems

Disks and flash devices exhibit a subtle and complex failure model: a few blocks of data could become inaccessible or be silently corrupted [8, 9, 32, 59]. Although such storage faults are rare compared to whole-machine failures, in large-scale distributed systems, even rare failures become prevalent [60, 62]. Thus, it is critical to reliably detect and recover from storage faults.

Storage faults occur due to several reasons: media errors [10], program/read disturbance [60], and bugs in firmware [9], device drivers [66], and file systems [27, 28]. Storage faults manifest in two ways: block *errors* and *corruption*. Block errors (or latent sector errors) arise when the device internally detects a problem with a block and throws an error upon access. Studies of both flash [33, 60] and hard drives [10, 59] show that block errors are common. Corruption could occur due to lost and misdirected writes that may not be detected by the device. Studies [9, 51] and anecdotal evidence [36, 37, 57] show the prevalence of data corruption in the real world.

Many local file systems, on encountering a storage fault, simply propagate the fault to applications [11, 54, 64]. For example, ext4 silently returns corrupted data if the underlying device block is corrupted. In contrast, a few file systems transform an underlying fault into a different one; for example, btrfs returns an error to applications if the accessed block is corrupted on the device. In either case, storage systems built atop local file systems should handle corrupted data and storage errors to preserve end-to-end data integrity.

One way to tackle storage faults is to use RAID-like storage to maintain multiple copies of data on each node. However, many distributed deployments would like to use inexpensive disks [22, 31]. Given that the data in a distributed system is inherently replicated, it is wasteful to store multiple copies on each node. Hence, it is important for distributed systems to use the inherent redundancy to recover from storage faults.

### 2.2 RSM-based Storage Systems

Our goal is to harden RSM systems to storage faults. In an RSM system, a set of nodes compute identical states by executing commands on a state machine (an in-memory data structure on each node) [58]. Typically, clients interact with a single node (the leader) to execute operations on the state machine. Upon receiving a command, the leader durably writes the command to an on-disk *log* and replicates it to the followers. When a majority of nodes have durably persisted the command in their logs, the leader applies the command to its state machine and returns the result to the client; at this point, the command is committed. The commands in the log have to be applied to the state machine *in-order*. Losing or overwriting committed commands violates the safety property of the state machine. The replicated log is kept consistent across nodes by a consensus protocol such as Paxos [41] or Raft [50].

Because the log can grow indefinitely and exhaust disk space, periodically, a *snapshot* of the in-memory state machine is written to disk and the log is garbage collected. When a node restarts after a crash, it restores the system state by reading the latest on-disk snapshot and the log. The node also recovers its critical metadata (e.g., log start index) from a structure called *metainfo*. Thus, each node maintains three critical persistent data structures: the *log*, the *snapshots*, and the *metainfo*.

These persistent data structures could be corrupted due to storage faults. Practical systems try to safely recover the data and remain available under such failures [15, 17]. However, as we will show, none of the current approaches correctly recover from storage faults, motivating the need for a new approach.

### 2.3 RSM Recovery Taxonomy

To understand the different possible ways to handling storage faults in RSM systems, we analyze a broad range of approaches. We perform this analysis by two means: first, we analyze practical systems including ZooKeeper, LogCabin, etcd [25], and a Paxos-based system [24] using a fault-injection framework we developed (§5); second, we analyze techniques proposed by prior research or used in proprietary systems [15, 17].

Through our analysis, we classify the approaches into two categories: *protocol-oblivious* and *protocol-aware*. The oblivious approaches do not use any protocol-level knowledge to perform recovery. Upon detecting a fault, these approaches take a recovery action locally on the faulty node; such actions interact with the distributed protocols in unsafe ways, leading to data loss. The protocol-aware approaches use some RSM-specific knowledge to recover; however, they do not use this knowledge correctly, leading to undesirable outcomes. Our taxonomy is *not* complete in that there may be other techniques; however, to the best of our knowledge, we have not observed other approaches apart from those in our taxonomy.
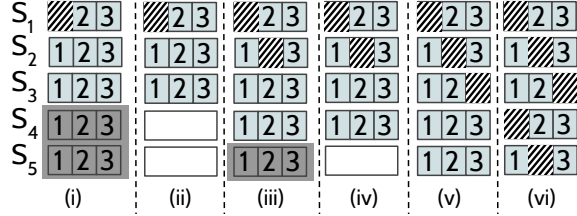
Figure 1: **Sample Scenarios.** *The figure shows sample scenarios in which current approaches fail. Faulty entries are striped. Crashed and lagging nodes are shown as gray and empty boxes, respectively.*

| Class | Approach | Safety | Availability | Performance | No Intervention | No extra nodes | Fast Recovery | Low Complexity | (i) | (ii) | (iii) | (iv) | (v) | (vi) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Protocol Oblivious | *NoDetection* | × | √ | √ | √ | √ | na | √ | E | E | E | E | E | E |
| | *Crash* | √ | × | √ | × | √ | na | √ | U | C | U | C | U | U |
| | *Truncate* | × | √ | √ | √ | √ | × | √ | C | L | C | L | L | L |
| | *DeleteRebuild* | × | √ | √ | × | √ | × | √ | C | L | C | L | L | L |
| Protocol Aware | *MarkNonVoting* | × | × | √ | √ | √ | × | √ | U | C | U | C | U | U |
| | *Reconfigure* | √ | × | √ | × | × | × | √ | U | C | U | C | U | U |
| | *Byzantine FT* | √ | × | × | √ | × | na | × | U | C | U | U | U | U |
| | CTRL | √ | √ | √ | √ | √ | √ | √ | C | C | C | C | C | C |

E- Return Corrupted, L- Data Loss, U- Unavailable, C- Correct

Table 1: **Recovery Taxonomy.** *The table shows how different approaches behave in Figure 1 scenarios. While all approaches are unsafe or unavailable, CTRL ensures safety and high availability.*

To illustrate the problems, we use Figure 1. In all cases, log entries[†] 1, 2, and 3 are committed; losing these items will violate safety. Table 1 shows how each approach behaves in Figure 1's scenarios. As shown in the table, all current approaches lead to safety violation (e.g., data loss), low availability, or both. A recovery mechanism that effectively uses redundancy should be safe and available in all cases. Table 1 also compares the approaches along other axes such as performance, maintenance overhead (intervention and extra nodes), recovery time, and complexity. Although Figure 1 shows only faults in the *log*, the taxonomy applies to other structures including the snapshots and the metainfo.

**NoDetection.** The simplest reaction to storage faults is none at all: to trust every layer in the storage stack to work reliably. For example, a few prototype Paxos-based systems [24] do not use checksums for their on-disk data; similarly, LogCabin does not protect its snapshots with checksums. *NoDetection* trivially violates safety; corrupted data can be obliviously served to clients. However, deployed systems do use checksums and other integrity strategies for most of their on-disk data.

**Crash.** A better strategy is to use checksums and handle I/O errors, and crash the node on detecting a fault. *Crash* may seem like a good strategy because it intends to prevent any damage that the faulty node may inflict on the system. Our experiments show that the *Crash* approach is common: LogCabin, ZooKeeper, and etcd crash sometimes when their logs are faulty. Also, ZooKeeper crashes when its snapshots are corrupted.

Although *Crash* preserves safety, it suffers from severe unavailability. Given that nodes could be unavailable due to other failures, even a single storage fault results in unavailability, as shown in Figure 1(i). Similarly, a single fault even in different portions of data on a majority (e.g., Figure 1(v)) renders the system unavailable. Note that simply restarting the node does not help; storage faults, unlike other faults, could be persistent: the node will encounter the same fault and crash again until manual intervention, which is error-prone and may cause a data loss. Thus, it is desirable to recover automatically.

**Truncate.** A more sophisticated action is to truncate

---

[†]A log entry contains a state-machine command and data.

(possibly faulty) portions of data and continue operating. The intuition behind *Truncate* is that if the faulty data is discarded, the node can continue to operate (unlike *Crash*), improving availability.

However, we find that *Truncate* can cause a safety violation (data loss). Consider the scenario shown in Figure 2 in which entry 1 is corrupted on $S_1$; $S_4$, $S_5$ are lagging and do not have any entry. Assume $S_2$ is the leader. When $S_1$ reads its log, it detects the corruption; however, $S_1$ truncates its log, losing the corrupted entry and all subsequent entries (Figure 2(ii)). Meanwhile, $S_2$ (leader) and $S_3$ crash. $S_1$, $S_4$, and $S_5$ form a majority and elect $S_1$ the leader. Now the system does not have any knowledge of committed entries 1, 2, and 3, resulting in a *silent data loss*. The system also commits new entries $x$, $y$, and $z$ in the place of 1, 2, and 3 (Figure 2(iii)). Finally, when $S_2$ and $S_3$ recover, they follow $S_1$'s log (Figure 2(iv)), completely removing entries 1, 2, and 3.

In summary, although the faulty node detects the corruption, it truncates its log, losing the data locally. When this node forms a majority along with other nodes that are lagging, data is silently lost, violating safety. We find this safety violation in ZooKeeper and LogCabin.

Further, *Truncate* suffers from *inefficient recovery*. For instance, in Figure 1(i), $S_1$ truncates its log after a fault, losing entries 1, 2, and 3. Now to fix $S_1$'s log, the leader needs to transfer *all* entries, increasing $S_1$'s recovery time and wasting network bandwidth. ZooKeeper and LogCabin suffer from this slow recovery problem.

