

Towards Robust File System Checkers

Om Rameshwar Gatla, Muhammad Hameed, and Mai Zheng, New Mexico State University; Viacheslav Dubeyko, Adam Manzanares, Filip Blagojevic, Cyril Guyot, and Robert Mateescu, Western Digital Research

https://www.usenix.org/conference/fast18/presentation/gatla

This paper is included in the Proceedings of the 16th USENIX Conference on File and Storage Technologies. February 12–15, 2018 • Oakland, CA, USA

ISBN 978-1-931971-42-3

Open access to the Proceedings of the 16th USENIX Conference on File and Storage Technologies is sponsored by USENIX.

Towards Robust File System Checkers

Om Rameshwar Gatla[†], Muhammad Hameed[†], Mai Zheng[†], Viacheslav Dubeyko[‡], Adam Manzanares[‡], Filip Blagojevic[‡], Cyril Guyot[‡], Robert Mateescu[‡]

†New Mexico State University

[‡]Western Digital Research

Abstract

File systems may become corrupted for many reasons despite various protection techniques. Therefore, most file systems come with a checker to recover the file system to a consistent state. However, existing checkers are commonly assumed to be able to complete the repair without interruption, which may not be true in practice.

In this work, we demonstrate via fault injection experiments that checkers of widely used file systems may leave the file system in an uncorrectable state if the repair procedure is interrupted unexpectedly. To address the problem, we first fix the ordering issue in the undo logging of e2fsck, and then build a general logging library (i.e., rfsck-lib) for strengthening checkers. To demonstrate the practicality, we integrate rfsck-lib with existing checkers and create two new checkers: (1) rfsck-ext, a robust checker for Ext-family file systems, and (2) rfsck-xfs, a robust checker for XFS file system, both of which require only tens of lines of modification to the original versions. Both rfsck-ext and rfsck-xfs are resilient to faults in our experiments. Also, both checkers incur reasonable performance overhead (i.e., up to 12%) comparing to the original unreliable versions. Moreover, rfsck-ext outperforms the patched e2fsck by up to nine times while achieving the same level of robustness.

Introduction

Achieving data integrity is critical for computer systems ranging from a single desktop to large-scale distributed storage clusters [21]. In order to make sense of the ever increasing amount of data stored, it is common to use local (e.g., Ext4 [4], XFS [70], F2FS [49]) and multinode file systems (e.g., HDFS [66], Ceph [74], Lustre [9]) to organize the data on top of storage devices. Although file systems are designed to maintain the data integrity [36, 38, 45, 60, 72, 75], situations arise when the file system metadata needs to be checked for integrity. Such situations may be caused by power outages, server crashes, latent sector errors, software bugs, etc [19, 20, 31, 51, 54].

File system checkers, such as e2fsck for Ext-family file systems [3], serve as the last line of defense to recover a corrupted file system back to a healthy state [54]. They contain intimate knowledge of file system metadata structures, and are commonly assumed to be able to complete the repair without interruption.

Unfortunately, the same issues that lead to file system inconsistencies (e.g., power outages or crashes), can also occur during file system repair. One real-world example happened at the High Performance Computing Center in Texas [17]. In this accident, multiple Lustre file systems suffered severe data loss after power outages: the first outage triggered the Lustre checker (lfsck [6]) after the cluster was restarted, while another outage interrupted 1fsck and led to the downtime and data loss. Because Lustre is built on top of a variant of Ext4 (ldiskfs [9]), and lfsck relies on e2fsck to fix local inconsistencies on each node, the checking and repairing is complicated (e.g., requiring several days [17]). As of today, it is still unclear which step of lfsck/e2fsck caused the uncorrectable corruptions. With the trend of increasing the storage capacity and scaling to more and more nodes, checking and repairing file systems will likely become more time-consuming and thus more vulnerable to faults. Such accidents and observation motivate us to remove the assumption that file system checkers can always finish normally without interruption.

Previous research has demonstrated that file system checkers themselves are error-prone [27, 42]. File system specific approaches have also been developed that use higher level languages to elegantly describe file system repair tasks [42]. In addition, efforts have also been made to speed up the repair procedure, which leads to a smaller window of potential data loss due to an interruption [54]. Although these efforts improve file system checkers, they do not address the fundamental issue of improving the resilience of checkers in the face of unexpected interruptions.

In this work, we first demonstrate that the checkers of widely used file systems (i.e., e2fsck [3] and xfs_repair [14]) may leave the file system in an uncorrectable state if the repair procedure is unexpectedly interrupted. We collect corrupted file system images from file system developers and additionally generate test images to trigger the repair procedure. Moreover, we develop rfsck-test, an automatic fault injection tool, to systematically inject faults during the repair, and thus manifest the vulnerabilities.

To address the problem exposed in our study, we analyze the undo logging feature of e2fsck in depth, and identify an ordering issue which jeopardizes its effectiveness. We fix the issue and create a patched version called e2fsck-patch which is truly resilient to faults.

However, we find that e2fsck-patch is inherently suboptimal as it requires extensive sync operations. To address the limitation, and to improve the checkers of other file systems, we design and implement rfsck-lib, a general logging library with a simple interface. Based on the similarities among checkers, rfsck-lib decouples the logging from the repairing, and provides an interface to log the repairing writes in fine granularity.

To demonstrate the practicality, we integrate rfsck-lib with existing checkers and create two new checkers: (1) rfsck-ext, a robust checker for Extfamily file systems, which adds 50 lines of code (LoC) to e2fsck; and (2) rfsck-xfs, a robust checker for XFS file system, which adds 15 LoC to xfs_repair.¹ Both rfsck-ext and rfsck-xfs are resilient to faults in our experiments. Also, both checkers incur reasonable performance overhead (i.e., up to 12%) compared to the original unreliable versions. Moreover, rfsck-ext outperforms e2fsck-patch by up to nine times while achieving the same level of fault resilience.

The rest of the paper is organized as follows. First, we introduce the background of file system checkers (§2). Next, we describe rfsck-test and study e2fsck and xfs_repair (§3). We analyze the ordering issue of the undo logging of e2fsck in §4. Then, we introduce rfsck-lib and integrate it with existing checkers (§5). We evaluate rfsck-ext and rfsck-xfs in §6, and discuss several issues in §7. Finally, we discuss related work ($\S 8$) and conclude ($\S 9$).

2 Background

Most file systems employ checkers to check and repair inconsistencies. The checkers are usually file system specific, and they examine different consistency rules depending on the metadata structures. We use two representative checkers as concrete examples to illustrate the complexity as well as the potential vulnerabilities of checkers in this section.

2.1 Workflow of e2fsck

e2fsck is the checker of the widely used Ext-family file systems. It first replays the journal (in case of Ext3 and Ext4) and then restarts itself. Next, e2fsck runs the following five passes in order:

Pass-1: Scan the file system and check inodes. e2fsck scans the entire volume and stores information of all inodes into a set of bitmaps. In addition, it performs four sub-passes to generate a list of duplicate blocks and their owners, check the integrity of extent trees, etc.

Pass-2: Check directory structure. Based on the bitmap information, e2fsck iterates through all directory inodes and checks a set of rules for each directory.

Pass-3: Check directory connectivity. e2fsck first checks if a root directory is available; if not, a new root directory is created and is marked "done". Then it traverses the directory tree, checks the reachability of each directory inode, breaks directory loops, etc.

Pass-4: Check reference counts. e2fsck iterates over all inodes to validate the inode link counts. Also, it checks the connectivity of the extended attribute blocks and reconnects them if necessary.

Pass-5: Recalculate checksums and flush updates. Finally, e2fsck checks the repaired in-memory data structures against on-disk data structures and flushes necessary updates to the disk.