**DeleteRebuild.** Another commonly employed action is to manually delete all data on the faulty node and restart the node. Unfortunately, similar to *Truncate*, *DeleteRebuild* can violate safety; specifically, a node whose data is deleted could form a majority along with the lagging nodes, leading to a silent data loss. Surprisingly, administrators often use this approach hoping that the faulty

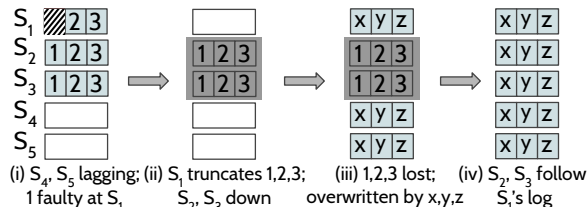S₁ / 2 3 → (blank) → x y z → x y z
(The figure shows log states — I'll render the caption)

**Figure 2: Safety Violation Example.** *The figure shows the sequence of events which exposes a safety violation in Truncate.*

(i) S₄, S₅ lagging; 1 faulty at S₁; (ii) S₁ truncates 1,2,3; S₂, S₃ down; (iii) 1,2,3 lost; overwritten by x,y,z; (iv) S₂, S₃ follow S₁'s log

node will be "simply fixed" by fetching the data from other nodes [63, 65, 73]. *DeleteRebuild* also suffers from the slow recovery problem similar to *Truncate*.

**MarkNonVoting.** In this approach, used by a Paxos-based system at Google [17], a faulty node deletes all its data on a fault and marks itself as a non-voting member; the node does not participate in elections until it observes one round of consensus and rebuilds its data from other nodes. By marking a faulty node as non-voting, safety violations such as the one in Figure 2 are avoided. However, *MarkNonVoting* can sometimes violate safety as noted by prior work [70]. The underlying reason for unsafety is that a corrupted node deletes all its state including the promises<sup>†</sup> given to leaders. Once a faulty node has lost its promise given to a new leader, it could accept an entry from an old leader (after observing a round of consensus on an earlier entry). The new leader, however, still believes that it has the promise from the faulty node and so can overwrite the entry, previously committed by the old leader.

Further, this approach suffers from unavailability. For example, when only a majority of nodes are alive, a single fault can cause unavailability because the faulty node cannot vote; other nodes cannot now elect a leader.

**Reconfigure.** In this approach, a faulty node is removed and a new node is added. However, to change the configuration, a configuration entry needs to be committed by a majority. Hence, the system remains unavailable in many cases (for example, when a majority are alive but one node's data is corrupted). Although *Reconfigure* is not used in practical systems to tackle storage faults, it has been suggested by prior research [15, 44].

**BFT.** An extreme approach is to use a Byzantine-fault-tolerant algorithm which should theoretically tolerate storage faults. However, *BFT* is expensive to be used in practical storage systems; specifically, *BFT* can achieve only half the throughput of what a crash-tolerant protocol can achieve [21]. Moreover, *BFT* requires $3f + 1$ nodes to tolerate $f$ faults [2], thus remaining unavailable in most scenarios in Figure 1.

**Taxonomy Summary.** None of the current approaches effectively use redundancy to recover from storage faults.

---

<sup>†</sup>In Paxos, a promise for a proposal numbered $p$ is a guarantee given by a follower (acceptor) to the leader (proposer) that it will not accept a proposal numbered less than $p$ in the future [41].

Most approaches do not use any protocol-level knowledge to recover; for example, *Truncate* and *DeleteRebuild* take actions locally on the faulty node and so interact with the distributed protocol in unsafe ways, causing a global data loss. Although some approaches (e.g., *MarkNonVoting*) use some RSM-specific knowledge, they do not do so correctly, causing data loss or unavailability. Thus, to ensure safety and high availability, a recovery approach should effectively use redundancy by exploiting protocol-specific knowledge. Further, it is beneficial to avoid other problems such as manual intervention and slow recovery. Our protocol-aware approach, CTRL, aims to achieve these goals.

# 3 Corruption-Tolerant Replication

Designing a correct recovery mechanism needs a careful understanding of the underlying protocols of the system. For example, the recovery mechanism should be cognizant of how updates are performed on the replicated data and how the leader is elected. We base CTRL's design on the following important protocol-level observations common to most RSM systems.

**Leader-based.** A single node acts as the leader; all data updates flow only through the leader.

**Epochs.** RSM systems partition time into logical units called *epochs*. For any given epoch, only one leader is guaranteed to exist. Every data item is associated with the epoch in which it was appended and its *index* in the log. Since the entries could only be proposed by the leader and only one leader could exist for an epoch, an ⟨*epoch, index*⟩ pair uniquely identifies a log entry.

**Leader Completeness.** A node will not vote for a candidate if it has more up-to-date data than the candidate. Since committed data is present at least in a majority of nodes and a majority vote is required to win the election, the leader is guaranteed to have all the committed data.

CTRL exploits these protocol-level attributes common to RSM systems to correctly recover from storage faults. CTRL divides the recovery responsibility between two components: the *local storage layer* and the *distributed recovery protocol*; while the storage layer reliably detects faulty data on a node, the distributed protocol recovers the data from redundant copies. Both the components use RSM-specific knowledge to perform their functions.

In this section, we first describe CTRL's fault model (§3.1) and safety and availability guarantees (§3.2). We then describe the local storage layer (§3.3). Finally, we describe CTRL's distributed recovery in two parts: first, we show how faulty *logs* are recovered (§3.4) and then we explain how faulty *snapshots* are recovered (§3.5).

## 3.1 Fault Model

Our fault model includes the standard failure assumptions made by crash-tolerant RSM systems: nodes could

| | Fault Outcome | Possible Causes |
|---|---|---|
| **Data** | corrupted data | misdirected and lost writes in ext |
| | inaccessible data | LSE, corruptions in ZFS and btrfs |
| **FS Metadata** | missing files/directories | directory entry corrupted, fsck may remove a faulty inode |
| | unopenable files/directories | sanity check fails after inode corruption, permission bits corrupted |
| | files with more or fewer bytes | $i\_size$ field in the inode corrupted |
| | file system read-only | journal corrupted; fsck not run |
| | file system unmountable | superblock corrupted; fsck not run |

**Table 2:** **Storage Fault Model.** *The table shows storage faults included in our model and possible causes that lead to a fault outcome.*

crash at any time and recover later, and nodes could be unreachable due to network failures [21, 42, 50]. Our model adds another realistic failure scenario where persistent data on the individual nodes could be corrupted or inaccessible. Table 2 shows a summary of our storage fault model. Our model includes faults in both user data and the file-system metadata blocks.

User data blocks in the files that implement the system's persistent structures could be affected by errors or corruption. A number of (possibly contiguous) data blocks could be faulty as shown by studies [12,59]. Also, a few bits/bytes of a block could be corrupted. Depending on the local file system in use, corrupted data may be returned obliviously or transformed into errors.

File-system metadata blocks can also be affected by faults; for example, the inode of a log file could be corrupted. Our fault model considers the following outcomes that can be caused by file-system metadata faults: files/directories may go missing, files/directories may be unopenable, a file may appear with fewer or more bytes, the file system may be mounted read-only, and in the worst case, the file system may be unmountable. Some file systems such as ZFS may mask most of the above outcomes from applications [72]; however, our model includes these faulty outcomes because they could realistically occur on other file systems that provide weak protection against corruption (e.g., ext2/3/4). Through fault-injection tests, we have verified that the metadata fault outcomes shown in Table 2 do occur on ext4.

## 3.2 Safety and Availability Guarantees

CTRL guarantees that if there exists at least one correct copy of a *committed* data item, it will be recovered or the system will wait for that item to be fixed; committed data will never be lost. In unlikely cases where all copies of a committed item are faulty, the system will correctly remain unavailable. CTRL also guarantees that the system will make a decision about an *uncommitted* faulty item as early as possible, ensuring high availability.

## 3.3 CTRL Local Storage Layer

To reliably recover, the storage layer (CLSTORE) needs to satisfy three key requirements. First, CLSTORE must be able to reliably detect a storage fault. Second,
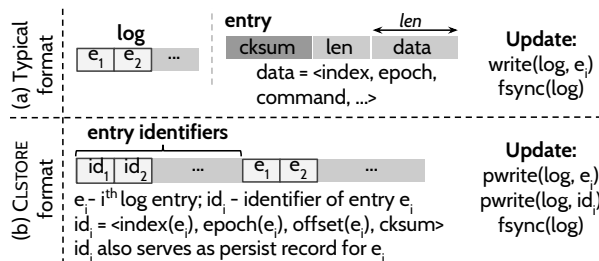


**Figure 3:** **Log Format.** *(a) shows the format and update protocol of a typical RSM log; (b) shows the same for* CLSTORE.

CLSTORE must correctly distinguish crashes from corruptions; safety can be violated otherwise. Third, CLSTORE must identify which pieces of data are faulty; only if CLSTORE identifies which pieces have been affected, can the distributed protocol recover those pieces.

### 3.3.1 Persistent Structures Overview

As we discussed, RSM systems maintain three persistent structures: the log, the snapshots, and the metainfo. CLSTORE uses RSM-specific knowledge of how these structures are used and updated, to perform its functions. For example, CLSTORE detects faults at a different granularity depending on the RSM data structure: faults in the log are detected at the granularity of individual entries, while faults in the snapshot are detected at the granularity of chunks. Similarly, CLSTORE uses the RSM-specific knowledge that a log entry is uniquely qualified by its $\langle epoch, index \rangle$ pair to identify faulty log entries.

**Log.** The log is a set of files containing a sequence of entries. The format of a typical RSM log is shown in Figure 3(a). The log is updated synchronously in the critical path; hence, changes made to the log format should not affect its update performance. CLSTORE uses a modified format as shown in Figure 3(b) which achieves this goal. A corrupted log is recovered at the granularity of individual entries.

**Snapshots.** The in-memory state machine is periodically written to a snapshot. Since snapshots can be huge, CLSTORE splits them into chunks; a faulty snapshot is recovered at the granularity of individual chunks.

**Metainfo.** The metainfo is special in that faulty metainfo *cannot* be recovered from other nodes. This is because the metainfo contains information unique to a node (e.g., its current epoch); recovering metainfo obliviously from other nodes could violate safety. CLSTORE uses this knowledge correctly and so maintains two copies of the metainfo locally; if one copy is faulty, the other copy is used. Fortunately, the metainfo is only a few tens of bytes in size and is updated infrequently; therefore, maintaining two copies does not incur significant overheads.

### 3.3.2 Detecting Faulty Data

CLSTORE uses well-known techniques for detection: *inaccessible* data is detected by catching return codes (e.g.,

EIO) and *corrupted* data is detected by a checksum mismatch. CLSTORE assumes that if an item and its checksum agree, then the item is not faulty. In the log, each entry is protected by a checksum; similarly, each chunk in a snapshot and the entire metainfo are checksummed.

CLSTORE also handles file-system metadata faults. Missing and unopenable files/directories are detected by handling error codes upon `open`. Log and metainfo files with fewer or more bytes are detected easily because these files are preallocated and are of a fixed size; snapshot sizes are stored separately, and CLSTORE cross-checks the stored size with the file-system reported size to detect discrepancies. A read-only/unmountable file system is equivalent to a missing data directory. In most cases of file-system metadata faults, CLSTORE crashes the nodes. Crashing reliably on a metadata fault preserves safety but compromises on availability. However, we believe this is an acceptable behavior because there are far more data blocks than metadata blocks; therefore, the probability of faults is significantly less for metadata than data blocks.

### 3.3.3 Disentangling Crashes and Corruption in Log

An interesting challenge arises when detecting corruptions in the log. A checksum mismatch for a log entry could occur due to two different situations. First, the system could have *crashed* in the middle of an update; in this case, the entry would be partially written and hence cause a mismatch. Second, the entry could be safely persisted but *corrupted* at a later point. Most log-based systems conflate these two cases: they treat a mismatch as a crash [30]. On a mismatch, they discard the corrupted entry and all subsequent entries, losing the data. Discarding entries due to such conflation introduces the possibility of a global data loss (as shown earlier in Figure 2).

Note that if the mismatch were really due to a crash, it is safe to discard the partially written entry. It is safe because the node would not have acknowledged to any external entity that it has written the entry. However, if an entry is *corrupted*, the entry cannot be simply discarded since it could be globally committed. Further, if a mismatch can be correctly attributed to a crash, the faulty entry can be quickly discarded locally, avoiding the distributed recovery. Hence, it is important for the local storage layer to distinguish the two cases.

To denote the completion of an operation, many systems write a commit record [13, 18]. Similarly, CLSTORE writes a persist record, $p_i$, after writing an entry $e_i$. For now, assume that $e_i$ is ordered before $p_i$, i.e., the sequence of steps to append an entry $e_i$ is $write(e_i), fsync(), write(p_i), fsync()$. On a checksum mismatch for $e_i$, if $p_i$ is not present, we can conclude that the system crashed during the update. Conversely, if $p_i$ is present, we can conclude that the mismatch was caused due to a corrup-

tion and *not* due to a crash. $p_i$ is checksummed and is very small; it can be atomically written and thus cannot be "corrupted" due to a crash. If $p_i$ is corrupted in addition to $e_i$, we can conclude that it is a corruption and not a crash.

The above logic works when $e_i$ is ordered before $p_i$. However, such ordering requires an (additional) expensive *fsync* in the critical path, affecting log-update performance. For this reason, CLSTORE does not order $e_i$ before $p_i$; thus, the append protocol is $\mathbf{t_1}$:$write(e_i)$, $\mathbf{t_2}$:$write(p_i)$, $\mathbf{t_3}$:$fsync()$.[†] Given this update sequence, assume a checksum mismatch occurs for $e_i$. If $p_i$ is *not present*, CLSTORE can conclude that it is a crash (before $\mathbf{t_2}$) and discard $e_i$. Contrarily, if $p_i$ is *present*, there are two possibilities: either $e_i$ could be affected by a corruption after $\mathbf{t_3}$ or a crash could have occurred between $\mathbf{t_2}$ and $\mathbf{t_3}$ in which $p_i$ hit the disk while $e_i$ was only partially written. The second case is possible because file systems can reorder writes between two *fsync* operations and $e_i$ could span multiple sectors [3, 19, 52, 53]. CLSTORE can still conclude that it is a corruption if $e_{i+1}$ or $p_{i+1}$ is present. However, if $e_i$ is the *last entry*, then we cannot determine whether it was a crash or a corruption.[*]

The inability to disentangle the last entry when its persist record is present is not specific to CLSTORE, but rather a fundamental limitation in log-based systems. For instance, in ext4's journal_async_commit mode (where a transaction is not ordered before its commit record), a corrupted last transaction is assumed to be caused due to a crash, possibly losing data [38, 69]. Even if crashes and corruptions can be disentangled, there is little a single-machine system can do to recover the corrupted data. However, in a distributed system, redundant copies can be used to recover. Thus, when the last entry cannot be disentangled, CLSTORE safely marks the entry as *corrupted* and leaves it to the distributed recovery to fix or discard the entry based on the global commitment.

The entanglement problem does not arise for snapshots or metainfo. These files are first written to a temporary file and then atomically renamed. If a crash happens before the rename, the partially written temporary file is discarded. Thus, the system will never see a corrupted snapshot or metainfo due to a crash; if these structures are corrupted, it is because of a storage corruption.

### 3.3.4 Identifying Faulty Data

Once a faulty item is detected, it has to be *identified*; only if CLSTORE can identify a faulty item, the distributed layer can recover the item. For this purpose, CLSTORE redundantly stores an *identifier* of an item apart from the item itself; duplicating only the identifier instead of the whole item obviates the ($2\times$) storage and performance

---

[†]The final *fsync* is required for durability.
[*]The proof of this claim is available [1].

overhead. However, storing the identifier near the item is less useful; a misdirected write can corrupt both the item and its identifier [9, 10]. Hence, identifiers are physically separated from the items they identify.

The $\langle epoch, index \rangle$ pair serves as the identifier for a log entry and is stored separately at the head of the log, as shown in Figure 3(b). The offset of an entry is also stored as part of the identifier to enable traversal of subsequent entries on a fault. The identifier of a log entry also conveniently serves as its persist record. Similarly, for a snapshot chunk, the $\langle snap\text{-}index, chunk\# \rangle$ pair serves as the identifier; the *snap-index* and the snapshot size are stored in a separate file than the snapshot file. The identifiers have a nominal storage overhead (32 bytes for log entries and 12 bytes for snapshots), can be atomically written, and are also protected by a checksum.

It is highly unlikely an item and its identifier will both be faulty since they are physically separated [9, 10, 12, 59]. In such unlikely and unfortunate cases, CLSTORE crashes the node to preserve safety. Table 3 (second column) summarizes CLSTORE's key techniques.

## 3.4 CTRL Distributed Log Recovery

The local storage layer detects faulty data items and passes on their identifiers to the distributed recovery layer. We now describe how the distributed layer recovers the identified faulty items from redundant copies using RSM-specific knowledge. We first describe how *log entries* are recovered and subsequently describe *snapshot* recovery. As we discussed, metainfo files are recovered locally and so we do not discuss them any further. We use Figure 4 to illustrate how log recovery works.

**Naive Approach: Leader Restriction.** RSM systems do not allow a node with an incomplete log to become the leader. A naive approach to recovering from storage faults could be to impose an additional constraint on the election: *a node cannot be elected the leader if its log contains a faulty entry*. The intuition behind the naive approach is as follows: since the leader is guaranteed to have all committed data and our new restriction ensures that the leader is not faulty, faulty log entries on other nodes could be fixed using the corresponding entries on the leader. Cases (a)(i) and (a)(ii) in Figure 4 show scenarios where the naive approach could elect a leader. In (a)(i), only $S_1$ can become the leader because other nodes are either lagging or have at least one faulty entry. Assume $S_1$ is the leader also in case (a)(ii).

**Fixing Followers' Logs.** When the leader has no faulty entries, fixing the followers is straightforward. For example, in case (a)(i), the followers inform $S_1$ of their faulty entries; $S_1$ then supplies the correct entries. However, sometimes the leader might not have any knowledge of an entry that a follower is querying for. For instance, in case (a)(ii), $S_5$ has a faulty entry at index 3 but

with a *different epoch*. This situation is possible because $S_5$ could have been the leader for epoch 2 and crashed immediately after appending an entry. As discussed earlier, an entry is uniquely identified by its $\langle epoch, index \rangle$; thus, when querying for faulty entries, a node needs to specify the epoch of the entry in addition to its index. Thus, $S_5$ informs the leader that its entry $\langle epoch{:}2, index{:}3 \rangle$ is faulty. However, $S_1$ does not have such an entry in its log. If the leader does not have an entry that the follower has, then the entry *must be an uncommitted entry* because the leader is guaranteed to have all committed data; thus, the leader instructs $S_5$ to truncate the faulty entry and also replicates the correct entry.

Although the naive approach guarantees safety, it has availability problems. The system will be unavailable in cases such as the ones shown in (b): a leader cannot be elected because the logs of the alive nodes are either faulty or lagging. Note that even a single storage fault can cause an unavailability as shown in (b)(i). It is possible for a carefully designed recovery protocol to provide better availability in these cases. Specifically, since at least one intact copy of all committed entries exists, it is possible to collectively reconstruct the log.