2.2 Workflow of xfs_repair

xfs_repair is the checker of the popular XFS file system.² Similar to e2fsck, xfs_repair fixes inconsistencies in seven passes (or phases), including: Pass-1, superblock verification; Pass-2, replay logs, validate maps and the root inode; Pass-3, check inodes in each allocation group; Pass-4, check duplicate block allocations; Pass-5, rebuild the allocation group structure and superblock; Pass-6, check inode connectivity; Pass-7, verify and correct link counts.

Unlike e2fsck which is single-threaded, xfs_repair employs multi-threading in passes 2, 3, 6 and 7 to improve the performance. Nevertheless, we can see that both checkers are complicated and may be vulnerable to faults. For example, later passes may depend on previous passes, and there is no atomicity guarantee for related updates. We describe our method for systematically exposing the vulnerabilities in §3.

¹The prototypes of rfsck-test, e2fsck-patch, rfsck-lib, rfsck-ext, and rfsck-xfs are publicly available [10].

²There is another utility called xfs_check [14], which checks the consistency without repairing; we do not evaluate it in this work as it is impossible for the read-only utility to introduce additional corruption.

2.3 The Logging Support of Checkers

Some file system developers have envisioned the potential need of reverting the changes done to the file system. For example, the "undo io manager" has been added to the utilities of Ext-family file systems since 2007 [3, 15]. It can save the content of the location being overwritten to an undo log before committing the overwrite.

However, due to the degraded performance as well as the log format issues [2, 16], the undo feature has not been integrated into e2fsck until recently. Starting from v1.42.12, e2fsck includes a "-z" option to allow the user to specify the path of the log file and enable logging [3]. When enabled, e2fsck maintains an undo log during the checking and repairing, and writes an undo block to the log before updating any block of the image. If e2fsck fails unexpectedly, the undo log can be replayed via e2undo [3] to revert the undesired changes.

Given the undo logging, one might expect that an interrupted e2fsck will not cause any issue. As we will see in the next section, however, this is not true.

Are the Existing Checkers Resilient to

In this section, we first describe our method for analyzing the fault resilience of file system checkers (§3.1 - $\S 3.3$), and then present our findings on e2fsck ($\S 3.4$) and $xfs_repair(\S 3.5)$.

Generating Corrupted Test Images

File system checkers are designed to repair corrupted file systems, so the first step of testing checkers is to generate a set of corrupted file system images to trigger the target checker. We call this set of images as test images.

To generate test images, we use two methods. First, some file system developers may provide test images to perform regression testing of their checkers, which usually cover the most representative corruption scenarios as envisioned by the developers [3]. We collect such default test images to trigger the target checker if they are available. Additionally, we create test images by ourselves using the debug tools provided by the file system developers (e.g., debugfs [3] and xfs_db [14]). These tools allow "trashing" specific metadata structures with random bits, which may cover corruption scenarios beyond the default test images.

In both cases, the test images are generated as regular files instead of real physical disks, which makes the testing more efficient.

3.2 Interrupting Checkers

Generating corrupted test images solves only one part of the problem. Another challenge in evaluating the fault resilience is how to interrupt checkers in a systematic and controllable way. To this end, we emulate the effect of faults using software.

To make the emulation precise and reasonable, we follow the "clean power fault" model [80], which assumes that there is a minimal atomic unit of write operations (e.g., 512B or 4KB). Under this model, the size of data written to the on-disk file system is always an integer multiple of the minimal atomic block. A fault can occur at any point during the repair procedure of the checker; once a fault happens, all atomic blocks committed before the fault are durable without corruption, and all blocks after the fault have no effect on the media.

Apparently, this is an idealized model under power outages or system crashes. More severe damage (e.g., reordering or corruption of committed blocks) may happen in practice [61, 73, 77, 81, 82]. However, such clear model can serve as a conservative lower bound of the failure impact. In other words, file system checkers must be able to handle this fault model gracefully before addressing more aggressive fault models.

Based on the fault model, we build a fault injection tool called rfsck-test using a customized driver [8], which has two modes of operation as follows:

Basic mode: This is used for testing a checker without logging support. In this mode, the target checker writes to the test image and generates I/O commands through the customized driver. rfsck-test records the I/O commands generated during the execution of the checker in a command history file, and replays a prefix of the command history (i.e., partial commands) to a copy of the initial test image, which effectively generates the effect of an interrupted checker on the test image. For each command history, we exhaustively replay all possible prefixes, and thus generate a set of interrupted images which correspond to injecting faults at different points during the execution of the checker.

Advanced mode: This is used for testing a checker with logging support. In this mode, the target checker writes to the test image as well as its log file. rfsck-test records the commands sent to both the image and the log in the command history. During the replay, rfsck-test selects a prefix of the command history, and replays the partial commands either to a copy of the initial test image or to a copy of the initial log, depending on the original destination of the commands. In this way, rfsck-test generates the effect of an interrupted checker on both the test image and the log. Moreover, rfsck-test replays the log to the test image, which is necessary for the logging to take effect.

3.3 Summary of Testing Framework

Putting it all together, we summarize our framework for testing the fault resilience of checkers with and without

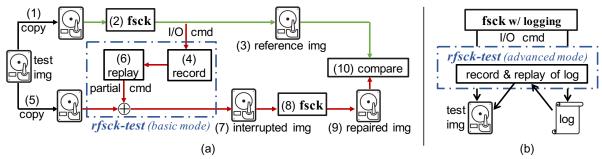


Figure 1: (a) Testing the fault resilience of a file system checker (fsck) without logging support. There are ten steps: (1) make a copy of the test image which contains a corrupted file system; (2) run fsck on the test image copy; (3) store the image generated in step 2 as the reference image; (4) record the I/O commands generated during the fsck; (5) make another copy of the test image; (6) replay partial commands to emulate the effect of an interrupted fsck; (7) store the image generated in step 6 as the interrupted image; (8) run fsck on the interrupted image; (9) store the image generated in step 8 as the repaired image; (10) compare the repaired image with the reference image to identify mismatches. (b) Testing fsck with logging support. The workflow is similar except that rfsck-test interrupts the I/O commands sent to both the test image and the log, and the log is replayed between steps 7 and 8.

logging support as follows:

Testing checkers without logging support: As shown in Figure 1a, there are ten steps: (1) we make a copy of the test image which contains a corrupted file system; (2) the target checker (i.e., fsck) is executed to check and repair the original corruption on the copy of the test image; (3) after fsck finishes normally in the previous step, the resulting image is stored as the reference image;³ (4) during the checking and repairing of fsck, the fault injection tool rfsck-test operates in the basic mode, which records the I/O commands generated by fsck in a command history file; (5) we make another copy of the original test image; (6) rfsck-test replays partial commands recorded in step 4 to the new copy of the test image, which emulates the effect of an interrupted fsck; (7) the image generated in step 6 is stored as the interrupted image; (8) fsck is executed again on the interrupted image to fix any repairable issues; (9) the image generated in step 8 is stored as the repaired image; (10) finally, we compare the file system on the repaired image with that on the reference image to identify any mismatches.

The comparison in step 10 is first performed via the diff command. If a mismatch is reported, we further verify it manually. Note that in step 8 we have run fsck without interruption, so a mismatch implies that there is some corruption which cannot be recovered by fsck.

Testing checkers with logging support: The workflow of testing a checker with logging support is similar. As shown in Figure 1b, rfsck-test operates in the advanced mode, which records the I/O commands sent to both the test image and the log and emulates the effect of interruption on both places. Also, between steps 7 and 8, the (interrupted) log is replayed to the test image to make the logging take effect. The other steps are the same.

Case Study I: e2fsck

In this section, we apply the testing framework to study e2fsck. As discussed in §2.3, e2fsck has recently added the undo logging support. For clarity, we name the original version without undo logging as e2fsck, and the version with undo logging as e2fsck-undo.