### 3.4.1 Removing the Restriction Safely

To recover from scenarios such as those in Figure 4(b), we remove the additional constraint on the election. Specifically, any node that has a more up-to-date log can now be elected the leader even if it has faulty entries. This relaxation improves availability; however, two key questions arise: first, when can the faulty leader proceed to accept new commands? second, and more importantly, is it safe to elect a faulty node as the leader?

To accept a new command, the leader has to append the command to its log, replicate it, and apply it to the state machine. However, before applying the new command, *all* previous commands must be applied. Specifically, faulty commands cannot be skipped and later applied when they are fixed; such out-of-order application would violate safety. Hence, it is required for the leader to fix its faulty entries before it can accept new commands. Thus, for improved availability, the leader needs to fix its faulty entries as early as possible.

The crucial part of the recovery to ensure safety is to fix the leader's log using the redundant copies on the followers. In simple cases such as (b)(i) and (b)(ii), the leader $S_1$ could fix its faulty entry $\langle epoch{:}1, index{:}1 \rangle$ using the correct entries from the followers and proceed to normal operation. However, in several scenarios, the leader cannot immediately recover its faulty entries; for example, none of the reachable followers might have any knowledge of the entry to be recovered or the entry to be recovered could also be faulty on the followers.
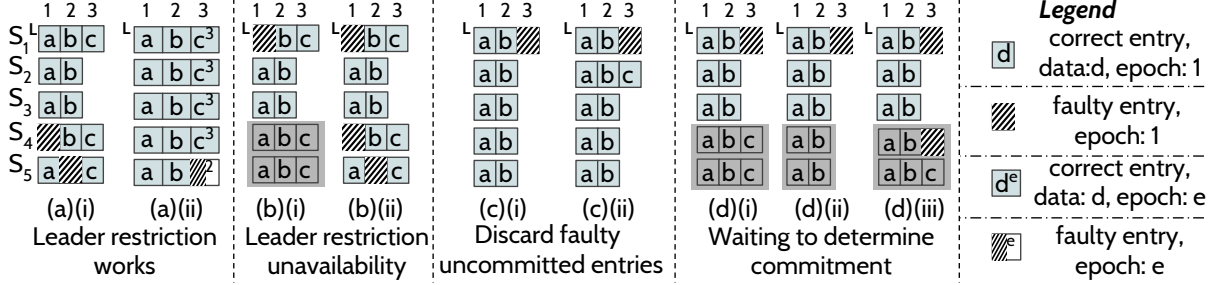
Figure 4: **Distributed Log Recovery.** *The figure shows how* CTRL's *log recovery operates. All entries are appended in epoch 1 unless explicitly mentioned. For entries appended in other epochs, the epoch number is shown in the superscript. Entries shown as striped boxes are faulty. A gray box around a node denotes that it is down or extremely slow. The leader is marked with L on the left. Log indexes are shown at the top.*

### 3.4.2 Determining Commitment

The main insight to fix the leader's faulty log safely and quickly is to distinguish *uncommitted* entries from *possibly committed* ones; while recovering the committed entries is necessary for safety, uncommitted entries can be safely discarded. Further, discarding uncommitted faulty entries immediately is crucial for availability. For instance, in case (c)(i), the faulty entry on $S_1$ cannot be fixed since there are no copies of it; waiting to fix that entry results in indefinite unavailability. Sometimes, an entry could be partially replicated but remain uncommitted; for example, in case (c)(ii), the faulty entry on $S_1$ is partially replicated but is not committed. Although there is a possibility of recovering this entry from the other node ($S_2$), this is *not* necessary for safety; it is completely safe for the leader to discard this uncommitted entry.

To determine the commitment of a faulty entry, the leader queries the followers. If a majority of the followers respond that they do *not* have the entry (negative acknowledgment), then the leader concludes that the entry is uncommitted. In this case, the leader safely discards that and all subsequent entries; it is safe to discard the subsequent entries because entries are committed in order. Conversely, if the entry were committed, at least one node in this majority would have that entry and inform the leader of it; in this case, the leader can fix its faulty entry using that response.

**Waiting to Determine Commitment**. Sometimes, it may be impossible for the leader to quickly determine commitment. For instance, consider the cases in Figure 4(d) in which $S_4$ and $S_5$ are down or slow. $S_1$ queries the followers to recover its entry $\langle epoch{:}1, index{:}3 \rangle$. $S_2$ and $S_3$ respond that they do not have such an entry (negative acknowledgment). $S_4$ and $S_5$ do not respond because they are down or slow. The leader, in this case, has to wait for either $S_4$ or $S_5$ to respond; discarding the entry without waiting for $S_4$ or $S_5$ could violate safety. However, once $S_4$ or $S_5$ responds, the leader will make a decision immediately. In (d)(i), $S_4$ or $S_5$ would respond with the correct entry, fixing the leader. In (d)(ii), $S_4$ or $S_5$ would respond that it does not have the entry, accu-

mulating three (a majority out of five) negative acknowledgments; hence, the leader can conclude that the entry is uncommitted, discard it, and continue to normal operation. In (d)(iii), $S_4$ would respond that it has the entry but is faulty in its log too. In this case, the leader has to wait for the response from $S_5$ to determine commitment. In the unfortunate and unlikely case where all copies of an entry are faulty, the system will remain unavailable.

### 3.4.3 The Complete Log Recovery Protocol

We now assemble the pieces of the log recovery protocol. First, fixing faulty followers is straightforward; the committed faulty entries on the followers can be eventually fixed by the leader because the leader is guaranteed to have all committed data. Faulty entries on followers that the leader does not know about are uncommitted; hence, the leader instructs the followers to discard such entries.

The main challenge is thus fixing the leader's log. The leader queries the followers to recover its entry $\langle epoch{:}e, index{:}i \rangle$. Three types of responses are possible:

Response 1: **have** – a follower could respond that it has the entry $\langle epoch{:}e, index{:}i \rangle$ and is not faulty in its log.

Response 2: **dontHave** – a follower could respond that it does not have the entry $\langle epoch{:}e, index{:}i \rangle$.

Response 3: **haveFaulty** – a follower could respond that it has $\langle epoch{:}e, index{:}i \rangle$ but is faulty in its log too.

Once the leader collects these responses from the followers, it takes the following possible actions:

Case 1: if it gets a *have* response from at least one follower, it fixes the entry in its log.

Case 2: if it gets a *dontHave* response from a majority of followers, it confirms that the entry is uncommitted, discards that entry and all subsequent entries.

Case 3: if it gets a *haveFaulty* response from a follower, it waits for either Case 1 or Case 2 to happen.

Case 1 and Case 2 can happen in any order; both orderings are safe. Specifically, if the leader decides to discard the faulty entry (after collecting a majority *dontHave* responses), it is safe since the entry was uncommitted anyways. Conversely, there is no harm in accepting a correct entry (at least one *have* response) and replicating it. The

first to happen out of these two cases will take precedence over the other.

The leader proceeds to normal operation only after its faulty data is discarded or recovered. However, CTRL discards uncommitted data as early as possible and minimizes the recovery latency by recovering faulty data at a fine granularity (as we show in §5.2), ensuring that the leader proceeds to normal operation quickly.

The leader could crash or be partitioned while recovering its log. On a leader failure, the followers will elect a new leader and make progress. The partial repair done by the failed leader is harmless: it could have either fixed committed faulty entries or discarded uncommitted ones, both of which are safe.

## 3.5 CTRL Distributed Snapshot Recovery

Because the logs can grow indefinitely, periodically, the in-memory state machine is written to disk and the logs are garbage collected. Current systems including ZooKeeper and LogCabin do not handle faulty snapshots correctly (§2.3): they either crash or load corrupted snapshots obliviously. CTRL aims to recover faulty snapshots from redundant copies. Snapshot recovery is different from log recovery in that all data in a snapshot is committed and already applied to the state machine; hence, faulty snapshots cannot be discarded in any case (unlike uncommitted log entries which can be discarded safely).

### 3.5.1 Leader-Initiated Identical Snapshots

Current systems [43] have two properties with respect to snapshots. First, they allow new commands to be applied to the state machine while a snapshot is in progress. Second, they take index-consistent snapshots: a snapshot $S_i$ represents the state machine in which log entries exactly up to $i$ have been applied. One of the mechanisms used in current systems to realize the above two properties is to take snapshots in a `fork`-ed child process; while the child can write an index-consistent image to the disk, the parent can keep applying new commands to its copy of the state machine. CTRL should enable snapshot recovery while preserving the above two properties.

In current systems, every node runs the snapshot procedure independently, taking snapshots at different log indexes. Because the snapshots are taken at different indexes, snapshot recovery can be complex: a faulty snapshot on one node cannot be simply fetched from other nodes. Further, snapshots cannot be recovered at the granularity of chunks because they will be byte-wise non-identical; entire snapshots have to be transferred across nodes, slowing down recovery.

This complexity can be significantly alleviated if the nodes take the snapshot at the same index; identical snapshots also enable chunk-based recovery.

However, coordinating a snapshot operation across nodes can, in general, affect the common-case perfor-

| | **Local Storage** | **Distributed Recovery** |
|---|---|---|
| *Log* | granularity: entry; identifier: $\langle epoch, index \rangle$; crash-corruption disentanglement | global-commitment determination to fix leader, leader fixes followers |
| *Snapshot* | granularity: chunk; identifier: $\langle snap\text{-}index, chunk\# \rangle$; no entanglement | leader-initiated identical snapshots, chunk-based recovery |
| *Metainfo* | granularity: file; identifier: n/a; no entanglement | none (only internal redundancy) |

Table 3: **Techniques Summary.** *The table shows a summary of techniques employed by* CTRL*'s storage layer and distributed recovery.*

mance. For example, one naive way to realize identical snapshots is for the leader to produce the snapshot, insert it into the log as yet another entry, and replicate it. However, such an approach will affect update performance since the snapshot could be huge and all client commands must wait while the snapshot commits [49]. Moreover, transferring the snapshot to the followers wastes network bandwidth.

CTRL takes a different approach to identical snapshots that preserves common-case performance. The leader initiates the snapshot procedure by first deciding the index at which a snapshot will be taken and informing the followers of the index. Once a majority agree on the index, all nodes independently take a snapshot at the index. When the leader learns that a majority (including itself) have taken a snapshot at an index $i$, it garbage collects its log up to $i$ and instructs the followers to do the same.

CTRL implements the above procedure using the log. When the leader decides to take a snapshot, it inserts a special marker called `snap` into the log. When the `snap` marker commits, and thus when a node applies the marker to the state machine, it takes a snapshot (i.e., the snapshot corresponds to the state where commands exactly up to the marker have been applied). Within each node, we reuse the same mechanism used by the original system (e.g., a `fork`-ed child) to allow new commands to be applied while a snapshot is in progress. Notice that the snapshot operation happens independently on all nodes but the operation will produce identical snapshots because the marker will be seen at the same log index by all nodes when it is committed. When the leader learns that a majority of nodes (including itself) have taken a snapshot at an index $i$, it appends another marker called `gc` for $i$; when the `gc` marker is committed and applied, the nodes garbage collect their log entries up to $i$.

### 3.5.2 Recovering Snapshot Chunks

With the identical-snapshot mechanism, snapshot recovery becomes easier. Once a faulty snapshot is detected, the local storage layer provides the distributed protocol the snapshot index and the chunk that is faulty. The distributed protocol recovers the faulty chunk from other

**Table 4 (a) Targeted Corruptions**

| System | Recovery Scenario | Total Test Cases | Original Approach | Original Outcomes Unavailable | Unsafe | Correct | CTRL Outcomes Unavailable | Unsafe | Correct |
|---|---|---|---|---|---|---|---|---|---|
| LogCabin | Possible | 2401 | truncate | 0 | 2355 | 46 | 0 | 0 | 2401 |
| | | | crash | 2355 | 0 | 46 | 0 | 0 | 2401 |
| | Not possible | 1695 | truncate | 0 | 1695 | 0 | 1695 | 0 | 0 |
| | | | crash | 1695 | 0 | 0 | 1695 | 0 | 0 |
| ZooKeeper | Possible | 2401 | truncate | 0 | 2355 | 46 | 0 | 0 | 2401 |
| | | | crash | 2355 | 0 | 46 | 0 | 0 | 2401 |
| | Not possible | 1695 | truncate | 0 | 1695 | 0 | 1695 | 0 | 0 |
| | | | crash | 1695 | 0 | 0 | 1695 | 0 | 0 |

**Table 4 (b) Random Block Corruptions and Errors**

| System | Experiment | Total Test Cases | Original Outcomes Unavailable | Unsafe | Correct | CTRL Outcomes Unavailable | Unsafe | Correct |
|---|---|---|---|---|---|---|---|---|
| LogCabin | Corruptions | 5000 | 738 | 793 | 3469 | 0 | 0 | 5000 |
| | Errors | 5000 | 2497 | 0 | 2503 | 0 | 0 | 5000 |
| ZooKeeper | Corruptions | 5000 | 807 | 656 | 3537 | 0 | 0 | 5000 |
| | Errors | 5000 | 2469 | 0 | 2531 | 0 | 0 | 5000 |

**Table 4 (c) Corruptions with Lagging Nodes**

| System | Total Test Cases | Original Outcomes Unavailable | Unsafe | Correct | CTRL Outcomes Unavailable | Unsafe | Correct |
|---|---|---|---|---|---|---|---|
| LogCabin | 5000 | 4194 | 141 | 665 | 0 | 0 | 5000 |
| ZooKeeper | 5000 | 1306 | 1806 | 1888 | 0 | 0 | 5000 |

(a) Targeted Corruptions  (b) Random Block Corruptions and Errors  (c) Corruptions with Lagging Nodes

Table 4: **Log Recovery.** *(a) shows results for targeted corruptions; we trigger two policies (truncate and crash) in the original systems. (b) shows results for random block corruptions and errors. (c) shows results for random corruptions with crashed and lagging nodes.*

nodes. First, the leader recovers its faulty chunks from the followers and then fixes the faulty snapshots on followers. Three cases arise during snapshot recovery.

First, the log entries for a faulty snapshot may not be garbage collected yet; in this case, the snapshot is recovered locally from the log (after fixing the log if needed).

Second, if the log is garbage collected, then a faulty snapshot has to be recovered from other nodes. However, if the log entries for a snapshot are garbage collected, then at least a majority of the nodes must have taken the same snapshot. This is true because the `gc` marker is inserted only after a majority of nodes have taken the snapshot. Thus, faulty garbage-collected snapshots are recovered from those redundant copies.

Third, sometimes, the leader may not know a snapshot that a follower is querying for (for example, if a follower took a snapshot and went offline for a long time and the leader replaced that snapshot with an advanced one); in this case, the leader supplies the full advanced snapshot.

### 3.6 CTRL Summary

The storage layer detects and identifies faulty data. Atop the storage layer, the distributed protocol recovers the faulty items from redundant copies. Both the layers exploit RSM-specific knowledge to correctly perform their functions. A summary of CTRL's local storage and distributed recovery techniques is shown in Table 3.

## 4 Implementation

We implement CTRL in two different RSM systems, LogCabin (v1.0) and ZooKeeper (v3.4.8); while LogCabin is based on Raft, ZooKeeper is based on ZAB. Implementing CTRL's storage layer and distributed recovery took only a moderate developer effort; CTRL adds about 1500 lines of code to each of the base systems.

### 4.1 Local Storage Layer

We implemented CLSTORE by modifying the storage engines of LogCabin and ZooKeeper. In both systems, the log is a set of files, each of a fixed size and preallocated. The header of each file is reserved for the log-entry identifiers. The size of the reserved header is proportional to the file size. CLSTORE ensures that a log entry and its identifier are at least a few megabytes physically apart. Both systems batch many log entries to improve update performance. With batching, CLSTORE performs crash-corruption disentanglement as follows: the first faulty entry without an identifier and its subsequent entries are discarded; faulty entries preceding that point are marked as corrupted and passed on to the distributed layer.

In both systems, the state machine is a data tree. We modified both the systems to take index-consistent identical snapshots: when a `snap` marker is applied, the state machine (i.e., the tree) is serialized to the disk. The *snap-index* and snapshot size are stored separately. CLSTORE uses a chunk size of 4K, enabling fine-grained recovery.

In LogCabin, the metainfo contains the `currentTerm` and `votedFor` structures. Similarly, in ZooKeeper, structures such as `acceptedEpoch` and `currentEpoch` constitute the metainfo. CLSTORE stores redundant copies of metainfo and protects them using checksums.

Log entries, snapshot chunks, and metainfo are protected by a CRC32 checksum. CLSTORE detects inaccessible data items by catching errors (`EIO`); it then populates the item's in-memory buffer with zeros, causing a checksum mismatch. Thus, CLSTORE deals with both corruptions and errors as checksum mismatches.

### 4.2 Distributed Recovery

**LogCabin.** In Raft, *terms* are equivalent to epochs. Thus, a log entry is uniquely identified by its ⟨*term, index*⟩ pair. To fix the followers, we modified the Append-Entries RPC used by the leader to replicate entries [50]. The followers inform the leader of their faulty log entries and snapshot chunks in the responses of this RPC; the leader sends the correct entries and chunks in a subsequent RPC. A follower starts applying commands to its state machine once its faulty data is fixed. To fix the

leader, we added a new RPC which the leader issues to the followers. The leader does not proceed to normal operation until its faulty data is fixed. After a configurable recovery timeout, the leader steps down if it is unable to recover its faulty data (for example, due to a partition), allowing other nodes to become the leader. Several entries and chunks are batched in a single request/response, avoiding multiple round trips.

**ZooKeeper.** In ZAB, the epoch and index are packed into the *zxid* which uniquely identifies a log entry [5]. Followers discover and connect to the leader in Phase 1. We modified Phase 1 to send information about the followers' faulty data. The followers are synchronized with the leader in Phase 2. We modified Phase 2 so that the leader sends the correct data to the followers. The leader waits to hear from a majority during Phase 1 after which it sends a `newEpoch` message; we modified this message to send information about the leader's faulty data. The leader does not proceed to Phase 2 until its data is fixed.

# 5 Evaluation

We evaluate the correctness and performance of CTRL versions of LogCabin and ZooKeeper. We conducted our performance experiments on a three-node cluster on a 1-Gb network; each node is a 40-core Intel Xeon CPU E5-2660 machine with 128 GB memory running Linux 3.13, with a 500-GB SSD and a 1-TB HDD managed by ext4.

## 5.1 Correctness

To verify CTRL's safety and availability guarantees, we built a fault-injection framework that can inject storage faults (targeted corruptions and random block corruptions and errors). The framework can also inject crashes. By injecting crashes at different points in time, the framework simulates lagging nodes. After injecting faults, we issue reads from clients to determine whether the target system remains available and preserves safety.

We first exercise different log-recovery scenarios. Then, we test snapshot recovery, and finally file-system metadata fault recovery.