To trigger the checker, we collect 175 Ext4 test images from e2fsprogs v1.43.1 [3] as inputs. The sizes of these images range from 8MB to 128MB, and the file system block size is 1KB. To emulate faults on storage systems with different atomic units, we inject faults at two granularities: 512B and 4KB. In other words, we interrupt e2fsck/ e2fsck-undo after every 512B or 4KB of an I/O transfer command. Since the file system block is 1KB, we do not break file system blocks when injecting faults at the 4KB granularity.

First, we study e2fsck using the method in Figure 1a. As described in §3.3, for each fault injected (i.e., each interruption) we run e2fsck again and generate one repaired image. Because the repair procedure usually requires updating multiple file system blocks, it can often be interrupted at multiple points depending on the fault injection granularity. Therefore, we usually generate multiple repaired images from one test image.

For example, to fix the test image "f_dup" (block claimed by two inodes), e2fsck needs to update 16KB in total. At the fault injection granularity of 512B, we generate 32 interrupted images (and consequently 32 repaired images). The last fault is injected after all 16KB blocks, which leads to a repaired image equivalent to the

³It is possible that a checker may not be able to fully repair a corrupted file system even without interruption [27, 42]. So we simply use the result of an uninterrupted repair as a criterion in this work.

Fault injection	# of Ext4	# of repaired	# of images	reporting corruption
granularity	test images	images generated	test images	repaired images
512 B	175	25,062	34	240
4 KB	175	3,192	17	37

Table 1: Counts of images in testing e2fsck at two fault injection granularities. This table shows the number of repaired images (3rd column) generated from the 175 Ext4 test images when injecting faults at 512B/4KB granularities; the last two columns show the number of test images and repaired images reporting corruption respectively.

Corruption	test images		repaired images	
Type	512 B	4 KB	512 B	4 KB
cannot mount	20	1	41	3
data corruption	9	5	107	10
misplacement	9	11	82	23
others	1	1	10	1

Table 2: Classification of corruption. This table shows the number of test images and repaired images reporting different corruptions at two fault injection granularities.

reference image without interruption. Similarly, at the 4KB granularity, we generate 4 repaired images.

For every test image, we generate a number of repaired images and compare each of them with the corresponding reference image. If the comparison reports a mismatch, it implies that the repaired image contains uncorrectable corruption. We count the number of repaired images reporting such corruption. Moreover, if at least one repaired image contains uncorrectable corruption, we mark the test image as reporting corruption, too.

Table 1 summarizes the counts of images in testing e2fsck at the two fault injection granularities. The total number of repaired images generated from the 175 Ext4 test images is shown in the third column. We can see that at the 512B granularity there are more repaired images (25,062) because the repairing procedure is interrupted more frequently, while at the 4KB granularity only 3,192 repaired images are generated. Also, more test images report corruption at the 512B granularity (34 > 17). This is because the repair commands are broken into smaller pieces, and thus it is more challenging to maintain consistency when interrupted.

Table 2 further classifies the corruption into four types and shows the number of test images and repaired images reporting each type. Among the four types, data corruption (i.e., a file's content is corrupted) and misplacement (i.e., a file is either in the "lost+found" folder or completely missing) are the common ones. The most severe corruption is cannot mount (i.e., the whole file system volume becomes not mountable). Such corruption has been observed at both fault injection granularities. In other words, interrupting e2fsck may lead to an unmountable image, even when a fault cannot break the

Fault injection	# of images reporting corruption	
granularity	e2fsck	e2fsck-undo
512 B	34	34
4 KB	17	15

Table 3: Comparison of e2fsck and e2fsck-undo. This table compares the number of test images reporting corruption under e2fsck and e2fsck-undo.

superblock because the 4KB fault granularity is larger than the 1KB superblock.

Next, to see if the undo logging can avoid the corruption, we use the method in Figure 1b to study e2fsck-undo. We focus on the test images which report corruption when testing e2fsck (i.e., the 34 and 17 test images in Table 1).

Table 3 compares the number of test images reporting corruption under e2fsck and e2fsck-undo. Surprisingly, we observe a similar amount of corruption. For example, all 34 images which report corruption when testing e2fsck at the 512B granularity still report corruption under e2fsck-undo. We defer the analysis of the root cause to §4.

Case Study II: xfs_repair 3.5

We have also applied the testing framework to study xfs_repair. Since xfs_repair does not support logging, only the method in Figure 1a is used.

To generate test images, we create 20 clean XFS images first. Each image is 100MB, and the file system block size is 1KB (same as the Ext4 test images). We use the blocktrash command of xfs_db [14] to flip 2 random bits on the metadata area of each image. In this way, we generate 20 corrupted XFS test images in total.

Table 4 summarizes the total number of repaired images generated from the XFS test images at two fault injection granularities. We use 3 test images to inject faults at the 512B granularity, and 17 images for the 4KB granularity. Similar to the Ext4 case, the smaller granularity (i.e., 512B) leads to more repaired images (i.e., 3 test images lead to 1,127 repaired images). The table also shows the number of test images and repaired images reporting corruption. We can see that there are uncorrectable corruptions under both granularities, same as the Ext4 case.

Fault injection	# of XFS	# of repaired	# of images	reporting corruption
injection	test images	images generated	test images	repaired images
512 B	3	1,127	2	443
4 KB	17	1,409	12	737

Table 4: Counts of images in testing xfs_repair at two fault injection granularities. This table shows the number of repaired images (3rd column) generated from the XFS test images when injecting faults at 512B/4KB granularities; the last two columns show the number of test images and repaired images reporting corruption respectively.

```
/*open undo log*/
2.
    undo_open(...){
        open(...); /*no 0_SYNC*/
3.
4.
5.
    /*fix 1st inconsistency*/
6.
7.
    undo_write_blk64(...){
         /*write to undo log asynchronously*/
         undo_write_tdb(...){
9.
10.
             pwrite(...); /*no fsync()*/
11.
12.
         /*write to fs image asynchronously*/
13.
14.
         io_channel_write_blk64(...){...}
15. }
16. /*fix 2<sup>nd</sup>, 3<sup>rd</sup>, ... inconsistencies*/
17. ...
19. /*sync buffered writes to fs image*/
20. ext2fs_flush(...){...}
21. /*close undo log*/
22. undo_close(...){...}
```

Figure 2: Workflow of the undo logging in e2fsckundo. The writes to the log (line 9 -12) and the writes to the file system image (line 14) are asynchronous, and there is no ordering guarantee between the writes.

Why Does the Existing Undo Logging Not Work?

The study in §3 shows that even the checkers of some most popular file systems are not resilient to faults. This is consistent with other studies on the catastrophic failures of real-world systems [41, 44], which find that the recovery procedures themselves are often imperfect, and sometimes "the cure is worse than the disease" [44].

One way to handle the faults and provide crash consistency is write-ahead logging (WAL) [58], which has been widely used in databases [12] and journaling file systems [72] for transactional recovery. While it is perhaps not surprising that file system checkers without crash consistency support (e.g., e2fsck and xfs_repair) may introduce additional corruptions upon interruption, it is counterintuitive that e2fsck-undo, which has the undo logging support, still cannot prevent cascading damage.

To understand the root cause, we analyze the source code of e2fsck-undo as well as the runtime traces (e.g., system calls and I/O commands), and have found that there is no ordering guarantee between the writes to the undo log and the writes to the image being fixed, which essentially invalidates the WAL mechanism.

To better illustrate the issue, Figure 2 shows a simplified workflow of the undo logging in e2fsck-undo. At the beginning of checking (line 2-4), the undo log file is opened without the O_SYNC flag. To fix an inconsistency, e2fsck-undo first gets the original content of the block being repaired (not shown) and then writes it as an undo block to the log asynchronously (line 9-12). After the write to the log, it updates the file system image asynchronously (line 14). The same pattern (i.e., locate the block that needs to be repaired, copy the old content to the log, and update the file system image) is repeated for fixing every inconsistency. At the end, e2fsck-undo flushes all buffered writes of the image to the persistent storage (line 20) and closes the undo log (line 22).