### 5.1.1 Log Recovery

**Targeted Corruptions.** We initialize the cluster with four log entries, replicated to all three nodes. We exercise all combinations of entry corruptions across the three nodes ($(2^4)^3 = 4096$ combinations). Out of the 4096 cases, a correct recovery is possible in 2401 cases (at least one non-faulty copy of each entry exists). In the remaining 1695 cases, recovery is not possible because one or more entries are corrupted on *all* the nodes. We inject targeted corruptions into two different sets of on-disk structures. In the first set, on a corruption, the original systems invoke the *truncate* action (i.e., they truncate faulty data and continue). In the second set, the original systems invoke the *crash* action (i.e., node crashes

| System | Total Test Cases | Outcomes | | | | | | Total Test Cases | Outcomes | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Original | | | CTRL | | | | Original | | | CTRL | | |
| | | Unavailable | Unsafe | Correct | Unavailable | Unsafe | Correct | | Unavailable | Unsafe | Correct | Unavailable | Unsafe | Correct |
| Log-Cabin | 1000 | 297 | 257 | 446 | 0 | 0 | 1000 | 1000 | 405 | 36 | 559 | 434 | 0 | 566 |
| Zoo-Keeper | 1000 | 417 | 200 | 383 | 0 | 0 | 1000 | 1000 | 329 | 192 | 479 | 502 | 0 | 498 |

(a) Snapshot Recovery      (b) FS Metadata Faults

Table 5: **Snapshot and FS Metadata Faults.** *(a) and (b) show how* CTRL *recovers from snapshot and FS metadata faults, respectively.*

on detection). For example, while ZooKeeper *truncates* when the tail of a transaction is corrupted, it *crashes* the node if the transaction header is corrupted. CTRL always recovers the corrupted data from other replicas.

Table 4(a) shows the results. When recovery is possible, the original systems recover only in 46/2401 cases. In those 46 cases, no node or only one node is corrupted. In the remaining 2355 cases, the original systems are either unsafe (for *truncate*) or unavailable (for *crash*). In contrast, CTRL correctly recovers in all 2401 cases. When a recovery is not possible (all copies corrupted), the original systems are either unsafe or unavailable in all cases. CTRL, by design, correctly remains unavailable since continuing would violate safety.

**Random Block Corruptions and Errors.** We initialize the cluster by replicating a few entries to all nodes. We first choose a random set of nodes. In each such node, we then corrupt a randomly selected file-system block (from the files implementing the log). We repeat this process, producing 5000 test cases. We similarly inject block errors. Since we inject a fault into a block, several entries and their checksums within the block will be faulty.

Table 4(b) shows the results. For *block corruptions*, original LogCabin is unsafe or unavailable in about 30% ($(738 + 793)/5000$) of cases. Similarly, original ZooKeeper is incorrect in about 30% of cases. On a *block error*, original LogCabin and ZooKeeper simply crash the node, leading to unavailability in about 50% of cases. In contrast, CTRL correctly recovers in all cases.

**Faults with Crashed and Lagging Nodes.** In the previous experiments, all entries were committed and present on all nodes. In this experiment, we inject crashes at different points on a random set of nodes while inserting entries. Thus, in the resultant log states, nodes could be lagging, entries could be uncommitted, and have different epochs on different nodes for the same log index. $\langle S_1 : [a^1, \_, \_], S_2 : [b^2, c^3, \_], S_3 : [b^2, \_, \_] \rangle$ is an example state where $S_1$ appends $a$ at index 1 in epoch 1 (shown in superscript) and crashes, $S_2$ appends $b$ at index 1 in epoch 2, replicates to $S_3$, then $S_2, S_3$ crash and recover,
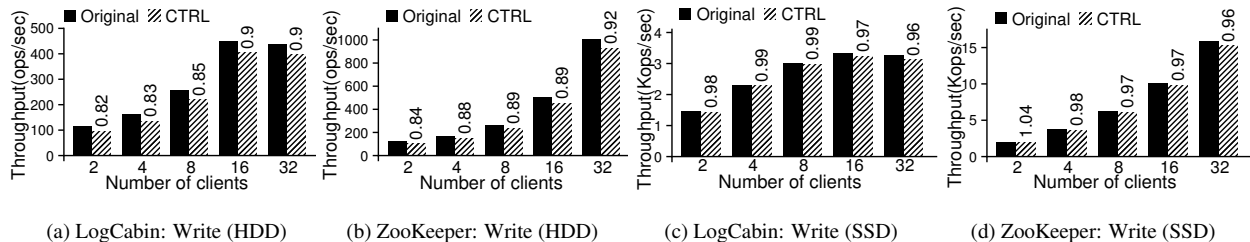
(a) LogCabin: Write (HDD)　(b) ZooKeeper: Write (HDD)　(c) LogCabin: Write (SSD)　(d) ZooKeeper: Write (SSD)

Figure 5: **Common-Case Write Performance.** *(a) and (b) show the write throughput in original and* CTRL *versions of LogCabin and ZooKeeper on an HDD. (c) and (d) show the same for SSD. The number on top of each bar shows the performance of* CTRL *normalized to that of original.*

$S_2$ appends $c$ in epoch 3 and crashes. From each such state, we corrupt different entries, generating 5000 test cases. For example, from the above state, we corrupt $a$ on $S1$ and $b, c$ on $S_2$. If $S_2$ is elected the leader, $S_2$ needs to fix $b$ from $S_3$ (since $b$ is committed), discard $c$ ($c$ is uncommitted and cannot be recovered), and also instruct $S1$ to discard $a$ ($a$ is uncommitted) and replicate correct entry $b$. As shown in Table 4(c), CTRL correctly recovers from all such cases, while the original versions are unsafe or unavailable in many cases.

**Model Checking.** We also model checked CTRL's log recovery since it involves many corner cases, using a python-based model that we developed. We explored over 2.5M log states all of in which CTRL correctly recovered. Also, when key decisions are tweaked, the checker finds a violation immediately: for example, the leader concludes that a faulty entry is uncommitted only after gathering $\lfloor N/2 \rfloor + 1$ *dontHave* responses; if this number is reduced, then the checker finds a safety violation. We have also added the specification of CTRL's log recovery to the TLA+ specification of Raft [23] and confirmed that it correctly recovers from corruptions, while the original specification violates safety.

### 5.1.2　Snapshot Recovery

We trigger the nodes to take a snapshot, crashing them at different points, producing three possible states for each node: $l$, $t$, and $g$, where $l$ is a state where the node has only the log (it has not taken a snapshot), $t$ is a snapshot for which garbage collection has not been performed yet, and $g$ is a snapshot which has been garbage collected. We produce all possible combinations of states across three nodes. On each such state, we randomly pick a set of nodes to inject faults, and corrupt a random combination of snapshots and log entries, generating 1000 test cases. For example, $\langle S_1 : t, S_2 : g, S_3 : l \rangle$ is a base state on which we corrupt snapshot $t$ and a few preceding log entries on $S_1$ and $g$ on $S_2$. In such a state, if $S_1$ becomes the leader, it has to fix its log from $S_3$, then has to locally recover its $t$ snapshot, after which it has to fix $g$ on $S_2$. $S_1$ also needs to install the snapshot on $S_3$. As shown in Table 5(a), CTRL correctly recovers from all such cases. Original LogCabin is incorrect in about half of the cases because it obliviously loads faulty snapshots sometimes and crashes sometimes. Original ZooKeeper crashes the

node if it is unable to locally construct the data from the snapshot and the log, leading to unavailability; unsafety results because a faulty log is truncated in some cases.

### 5.1.3　File-system Metadata Faults

To test how CTRL recovers from file-system metadata faults, we corrupt file-system metadata structures (such as inodes and directory blocks) resulting in unopenable files, missing files, and files with fewer or more bytes. We inject such faults in a randomly chosen file on one or two nodes at a time, creating 1000 test cases. Table 5(b) shows the results. In some cases, the faulty nodes in original versions crash because of a failed deserialization or assertion. However, sometimes original LogCabin and ZooKeeper do not detect the fault and continue operating, violating safety in 36 and 192 cases, respectively. In contrast, CTRL reliably crashes the node on a file-system metadata fault, preserving safety always.

## 5.2　Performance

We now compare the common-case performance of the CTRL versions against the original versions. In both Log-Cabin and ZooKeeper, reads are served from memory and the read paths are not affected by CTRL. Hence, we show only performance of write workloads. The workload runs for 300 seconds, inserting entries each of size 1K. Both systems batch writes to improve throughput. Snapshots are taken periodically during the updates. Numbers reported are the average over five runs.

Figure 5(a) and (b) show the throughput on an HDD for varying number of clients in LogCabin and ZooKeeper, respectively. CLSTORE physically separates the identifier from the entry; this separation induces a seek on disks in the update path. However, the seek cost is amortized when more requests are batched; CTRL has an overhead of 8%-10% for 32 clients on disks. Figure 5(c) and (d) show throughput on an SSD; CTRL adds very minimal overhead on SSDs (4% in the worst case). Note that our workload performs only writes and therefore shows CTRL's overheads in the worst case; for more realistic workloads that predominantly perform reads, the overheads should be even lower.

**Fast Log Recovery.** To show the potential reduction in log-recovery time, we insert 30K log entries (each of size 1K) and corrupt the first entry on one node. In origi-

nal LogCabin, the faulty node detects the corruption but truncates *all* entries; hence, the leader transfers all entries to bring the node up-to-date. CTRL fixes only the faulty entry, reducing recovery time. The faulty node is fixed in 1.24 seconds (32MB transferred) in the original system, while CTRL takes only 1.2 ms (7KB transferred). We see a similar reduction in log-recovery time in ZooKeeper.

# 6   Related Work

Our analysis of how RSM-based systems react to storage faults (§2.3) builds upon several fault-injection studies. Our design of CTRL (§3) builds upon several efforts on tolerating practical faults in distributed systems.