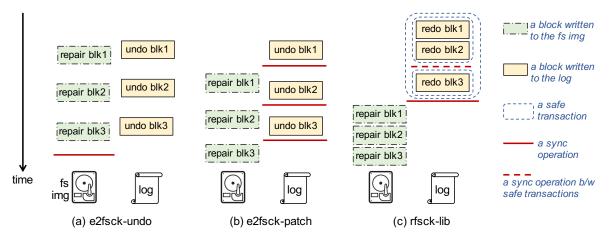
While the extensive asynchronous writes (line 6-17) is good for performance, it is problematic from the WAL's perspective. All asynchronous writes are buffered in memory. Since the dirty pages may be flushed by kernel threads due to memory pressure or timer expiry (e.g., dirty_writeback_centisecs), or by the internal flushing routine of the host file system, there is no strict ordering guarantee among the buffered writes. In other words, for every single fix, the writes to the log and the writes to the file system image may reach the persistent storage in an arbitrary order. Consequently, when e2fsck-undo is interrupted, the file system image may have been modified without the appropriate undo blocks recorded. Because the WAL mechanism works only if a log block reaches the persistent storage before the updated data block it describes, the lack of ordering guarantee between the writes to the log and the writes to the image invalidates the WAL mechanism. As a result, the existing undo logging does not work as expected.

5 **Robust File System Checkers**

In this section, we describe our method to address the problem exposed in §3 and §4.

First, we fix the ordering issue of e2fsck-undo by enforcing necessary sync operations. For clarity, we name the version with our patch as e2fsck-patch.

Next, we observe that although e2fsck-patch may provide the desired robustness, it inherently requires extensive synchronized I/O, which may hurt the perfor-



Comparison of different logging schemes. This figure compares different logging schemes using a Figure 3: sequence of blocks written to the file system image (i.e., "fs img") and the log: (a) e2fsck-undo is the logging scheme of e2fsck, which does not have the necessary ordering guarantee between the writes to the log and the writes to the file system image; (b) e2fsck-patch guarantees the correct ordering between each undo block (e.g., "undo blk1") and the corresponding repair block (e.g., "repair blk1") by enforcing a sync operation (i.e., the red line) after each write of an undo block; (c) rfsck-lib uses redo logging to eliminate the frequent sync required in e2fsck-patch, and only syncs after a safe transaction which includes a set of blocks constituting a consistent update.

mance severely. To address the limitation, and to provide a generic solution to the checkers of other file systems, we design and implement rfsck-lib, a general logging library with a simple interface. Different from e2fsck-patch which interleaves the writes to the log (i.e., log writes) and the writes to the image being repaired (i.e., repair writes), rfsck-lib makes use of the similarities among checkers to decouple the logging from the repairing of the file system, and provides finegrained control of logging.

To demonstrate the practicality, we use rfsck-lib to strengthen existing checkers and create two new checkers: (1) rfsck-ext, a robust checker for Ext-series file systems, and (2) rfsck-xfs, a robust checker for XFS file system, both of which require only a few lines of modification to the original versions.

5.1 Goals

While there are many desired objectives, rfsck-lib is designed to meet the three key goals as follows:

Robustness: Unlike existing checkers which may introduce uncorrectable corruptions when interrupted, we expect checkers integrated with rfsck-lib to be resilient to faults. We believe such robustness should be of prime concern for file system practitioners besides the heavily studied performance issue [54].

Performance: Guaranteeing the robustness may come at the cost of performance, because it almost inevitably requires additional operations. However, the performance overhead should be reduced to minimum, which is particularly important for production environments.

Compatibility: We expect rfsck-lib to be compatible to existing file systems and checkers. For example, no change to the existing on-disk layouts or repair rules is needed. While such compatibility may sacrifice some flexibility and optimization opportunities, it directly enables improving the robustness of many widely used systems in practice.

e2fsck-patch: Fixing the Ordering Issue in e2fsck-undo

As discussed in §4, e2fsck-undo does not guarantee the necessary ordering between log writes and repair writes. Figure 3a illustrates the scenario using a sequence of writes. In this example, three blocks are written to the file system image (i.e., "fs img") to repair inconsistencies (i.e., "repair blk1" to "repair blk3"). Meanwhile, three blocks are written to the undo log (i.e., "undo blk1" to "undo blk3") to save the original content of the blocks being overwritten, for the purpose of undoing changes in case the repair fails. Because all blocks are written asynchronously, the repair blocks may reach the persistent storage before the corresponding undo blocks, which essentially invalidate the undo logging scheme. Although there is a sync operation at the end to the file system image (i.e., the red solid line), it cannot prevent the previous buffered blocks from reaching the persistent storage out of the desired order.

A naive way to solve the issue is to use a synchronous write for each block. However, this is overkill. As long as an undo block (e.g., "undo blk1") becomes persistent, it is unnecessary for the corresponding repair block (e.g.,

"repair blk1") to be written synchronously. Therefore, we only enforce synchronized I/O for the undo log file.

Specifically, we add the O_SYNC flag when opening the undo log file, which is equivalent to adding an fsync call after each write to the log [7]. As shown in Figure 3b, the simple patch guarantees that a repair block is always written after the corresponding undo block becomes persistent. On the other hand, all repair blocks are still written asynchronously. In this way, e2fsck-patch fixes e2fsck-undo with minimum modification.

rfsck-lib: A General Library for **Strengthening File System Checkers**

While the logging scheme of e2fsck-patch may improve the fault resilience, it has two limitations. First, the log writes and the repair writes are interleaved. Consequently, it requires extensive synchronized I/O to maintain the correct ordering (e.g., three sync operations are required in Figure 3b), which may incur severe performance overhead. Second, as part of e2fsck, the logging feature is closely tied to Ext-family file systems, and thus it cannot benefit other file system checkers directly. We address the limitations by building a general logging library called rfsck-lib.

5.3.1 Similarities Among File System Checkers

Different file systems usually vary a lot in terms of ondisk layouts and consistency rules. However, there are similarities among different checkers, which makes designing a general and efficient solution possible.

First of all, as a user-level utility, file system checkers always repair corrupted images through a limited number of system calls, which are irrelevant to file systems' internal structures and consistency rules. Moreover, based on our survey on popular file system checkers (e.g., e2fsck, xfs_repair, fsck.f2fs), we find that they always use write system calls (e.g., pwrite and its variants) instead of other memory-based system calls (e.g., mmap, msync). Therefore, only a few writes may cause potential cascading corruptions under faults. In other words, by focusing on the writes, we may improve different checkers.

Second, there is natural locality in checkers. Similar to the cylinder groups of FFS [56], many modern file systems have a layout consisting of relatively independent areas with an identical structure (e.g., block groups of Ext4 [4], allocation groups of XFS [70], and cubes of IceFS [52]). Among others, such common design enables co-locating related files to mitigate file system aging [33, 68] while isolating unrelated files. From the checker's perspective, most consistency rules within each area may be checked locally without referencing other areas. Also, each type of metadata usually has its unique structure and consistency rules (e.g., the rec_len of each directory entry in an Ext4 inode should be within a range). These local consistency rules may be checked independently without cross-checking other metadata.

Due to the locality, checkers usually consist of relatively self-contained components. For example, e2fsck includes five passes for checking different sets of consistency rules (§2.1). Similarly, xfs_repair includes seven passes, and it forks multiple threads to check multiple allocation groups separately (§2.2). Such locality exists even without changing the file system layout or reordering the checking of consistency rules [54]. Therefore, it is possible to split an existing checker into several pieces and isolate their impact under faults.

Based on the observations above, we describe rfsck-lib's design in the following subsections.

5.3.2 Basic Redo Logging

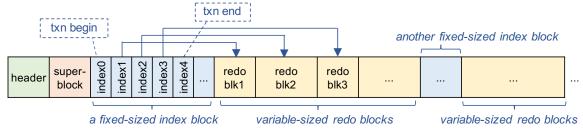
A corrupted file system image is repaired by a checker through a set of repair writes. If the checker finishes without interruption, the set of writes turn the image back to a consistent state. On the other hand, if the checker is interrupted, only a subset of writes changes the image, and the resulting state may become uncorrectable. Therefore, the key of preventing uncorrectable states is to maintain the atomicity of the checker's writes. To this end, rfsck-lib redirects the checker's writes to a log first, and then repairs the image based on the log. Essentially, it implements a redo logging scheme [58].

As shown in Figure 3c, all repair writes are issued to the redo log first (i.e., "redo blk1" to "redo blk3"). After the write of the last redo block (i.e., "redo blk3"), a sync operation (i.e., the red solid line) is issued to make the redo blocks persistent. After the sync operation returns, the image is repaired (i.e., "repair blk1" to "repair blk3") based on the redo log. Compared with e2fsck-patch in Figure 3b, the log writes and the repair writes are separated, and the required number of sync operations is reduced from three to one. Such improvement in terms of sync overhead can be more significant if more blocks on the image need to be repaired.

5.3.3 Fine-grained Logging with Safe Transactions

While the basic redo logging scheme reduces the ordering constraint to minimum, there is one limitation: if a fault happens before the final sync operation finishes, all checking and repairing effort may be lost. In some complicated cases where the checker may take hours to finish [54], the waste is undesirable. On the other hand, a checker may be split into relatively independent pieces due to the locality (§5.3.1). Therefore, rfsck-lib extends the basic redo logging with safe transactions.

A safe transaction is a set of repair writes which will not lead to uncorrectable inconsistencies if they are written to the file system image atomically. In the simplest case, the whole checker (i.e., the complete set of all re-



The log format of rfsck-lib. The log includes a header, a superblock, fixed-sized index blocks, and variable-sized redo blocks. Each index block includes a fixed number of indexes. Each index can either describe the beginning/end of a transaction (i.e., "txn begin"/"txn end"), or describe one variable-sized redo block. "index0" to "index4" describe one safe transaction with three redo blocks (i.e., "redo blk1" to "redo blk3") in this example.

pair writes) is one safe transaction. At a finer granularity, each pass of the checker (or the check of each allocation group) may be considered as one safe transaction. While a later pass may depend on the result of a previous pass, the previous pass is executed without any dependency on the later passes. Therefore, by guaranteeing the atomicity of each pass as well as the ordering among pass-based safe transactions, the repair writes may be committed in several batches without introducing uncorrectable inconsistencies. In the extreme case, the checking and repairing of each individual consistency rule may be considered as one safe transaction.

Figure 3c illustrates the safe transactions. In the simplest case, all three redo blocks (i.e., "redo blk1" to "redo blk3") constitute one safe transaction, and only one sync operation (i.e., the red solid line) is needed, same as the basic redo logging (§5.3.2). At a finer granularity, the first two redo blocks (i.e., "redo blk1" and "redo blk2") may constitute one safe transaction (e.g., updating an inode and the corresponding bitmap), and the third block itself (i.e., "redo blk3") may be another safe transaction (e.g., updating another inode). Another sync operation (i.e., the red dash line) is issued between the two transactions to guarantee the correct ordering. If a crash happens between the two sync operations, the first safe transaction (i.e., "redo blk1" and "redo blk2") is still valid. In this case, instead of re-calculating the rules and regenerating the blocks, the checker can directly replay the valid transaction from the log after restart.

In summary, a checker may be logged as one or more safe transactions. Compared to the basic redo logging, such fine-grained control avoids losing all recovery effort before the fault. On the other hand, maintaining the atomicity as well as the ordering requires additional sync operations. So there is a tradeoff between the transaction granularity and the transaction overhead. Since different systems may have different preferences, rfsck-lib simply provides an interface to define safe transactions, without restricting the number of transactions.

5.3.4 Log Format

To support the redo logging with safe transactions, rfsck-lib uses a special log format extended from e2fsck-undo. As shown in Figure 4, the log includes a header, a superblock, fixed-sized index blocks, and variable-sized redo blocks.

The header starts with a magic number to distinguish the log from other files. Besides, it includes the offsets of the superblock and the first index block, the total number of index blocks, a flag showing whether the log has been replayed, and a checksum of the header itself.

The superblock is copied from the file system image to be repaired, which is used to match the log with the image to avoid replaying an irrelevant log to the image.

The index block includes a fixed number of indexes. Each index can describe the beginning of a transaction (i.e., "txn begin"), the end of a transaction (i.e., "txn end"), or one variable-sized redo block. Therefore, a group of indexes can describe one safe transaction together. For example, in Figure 4 five indexes (i.e., "index0" to "index4") describe one safe transaction with three redo blocks (i.e., "redo blk1" to "redo blk3").

As shown in Table 5, an index has 16 bytes consisting of three fields. To describe one redo block, the first field (i.e., uint32_t cksum) stores a checksum of the redo block, the second field (i.e., uint32_t size) stores its size, and the third field (i.e., uint64_t fs_lba) stores its logical block address (LBA) in the file system image.

To describe "txn begin" or "txn end", the first field of the index is repurposed to store a transaction ID instead of a checksum, which marks the boundary of indexes belonging to the same transaction. The second field (size) is set to zero. Since a valid redo block must have a non-zero size, rfsck-lib can differentiate "txn begin" or "txn end" indexes from those describing redo blocks even if a transaction ID happens to collide with a checksum. In addition, the "txn begin" index uses the third field to denote whether the transaction has been replayed or not, and the "txn end" index uses the third field to store a checksum of all indexes in the transaction.

Field	Description
uint32_t cksum	checksum of the redo block
$\mathtt{uint32}_{-}\mathtt{t}\ \mathtt{size}$	size of the redo block
$uint64_t fs_lba$	LBA in the file system image

Table 5: The structure of an index.

Function	Description
rfsck_get_sb	get the superblock
rfsck_open	create a redo log
$rfsck_txn_begin$	begin a safe transaction
$rfsck_write$	write a redo block
$rfsck_txn_end$	end of a safe transaction
rfsck_replay	replay the redo log
rfsck_close	close the redo log

Table 6: **The interface of rfsck-lib**. rfsck_get_sb is a wrapper function for invoking file-system-specific procedure to get the superblock, while the others are filesystem agnostic.

For each write of the checker, rfsck-lib creates an index in the index block and then append the content of the write to the area after the index block as a redo block. Since the writes may have different sizes, the redo blocks may vary in size as well. However, since all other metadata blocks (i.e., header, superblock, index blocks) have known fixed sizes, the offset of a redo block in the log can be identified by accumulating the sizes of all previous blocks. In other words, there is no need to maintain the offsets of redo blocks in the log.

When an index block becomes full, another index block is allocated after the previous redo blocks (which are described by the previous index block). In this way, rfsck-lib can support various numbers of writes and transactions.

5.3.5 Interface

To enable easy integration with existing checkers, rfsck-lib provides a simple interface. As shown in Table 6, there are seven function calls in total. The first function (rfsck_get_sb) is a wrapper for invoking a file-system-specific procedure to get the superblock, which is written to the second part of the log (Figure 4). Since all checkers need to read the superblock anyway, rfsck_get_sb can wrap around the existing procedure.

rfsck_open is used to create a log file at a given path at the beginning of the checker procedure. Internally, rfsck-lib initializes the metadata blocks of the log.

rfsck_txn_begin is used to denote the beginning of a safe transaction, which creates a "txn begin" index in the log. Similarly, rfsck_txn_end denotes the end of a transaction, which generates a "txn end" index and sync all updates to the log. All writes between rfsck_txn_begin and rfsck_txn_end are replaced with rfsck_write, which creates a redo block and the corresponding index in the log.

rfsck_replay is used to replay logged transactions to the file system image. Besides, similar to the e2undo utility [3], the replay functionality is also implemented as an independent utility called rfsck-redo, which can replay an existing (potentially incomplete) log to a file system image (e.g., after the checker is interrupted). rfsck-redo first verifies if the log belongs to the image (based on the superblock). If yes, rfsck-redo further verifies the integrity of the log based on metadata and then replays valid transactions. Note that the additional verifications are only needed when the log is replayed by rfsck-redo. The rfsck_replay function can skip these verifications as it is invoked directly after the logging by the (uninterrupted) checker.