**Storage Faults**. Several studies on storage faults [34, 46, 48, 59, 60] motivated our work. Our previous work [29, 30] discovered fundamental reasons why distributed systems are not resilient to storage faults. However, the study did not uncover any safety or availability violations reported in §2.3; this is because the fault model in our previous study considers injecting only storage faults (precisely, a single storage fault on a single node at a time). In contrast, our fault model in this work considers crashes and network failures in addition to storage faults, exposing previously unknown safety and availability violations in RSM systems.

**Targeted Approaches**. Prior research describes two approaches [15, 17] to tackle storage faults in RSM systems. However, these approaches suffer from unavailability. Furthermore, the *MarkNonVoting* approach [17] can violate safety because important metainfo such as promises can be lost on a storage fault [70]. CTRL avoids such safety violations by storing two copies of metainfo on each node. Approaches that improve the reliability of other specific systems have also been proposed [68, 71].

**Generic Approaches**. Many generic approaches to handling practical faults other than crashes have been proposed. PASC [21] hardens systems to tolerate corruptions by maintaining two copies of the entire state on each node and assumes that both the copies will not be faulty at the same time. This approach does not work well for storage faults; having two copies of on-disk state incurs $2\times$ space overhead. Furthermore, in most cases, PASC crashes the node on a fault, causing unavailability. XFT [42] is designed to tolerate non-crash faults. However, it can tolerate only a total of $\lfloor (N-1)/2 \rfloor$ crash and non-crash faults. Similarly, UpRight [20] has an upper bound on the total faults to remain safe and available.

CTRL differs from the generic approaches through its special focus on storage faults. This focus brings two main advantages. First, CTRL attributes faults at a fine granularity: while the generic approaches consider a node as faulty if any of its data is corrupted, CTRL considers faults at the granularity of individual data items. Second, because of such fine-granular fault treatment,

CTRL can be available as long as a majority of nodes are up and at least one non-faulty copy of a data item exists even though portions of data on *all* nodes could be corrupted. CTRL cannot tolerate arbitrary non-crash faults [40] (e.g., memory errors). However, CTRL can augment the generic approaches: for example, a system can be hardened against memory faults using PASC while making it robust to storage faults using CTRL.

# 7   Conclusions

Recovering from storage faults in distributed systems is surprisingly hard. We introduce *protocol-aware recovery* (PAR), a new approach that exploits protocol-specific knowledge of the underlying distributed system to correctly recover from storage faults. We design CTRL, a protocol-aware recovery approach for RSM systems. We experimentally show that CTRL correctly recovers from a range of storage faults with little performance overhead.

Our work is only a first step in hardening distributed systems to storage faults: while we have successfully applied the PAR approach to RSM systems, other classes of systems (e.g., primary-backup, Dynamo-style quorums) still remain to be analyzed. We believe the PAR approach can be applied to such classes as well. We hope our work will lead to more work on building reliable distributed storage systems that are robust to storage faults.

# Acknowledgments

# References

[1] Crash-Corruption Disentanglement Proof. http://research.cs.wisc.edu/adsl/Publications/par/.

[2] Ittai Abraham, Gregory Chockler, Idit Keidar, and Dahlia Malkhi. Byzantine Disk Paxos: Optimal Resilience with Byzantine Shared Memory. *Distributed Computing*, 18(5):387–408, 2006.

[3] Ramnatthan Alagappan, Aishwarya Ganesan, Yuvraj Patel, Thanumalayan Sankaranarayana Pillai,

Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. Correlated Crash Vulnerabilities. In *Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation (OSDI '16)*, Savannah, GA, November 2016.

[4] Apache. ZooKeeper. `https://zookeeper.apache.org/`.

[5] Apache. ZooKeeper Guarantees, Properties, and Definitions. `https://zookeeper.apache.org/doc/r3.2.2/zookeeperInternals.html#sc_guaranteesPropertiesDefinitions`.

[6] Apache Cassandra. Cassandra Replication. `http://docs.datastax.com/en/cassandra/2.0/cassandra/architecture/architectureDataDistributeReplication_c.html`.

[7] Apache ZooKeeper. Applications and Organizations using ZooKeeper. `https://cwiki.apache.org/confluence/display/ZOOKEEPER/PoweredBy`.

[8] Remzi H. Arpaci-Dusseau and Andrea C. Arpaci-Dusseau. *Operating Systems: Three Easy Pieces*. Arpaci-Dusseau Books, 0.91 edition, May 2015.

[9] Lakshmi N. Bairavasundaram, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau, Garth R. Goodson, and Bianca Schroeder. An Analysis of Data Corruption in the Storage Stack. In *Proceedings of the 6th USENIX Symposium on File and Storage Technologies (FAST '08)*, San Jose, CA, February 2008.

[10] Lakshmi N. Bairavasundaram, Garth R. Goodson, Shankar Pasupathy, and Jiri Schindler. An Analysis of Latent Sector Errors in Disk Drives. In *Proceedings of the 2007 ACM SIGMETRICS Conference on Measurement and Modeling of Computer Systems (SIGMETRICS '07)*, San Diego, CA, June 2007.

[11] Lakshmi N. Bairavasundaram, Meenali Rungta, Nitin Agrawal, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau, and Michael M. Swift. Analyzing the Effects of Disk-Pointer Corruption. In *Proceedings of the International Conference on Dependable Systems and Networks (DSN '08)*, Anchorage, Alaska, June 2008.

[12] Lakshmi Narayanan Bairavasundaram. *Characteristics, Impact, and Tolerance of Partial Disk Failures*. PhD thesis, University of Wisconsin, Madison, 2008.

[13] Mahesh Balakrishnan, Dahlia Malkhi, Vijayan Prabhakaran, Ted Wobber, Michael Wei, and John D. Davis. CORFU: A Shared Log Design for Flash Clusters. In *Proceedings of the 9th Symposium on Networked Systems Design and Implementation (NSDI '12)*, San Jose, CA, April 2012.

[14] Andrew D. Birrell, Roy Levin, Michael D. Schroeder, and Roger M. Needham. Grapevine: An Exercise in Distributed Computing. *Commun. ACM*, 25(4):260–274, April 1982.

[15] William J. Bolosky, Dexter Bradshaw, Randolph B. Haagens, Norbert P. Kusters, and Peng Li. Paxos Replicated State Machines As the Basis of a High-performance Data Store. In *Proceedings of the 8th Symposium on Networked Systems Design and Implementation (NSDI '11)*, Boston, MA, April 2011.

[16] Mike Burrows. The Chubby Lock Service for Loosely-Coupled Distributed Systems. In *Proceedings of the 7th Symposium on Operating Systems Design and Implementation (OSDI '06)*, Seattle, WA, November 2006.

[17] Tushar D Chandra, Robert Griesemer, and Joshua Redstone. Paxos Made Live: An Engineering Perspective. In *Proceedings of the 26th ACM Symposium on Principles of Distributed Computing*, Portland, OR, August 2007.

[18] Vijay Chidambaram, Thanumalayan Sankaranarayana Pillai, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. Optimistic Crash Consistency. In *Proceedings of the 24th ACM Symposium on Operating Systems Principles (SOSP '13)*, Farmington, PA, November 2013.

[19] Vijay Chidambaram, Tushar Sharma, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. Consistency Without Ordering. In *Proceedings of the 10th USENIX Symposium on File and Storage Technologies (FAST '12)*, San Jose, CA, February 2012.

[20] Allen Clement, Manos Kapritsos, Sangmin Lee, Yang Wang, Lorenzo Alvisi, Mike Dahlin, and Taylor Riche. Upright Cluster Services. In *Proceedings of the 22nd ACM Symposium on Operating Systems Principles (SOSP '09)*, Big Sky, Montana, October 2009.

[21] Miguel Correia, Daniel Gómez Ferro, Flavio P. Junqueira, and Marco Serafini. Practical Hardening of Crash-Tolerant Systems. In *2012 USENIX Annual Technical Conference (USENIX ATC 12)*, Boston, MA, June 2012.

[22] Jeff Dean. Building Large-Scale Internet Services. `http://static.googleusercontent.com/media/research.google.com/en//people/jeff/SOCC2010-keynote-slides.pdf`.

[23] Diego Ongaro. Raft TLA+ Specification. `https://github.com/ongardie/raft.tla`.

[24] epaxos. epaxos source code. `https://github.com/efficient/epaxos`.

[25] etcd. etcd. `https://coreos.com/etcd`.

[26] Etcd. Etcd: Production users. `https://coreos.com/etcd/docs/latest/production-users.html`.

[27] Daniel Fryer, Dai Qin, Jack Sun, Kah Wai Lee, Angela Demke Brown, and Ashvin Goel. Checking the Integrity of Transactional Mechanisms. In *Proceedings of the 12th USENIX Symposium on File and Storage Technologies (FAST '14)*, Santa Clara, CA, February 2014.

[28] Daniel Fryer, Kuei Sun, Rahat Mahmood, TingHao Cheng, Shaun Benjamin, Ashvin Goel, and Angela Demke Brown. Recon: Verifying File System Consistency at Runtime. In *Proceedings of the 10th USENIX Symposium on File and Storage Technologies (FAST '12)*, San Jose, CA, February 2012.

[29] Aishwarya Ganesan, Ramnatthan Alagappan, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. Redundancy Does Not Imply Fault Tolerance: Analysis of Distributed Storage Reactions to File-System Faults. *ACM Trans. Storage*, 13(3):20:1–20:33, September 2017.

[30] Aishwarya Ganesan, Ramnatthan Alagappan, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. Redundancy Does Not Imply Fault Tolerance: Analysis of Distributed Storage Reactions to Single Errors and Corruptions. In *Proceedings of the 15th USENIX Conference on File and Storage Technologies (FAST '17)*, Santa Clara, CA, February 2017.

[31] Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung. The Google File System. In *Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP '03)*, Bolton Landing, New York, October 2003.