Finally, rfsck_close is used at the end of the checker to release all resources used by rfsck-lib and exist.

5.3.6 Limitations

The current prototype of rfsck-lib is not threadsafe. Therefore, if a checker is multi-threaded (e.g., xfs_repair), using rfsck-lib may require additional attention to avoid race conditions on logging. However, as we will demonstrate in §5.4 and §6, rfsck-lib can still be applied to strengthen xfs_repair.

In addition, rfsck-lib only provides an interface, which requires manual modification of the source code. Since the modification is simple, we expect the manual effort to be acceptable. Also, it is possible to use compiler infrastructures [11, 13] to automate the code instrumentation, which we leave as future work.

5.4 Integration with Existing Checkers

Strengthening an existing checker using rfsck-lib is straightforward given the simple interface (§5.3.5). To demonstrate the practicality, we first integrate rfsck-lib with e2fsck, and create a robust checker for Ext-family file systems (i.e., rfsck-ext).

There are potential writes to the file system image in each pass of e2fsck (including the first scanning pass), so we create a safe transaction for each pass. Moreover, within Pass-1 and Pass-2 (§2.1), there are a few places where e2fsck explicitly flushes the writes to the image and restarts scanning from the beginning (via goto statement). In other words, the restarted scanning (and subsequent passes) requires the previous writes to be visible on the image. In this case, we insert additional rfsck_txn_end and rfsck_replay before the goto statement to guarantee that previous writes are visible on the image for re-scanning. We add a "-R" option to allow the user to specify the log path via command line. In total, we add 50 LoC to e2fsck.

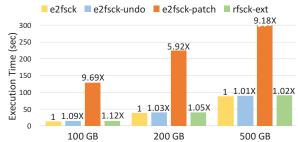


Figure 5: Performance comparison of e2fsck, e2fsckundo, e2fsck-patch, and rfsck-ext. This figure compares the execution time of e2fsck, e2fsck-undo, e2fsckpatch, and rfsck-ext. The y-axis shows the execution time in seconds. The x-axis shows file system sizes. The number above each bar indicates the normalized time (relative to e2fsck). Note: e2fsck and e2fsck-undo are not resilient to faults.

Similarly, we also integrate rfsck-lib with xfs_repair, and create a robust checker for XFS file system (i.e., rfsck-xfs). As mentioned in ($\S 2.2$), one feature of xfs_repair is multi-threading: it forks multiple threads to repair multiple allocation groups in parallel. The threads update in-memory structures concurrently, and the main thread writes all updates to the image at the end. Although it is possible to encapsulate each repair thread into one safe transaction, doing so requires additional concurrency control. To minimize the modification, we simply treat the whole repair procedure as one transaction. Since all writes are issued by the main thread, there is no race condition for rfsck-lib. We also add a "-R" command line option. In total, we add 15 LoC to xfs_repair.

Evaluation 6

In this section, we evaluate rfsck-ext and rfsck-xfs in terms of robustness ($\S6.1$) and performance ($\S6.2$).

Our experiments were conducted on a machine with a Intel Xeon 3.00GHz CPU, 8GB main memory, and two WD5000AAKS hard disks. The operating system is Ubuntu 16.04 LTS with kernel v4.4. To evaluate the robustness, we used the test images reporting corruption under e2fsck-undo (§3.4) and xfs_repair (§3.5). To evaluate the performance, we created another set of images with practical sizes, and measured the execution time of e2fsck, e2fsck-undo, e2fsck-patch, rfsck-ext, xfs_repair, and rfsck-xfs. For each checker, we report the average time of three runs.

In general, we demonstrate that both rfsck-ext and rfsck-xfs can survive fault injection experiments. Also, both checkers incur reasonable performance overhead (i.e., up to 12%) compared to the original unreliable versions. Moreover, rfsck-ext outper-

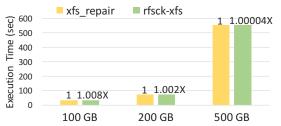


Figure 6: Performance comparison of xfs_repair and rfsck-xfs. This figure compares the execution time of xfs_repair and rfsck-xfs. The y-axis shows the execution time in seconds. The x-axis shows file system sizes. The number above each bar indicates the normalized time (relative to xfs_repair). Note: xfs_repair is not resilient to faults.

forms e2fsck-patch by up to 9 times while achieving the same level of robustness.

6.1 Robustness

As discussed in §3, when injecting faults at the 4KB granularity, 17 Ext4 test images report corruptions under e2fsck (Table 1), and 12 XFS test images report corruptions under xfs_repair (Table 4). We use these test images to trigger rfsck-ext and rfsck-xfs, respectively. Since both checkers have the logging support, we use the method in Figure 1b to evaluate them.

For rfsck-ext, all 17 test images report no corruptions. Similarly, for rfsck-xfs, all 12 test images report no corruptions. This result verifies that rfsck-lib can help improve the fault resilience of existing checkers.

Performance

The test images used in §3 are created as regular files, and they are small in size (i.e., 8MB to 128MB). Therefore, they are unsuitable for evaluating the execution time of checkers. So we create another set of Ext4 and XFS test images with practical sizes (i.e., 100G, 200GB, 500GB) on real hard disks. We first fill up the entire file system by running fs_mark [5] for five times. Each time fs_mark fills up 20% of the capacity by creating directories and files with a certain size. The file size is a random value between 4KB to 1MB, which is relatively small in order to maximize the number of inodes used. After filling up the entire file system, we inject 2 random bit corruptions to the metadata using either debugfs [1] (for Ext4) or blocktrash [14] (for XFS). We measure the execution time of checkers on corrupted images, and verify that the repair results of rfsck-ext and rfsck-xfs are the same as that of the original checkers.

Figure 5 compares the execution time of e2fsck, e2fsck-undo, e2fsck-patch, and rfsck-ext on different images. For each size of image, the bars represent the execution time in seconds (y-axis). Also, the number above each bar shows the normalized execution time (relative to e2fsck). We can see that rfsck-ext incurs up to 12% overhead, while e2fsck-patch may incur more than 8 times overhead due to extensive sync operations.

Also, we can see that as the size of file system increases, the overhead of rfsck-ext decreases. This is because the execution time of rfsck-ext is largely dominated by the scanning in Pass-1 (§2.1) which is proportional to the file system size, similar to e2fsck [54].

Similarly, Figure 6 compares the execution time of xfs_repair and rfsck-xfs. We can see that rfsck-xfs incurs up to 0.8% overhead, and the overhead also decreases as the file system size increases.

Note that our aging method is relatively simple compared to other aging techniques [33, 68]. Also, the 2random-bit corruption may not necessarily lead to extensive repair operations of checkers. Therefore, the execution time measured here may not reflect the complexity of checking and repairing real-world file systems (which may take hours [35, 34, 54, 69]). We leave generating more representative file systems as future work.

7 Discussion

Co-designing file systems and checkers. Recent work has demonstrated the benefits of co-designing file systems and checkers. For example, by co-designing rext3 and ffsck, ffsck may be 10 times faster than e2fsck [54]. In contrast, rfsck-lib is designed to be file system agnostic, which makes it directly applicable to existing systems. We believe checkers may be improved further in terms of both reliability and performance by co-designing, and we leave it as future work.

Other reliability techniques. There are other techniques which may mitigate the impact of an inconsistent file system image or the loss of an entire image (e.g., replication [39]). However, maintaining the consistency of local file systems and improving the checkers is still important for many reasons. For example, a consistent local file system is the building block of large-scale file systems, and the local checker may be the foundation of higher-level recovery procedures (e.g., lfsck [6]). Therefore, our work is orthogonal to these other efforts.

Robustness. We evaluate the robustness of checkers based on fault injection experiments in this work. The test images we use are limited, and may not cover all corruption scenarios or trigger all code paths of the checkers. There are other techniques (e.g., symbolic execution and formal verification) which might provide more coverage, and we leave it as future work.

Related Work

Reliability of file system checkers. Gunawi et al. [42] find that the Ext2 checker may create inconsistent or even insecure repairs; they then propose a more elegant checker (i.e., SQCK) based on a declarative guery language. Carreira et al. [27] propose a tool (i.e., SWIFT) to test checkers using a mix of symbolic and concrete execution; they tested five popular checkers and found bugs in all of them. Ma et al. [54] change the structure of Ext3 and co-design the checker, which enables faster checking and thus narrows the window of vulnerability. Generally, these studies consider the behavior of checkers during normal executions (i.e., no interruption). Complimentarily, we study checkers under faults.

Reliability of file systems. Great efforts have been put towards improving the reliability of file systems [23, 29, 32, 36, 43, 51, 55, 57, 62, 67, 77, 79]. For example, Prabhakaran et al. [62] analyze the failure policies of four file systems and propose the IRON file system which implements a family of novel recovery techniques. Fryer et al. [36] transform global consistency rules to local consistency invariants and enable fast runtime checking. CrashMonkey [55] provides a framework to automatically test the crash consistency of file systems. Overall, these research help understand and improve the reliability of file systems, which may reduce the need for checkers. However, despite these efforts, checkers remain a necessary component for most file systems.

Reliability of storage devices. In terms of storage devices, research efforts are also abundant [19, 20, 30, 47, 63, 64]. For example, Bairavasundaram et al. [19, 20] analyze the data corruption and latent sector errors in production systems containing a total of 1.53 million HDDs. Besides HDDs, more recent work has been focused on flash memory and solid state drives (SSDs) [18, 22, 24, 25, 26, 28, 37, 40, 46, 48, 50, 53, 59, 65, 71, 73, 76, 78, 81, 82]. These studies provide valuable insights for understanding file system corruptions caused by hardware.

9 Conclusion

We have studied the behavior of file system checkers under faults. We find that running the checker after an interrupted repair may not return the file system to a valid state. To address the issue, we have built a general logging library which can help strengthen existing checkers with little modification. We hope our work will raise the awareness of reliability vulnerabilities in storage systems, and facilitate building truly fault-resilient systems.

10 Acknowledgements

We thank the anonymous reviewers and Keith Smith (our shepherd) for their insightful feedback. We also thank Linux practitioners including Theodore Ts'o and Ric Wheeler for the invaluable discussion. This work was supported in part by NSF under grants CNS-1566554 and CCF-1717630. Any opinions, findings, and conclusions expressed in this material are those of the authors and do not necessarily reflect the views of NSF.

References

- [1] debugfs. http://man7.org/linux/man-pages/man8/ debugfs.8.html.
- [2] Discussion with Theodore Ts'o at Linux FAST Summit'17. https://www.usenix.org/conference/ linuxfastsummit17.
- [3] E2fsprogs: Ext2/3/4 Filesystem Utilities. http: //e2fsprogs.sourceforge.net/.
- [4] Ext4 File System. https://ext4.wiki.kernel.org/ index.php/Main_Page.
- [5] fs_mark: Benchmark file creation. https://github. com/josefbacik/fs_mark.
- [6] LFSCK: an online file system checker for Lustre. https://github.com/Xyratex/lustre-stable/blob/ master/Documentation/Ifsck.txt.
- [7] Linux Programmer's Manual: O_SYNC flag for open. http://man7.org/linux/man-pages/man2/ open.2.html.
- [8] Linux SCSI target framework (tgt). http://stgt. sourceforge.net/.
- [9] Lustre File System. http://opensfs.org/lustre/.
- [10] Prototypes of rfsck-test, e2fsck-patch, refsck-lib, refsck-ext, rfsck-xfs. https://www.cs.nmsu.edu/ ~mzheng/lab/lab.html.
- [11] ROSE Compiler Infrastructure. http://rosecompiler.
- [12] SQLite documents. http://www.sqlite.org/docs. html.
- [13] The LLVM Compiler Infrastructure. https://llvm. org/.
- [14] XFS File System Utilities. https: //access.redhat.com/documentation/en-US/Red_Hat_Enterprise_Linux/6/html/ Storage_Administration_Guide/xfsothers.html.
- [15] [PATCH 1/3] e2fsprogs: Add undo I/O manager. http://lists.openwall.net/linux-ext4/2007/07/ 25/2, 2007.
- [16] [PATCH 16/31] e2undo: ditch tdb file, write everything to a flat file. http://lists.openwall.net/linuxext4/2015/01/08/1, 2015.

- [17] High Performance Computing Center (HPCC) Power Outage Event. Email Annoucement by HPCC, Monday, January 11, 2016 at 8:50:17 AM CST. https://www.cs.nmsu.edu/~mzheng/ docs/failures/2016-hpcc-outage.pdf, 2016.
- [18] Nitin Agarwal, Vijayan Prabhakaran, Ted Wobber, John D. Davis, Mark Manasse, and Rina Panigrahy. Design tradeoffs for SSD performance, 2008.
- [19] 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. ACM Transactions on Storage, 4(3):8:1–8:28, November 2008.
- [20] 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 International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS'07), pages 289-300, 2007.
- [21] Luiz Andre Barroso and Urs Hoelzle. The Datacenter As a Computer: An Introduction to the Design of Warehouse-Scale Machines. Morgan and Claypool Publishers, 1st edition, 2009.
- [22] Hanmant P Belgal, Nick Righos, Ivan Kalastirsky, Jeff J Peterson, Robert Shiner, and Neal Mielke. A new reliability model for post-cycling charge retention of flash memories. In Proceedings of the 40th Annual Reliability Physics Symposium, pages 7–20. IEEE, 2002.
- [23] James Bornholt, Antoine Kaufmann, Jialin Li, Arvind Krishnamurthy, Emina Torlak, and Xi Wang. Specifying and checking file system crash-consistency models. Proceedings of the 21st International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'16), 51(4):83-98, 2016.
- [24] Adam Brand, Ken Wu, Sam Pan, and David Chin. Novel read disturb failure mechanism induced by FLASH cycling. In Proceedings of the 31st Annual Reliability Physics Symposium, pages 127-132. IEEE, 1993.
- [25] Yu Cai, Erich F. Haratsch, Onur Mutlu, and Ken Mai. Error Patterns in MLC NAND Flash Memory: Measurement, Characterization, and Analysis. In Proceedings of the Conference on Design, Automation and Test in Europe (DATE'12), pages 521-526, 2012.