[32] Matthias Grawinkel, Thorsten Schafer, Andre Brinkmann, Jens Hagemeyer, and Mario Porrmann. Evaluation of Applied Intra-disk Redundancy Schemes to Improve Single Disk Reliability. In *Proceedings of the 19th Annual Meeting of the IEEE International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS)*, Washington, DC, July 2011.

[33] Kevin M Greenan, Darrell DE Long, Ethan L Miller, Thomas Schwarz, and Avani Wildani. Building Flexible, Fault-Tolerant Flash-Based Storage Systems. In *The 5th Workshop on Hot Topics in System Dependability (HotDep '09)*, Lisbon, Portugal, June 2009.

[34] Laura M Grupp, Adrian M Caulfield, Joel Coburn, Steven Swanson, Eitan Yaakobi, Paul H Siegel, and Jack K Wolf. Characterizing Flash Memory: Anomalies, Observations, and Applications. In *Proceedings of the 42nd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO'09)*, New York, New York, December 2009.

[35] James Hamilton. On Designing and Deploying Internet-Scale Services. In *Proceedings of the 21st Annual Large Installation System Administration Conference (LISA '07)*, Dallas, Texas, November 2007.

[36] James Myers. Data Integrity in Solid State Drives. `http://intel.ly/2cF0dTT`.

[37] John Goerzen. Silent Data Corruption Is Real. `http://changelog.complete.org/archives/9769-silent-data-corruption-is-real`.

[38] Jonathan Corbet. Responding to ext4 journal corruption. `https://lwn.net/Articles/284037/`.

[39] Flavio P Junqueira, Benjamin C Reed, and Marco Serafini. Zab: High-Performance Broadcast for Primary-Backup Systems. In *Proceedings of the International Conference on Dependable Systems and Networks (DSN '11)*, Hong Kong, China, June 2011.

[40] Dmitrii Kuvaiskii, Rasha Faqeh, Pramod Bhatotia, Pascal Felber, and Christof Fetzer. HAFT: Hardware-assisted Fault Tolerance. In *Proceedings of the EuroSys Conference (EuroSys '16)*, London, United Kingdom, April 2016.

[41] Leslie Lamport. Paxos Made Simple. *ACM Sigact News*, 32(4):18–25, 2001.

[42] Shengyun Liu, Paolo Viotti, Christian Cachin, Vivien Quéma, and Marko Vukolic. XFT: Practical Fault Tolerance Beyond Crashes. In *Proceedings*

*of the 12th USENIX Conference on Operating Systems Design and Implementation (OSDI '16)*, Savannah, GA, November 2016.

[43] LogCabin. LogCabin. `https://github.com/logcabin/logcabin`.

[44] Jacob R Lorch, Atul Adya, William J Bolosky, Ronnie Chaiken, John R Douceur, and Jon Howell. The SMART Way to Migrate Replicated Stateful Services. In *Proceedings of the EuroSys Conference (EuroSys '06)*, Leuven, Belgium, April 2006.

[45] Parisa Jalili Marandi, Christos Gkantsidis, Flavio Junqueira, and Dushyanth Narayanan. Filo: Consolidated Consensus As a Cloud Service. In *2016 USENIX Annual Technical Conference (USENIX ATC 16)*, Denver, CO, June 2016.

[46] Justin Meza, Qiang Wu, Sanjev Kumar, and Onur Mutlu. A Large-Scale Study of Flash Memory Failures in the Field. In *Proceedings of the 2015 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS '15)*, Portland, Oregon, June 2015.

[47] MongoDB. MongoDB Replication. `https://docs.mongodb.org/manual/replication/`.

[48] Iyswarya Narayanan, Di Wang, Myeongjae Jeon, Bikash Sharma, Laura Caulfield, Anand Sivasubramaniam, Ben Cutler, Jie Liu, Badriddine Khessib, and Kushagra Vaid. SSD Failures in Datacenters: What? When? And Why? In *Proceedings of the 9th ACM International on Systems and Storage Conference (SYSTOR '16)*, Haifa, Israel, June 2016.

[49] Diego Ongaro. *Consensus: Bridging Theory and Practice*. PhD thesis, Stanford University, 2014.

[50] Diego Ongaro and John Ousterhout. In Search of an Understandable Consensus Algorithm. In *2014 USENIX Annual Technical Conference (USENIX ATC 14)*, Philadelphia, PA, June 2014.

[51] Bernd Panzer-Steindel. Data Integrity. *CERN/IT*, 2007.

[52] Thanumalayan Sankaranarayana Pillai, Ramnatthan Alagappan, Lanyue Lu, Vijay Chidambaram, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. Application Crash Consistency and Performance with CCFS. In *Proceedings of the 15th USENIX Conference on File and Storage Technologies (FAST '17)*, Santa Clara, CA, February 2017.

[53] Thanumalayan Sankaranarayana Pillai, Vijay Chidambaram, Ramnatthan Alagappan, Samer Al-Kiswany, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. All File Systems Are Not Created Equal: On the Complexity of Crafting Crash-consistent Applications. In *Proceedings of the 11th Symposium on Operating Systems Design and Implementation (OSDI '14)*, Broomfield, CO, October 2014.

[54] Vijayan Prabhakaran, Lakshmi N. Bairavasundaram, Nitin Agrawal, Haryadi S. Gunawi, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. IRON File Systems. In *Proceedings of the 20th ACM Symposium on Operating Systems Principles (SOSP '05)*, Brighton, UK, October 2005.

[55] Redis. Redis Replication. `http://redis.io/topics/replication`.

[56] Robert Ricci, Eric Eide, and CloudLab Team. Introducing CloudLab: Scientific infrastructure for advancing cloud architectures and applications. *USENIX ;login:*, 39(6), 2014.

[57] Robert Harris. Data corruption is worse than you know. `http://www.zdnet.com/article/data-corruption-is-worse-than-you-know/`.

[58] Fred B. Schneider. Implementing Fault-tolerant Services Using the State Machine Approach: A Tutorial. *ACM Comput. Surv.*, 22(4):299–319, December 1990.

[59] Bianca Schroeder, Sotirios Damouras, and Phillipa Gill. Understanding Latent Sector Errors and How to Protect Against Them. In *Proceedings of the 8th USENIX Symposium on File and Storage Technologies (FAST '10)*, San Jose, CA, February 2010.

[60] Bianca Schroeder, Raghav Lagisetty, and Arif Merchant. Flash Reliability in Production: The Expected and the Unexpected. In *Proceedings of the 14th USENIX Conference on File and Storage Technologies (FAST '16)*, Santa Clara, CA, February 2016.

[61] Michael D. Schroeder, Andrew D. Birrell, and Roger M. Needham. Experience with Grapevine: The Growth of a Distributed System. *ACM Trans. Comput. Syst.*, 2(1):3–23, February 1984.

[62] Thomas Schwarz, Ahmed Amer, Thomas Kroeger, Ethan L. Miller, Darrell D. E. Long, and Jehan-Franois Pris. RESAR: Reliable Storage at Exabyte

Scale. In *Proceedings of the 24th Annual Meeting of the IEEE International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS)*, London, United Kingdom, September 2016.

[63] Romain Slootmaekers and Nicolas Trangez. Arakoon: A Distributed Consistent Key-Value Store. In *SIGPLAN OCaml Users and Developers Workshop*, volume 62, 2012.

[64] Stackoverflow. Can ext4 detect corrupted file contents? `http://stackoverflow.com/questions/31345097/can-ext4-detect-corrupted-file-contents`.

[65] Stackoverflow. ZooKeeper Clear State. `http://stackoverflow.com/questions/17038957/org-apache-hadoop-hbase-pleaseholdexception-master-is-initializing`.

[66] Michael M. Swift, Brian N. Bershad, and Henry M. Levy. Improving the Reliability of Commodity Operating Systems. In *Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP '03)*, Bolton Landing, New York, October 2003.

[67] D. B. Terry, M. M. Theimer, Karin Petersen, A. J. Demers, M. J. Spreitzer, and C. H. Hauser. Managing Update Conflicts in Bayou, a Weakly Connected Replicated Storage System. In *Proceedings of the 15th ACM Symposium on Operating Systems Principles (SOSP '95)*, Copper Mountain Resort, CO, December 1995.

[68] Thanh Do, Tyler Harter, Yingchao Liu, Haryadi S. Gunawi, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau. HARDFS: Hardening HDFS with Selective and Lightweight Versioning. In *Proceedings of the 11th Conference on File and Storage Technologies (FAST '13)*, San Jose, California, February 2013.

[69] Theodore Ts'o. What to do when the journal checksum is incorrect. `https://lwn.net/Articles/284038/`.

[70] Robbert Van Renesse, Nicolas Schiper, and Fred B Schneider. Vive La Différence: Paxos vs. Viewstamped Replication vs. Zab. *IEEE Transactions on Dependable and Secure Computing*, 12(4):472–484, 2015.

[71] Yang Wang, Manos Kapritsos, Zuocheng Ren, Prince Mahajan, Jeevitha Kirubanandam, Lorenzo Alvisi, and Mike Dahlin. Robustness in the Salus Scalable Block Store. In *Proceedings of the 10th Symposium on Networked Systems Design and Implementation (NSDI '13)*, Lombard, IL, April 2013.

[72] Yupu Zhang, Abhishek Rajimwale, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. End-to-end Data Integrity for File Systems: A ZFS Case Study. In *Proceedings of the 8th USENIX Symposium on File and Storage Technologies (FAST '10)*, San Jose, CA, February 2010.

[73] ZooKeeper Jira Issues. Unable to load database on disk when restarting after node freeze. `https://issues.apache.org/jira/browse/ZOOKEEPER-1546`.