- [26] Yu Cai, Gulay Yalcin, Onur Mutlu, Erich F Haratsch, Osman Unsal, Adrian Cristal, and Ken Mai. Neighbor-cell assisted error correction for MLC NAND flash memories. In ACM SIGMET-RICS Performance Evaluation Review, volume 42, pages 491–504. ACM, 2014.
- [27] João Carlos Menezes Carreira, Rodrigo Rodrigues, George Candea, and Rupak Majumdar. Scalable testing of file system checkers. In Proceedings of the 7th ACM European Conference on Computer Systems (EuroSys'12), pages 239–252, 2012.
- [28] Feng Chen, David A. Koufaty, and Xiaodong Zhang. Understanding intrinsic characteristics and system implications of flash memory based solid state drives. In Proceedings of the ACM Joint International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS'09), 2009.
- [29] Haogang Chen, Daniel Ziegler, Tej Chajed, Adam Chlipala, M Frans Kaashoek, and Nickolai Zeldovich. Using crash hoare logic for certifying the fscq file system. In *Proceedings of the 25th Sympo*sium on Operating Systems Principles (SOSP'15), pages 18-37. ACM, 2015.
- [30] Peter M. Chen, Edward K. Lee, Garth A. Gibson, Randy H. Katz, and David A. Patterson. RAID: high-performance, reliable secondary storage. ACM Computing Surveys, 26(2):145-185, June 1994.
- [31] 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.
- [32] Vijay Chidambaram, Tushar Sharma, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. Consistency Without Ordering. In Proceedings of the 10th Conference on File and Storage Technologies (FAST'12), February 2012.
- [33] Alex Conway, Ainesh Bakshi, Yizheng Jiao, William Jannen, Yang Zhan, Jun Yuan, Michael A. Bender, Rob Johnson, Bradley C. Kuszmaul, Donald E. Porter, and Martin Farach-Colton. File systems fated for senescence? nonsense, says science! In Proceedings of the 15th USENIX Conference on File and Storage Technologies (FAST'17), pages 45-58, 2017.
- e2fsck is taking forever. [34] GParted Forum. http://gparted-forum.surf4.info/viewtopic.php?id= 13613, 2009.

- [35] JaguarPC Forum. How long does it take FSCK to run?! http://forums.jaguarpc.com/hostingtalk-chit-chat/14217-how-long-does-take-fsckrun.html, 2006.
- [36] 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 Conference on File and Storage Technologies (FAST'12), February 2012.
- [37] Ryan Gabrys, Eitan Yaakobi, Laura M. Grupp, Steven Swanson, and Lara Dolecek. Tackling intracell variability in TLC flash through tensor product codes. In Proceedings of IEEE International Symposium of Information Theory, pages 1000-1004, 2012.
- [38] Gregory R Ganger, Marshall Kirk McKusick, Craig AN Soules, and Yale N Patt. Soft updates: a solution to the metadata update problem in file systems. ACM Transactions on Computer Systems (TOCS'00), 18(2):127–153, 2000.
- [39] Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung. The Google File System. In Proceedings of the 9th ACM Symposium on Operating Systems Principles (SOSP'03), pages 29-43, 2003.
- [40] 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 (MI-CRO'09), pages 24–33, 2009.
- [41] Haryadi S Gunawi, Mingzhe Hao, Riza O Suminto, Agung Laksono, Anang D Satria, Jeffry Adityatama, and Kurnia J Eliazar. Why does the cloud stop computing? lessons from hundreds of service outages. In Proceedings of ACM Symposium on Cloud Computing (SoCC'16), pages 1-16, 2016.
- [42] Haryadi S. Gunawi, Abhishek Rajimwale, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. Sqck: A declarative file system checker. In Proceedings of the 8th USENIX Conference on Operating Systems Design and Implementation (OSDI'08), pages 131–146, 2008.
- [43] Harvadi S Gunawi, Cindy Rubio-González, Andrea C Arpaci-Dusseau, Remzi H Arpaci-Dusseau, and Ben Liblit. Eio: Error handling is occasionally correct. In Proceedings of the 6th USENIX Conference on File and Storage Technologies (FAST'08), volume 8, pages 1-16, 2008.

- [44] Zhenyu Guo, Sean McDirmid, Mao Yang, Li Zhuang, Pu Zhang, Yingwei Luo, Tom Bergan, Madan Musuvathi, Zheng Zhang, and Lidong Zhou. Failure Recovery: When the Cure Is Worse Than the Disease. In Proceedings of the 14th Workshop on Hot Topics in Operating Systems (HotOS'13), 2013.
- [45] Cheng Huang, Huseyin Simitci, Yikang Xu, Aaron Ogus, Brad Calder, Parikshit Gopalan, Jin Li, and Sergey Yekhanin. Erasure coding in windows azure storage. In Proceedings of the 2012 USENIX Conference on Annual Technical Conference (USENIX ATC'12), pages 15-26, 2012.
- [46] Xavier Jimenez, David Novo, and Paolo Ienne. Wear unleveling: improving nand flash lifetime by balancing page endurance. In Proceedings of the 12th USENIX Conference on File and Storage Technologies (FAST'14), pages 47-59, 2014.
- [47] Andrew Krioukov, Lakshmi N Bairavasundaram, Garth R Goodson, Kiran Srinivasan, Randy Thelen, Andrea C Arpaci-Dusseau, and Remzi H Arpaci-Dusseau. Parity lost and parity regained. In Proceedings of the 13th USENIX Conference on File and Storage Technologies (FAST'08), volume 8, pages 1–15, 2008.
- [48] H Kurata, K Otsuga, A Kotabe, S Kajiyama, T Osabe, Y Sasago, S Narumi, K Tokami, S Kamohara, and O Tsuchiya. The impact of random telegraph signals on the scaling of multilevel flash memories. In Proceedings of the 2006 Symposium on VLSI *Circuits*, pages 112–113. IEEE, 2006.
- [49] Changman Lee, Dongho Sim, Joo-Young Hwang, and Sangyeun Cho. F2fs: A new file system for flash storage. In Proceedings of the 13th USENIX Conference on File and Storage Technologies (FAST'15), pages 273-286, 2015.
- [50] Jiangpeng Li, Kai Zhao, Xuebin Zhang, Jun Ma, Ming Zhao, and Tong Zhang. How much can data compressibility help to improve nand flash memory lifetime? In Proceedings of the 13th USENIX Conference on File and Storage Technologies (FAST'15), pages 227–240, 2015.
- [51] Lanyue Lu, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau, and Shan Lu. A study of linux file system evolution. In Proceedings of the 11th USENIX Conference on File and Storage Technologies (FAST'13), pages 31-44, 2013.
- [52] Lanyue Lu, Yupu Zhang, Thanh Do, Samer Al-Kiswany, Andrea C. Arpaci-Dusseau, and

- Remzi H. Arpaci-Dusseau. Physical disentanglement in a container-based file system. In Proceedings of the 11th USENIX Symposium on Operating Systems Design and Implementation (OSDI'14), pages 81–96, 2014.
- [53] Youyou Lu, Jiwu Shu, Weimin Zheng, et al. Extending the lifetime of flash-based storage through reducing write amplification from file systems. In Proceedings of the 11th USENIX Conference on File and Storage Technologies (FAST'13), volume 13, 2013.
- [54] Ao Ma, Chris Dragga, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. ffsck: The fast file system checker. In Proceedings of the 11th USENIX Conference on File and Storage Technologies (FAST'13), pages 1–15, 2013.
- [55] Ashlie Martinez and Vijay Chidambaram. Crash-Monkey: A Framework to Automatically Test File-System Crash Consistency. In Proceedings of the 9th USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage'17), 2017.
- [56] Marshall K. McKusick, William N. Joy, Samuel J. Leffler, and Robert S. Fabry. A fast file system for unix. Proceedings of the ACM Transactions on Computer Systems (TOCS'84), 2(3):181–197, August 1984.
- [57] Changwoo Min, Sanidhya Kashyap, Byoungyoung Lee, Chengyu Song, and Taesoo Kim. Crosschecking semantic correctness: The case of finding file system bugs. In Proceedings of the 25th Symposium on Operating Systems Principles (SOSP'15), pages 361-377. ACM, 2015.
- [58] C. Mohan, Don Haderle, Bruce Lindsay, Hamid Pirahesh, and Peter Schwarz. ARIES: A Transaction Recovery Method Supporting Fine-granularity Locking and Partial Rollbacks Using Write-ahead Logging. ACM Transactions on Database Systems (TODS'92), 1992.
- [59] T. Ong, A. Frazio, N. Mielke, S. Pan, N. Righos, G. Atwood, and S. Lai. Erratic Erase In ETOX/sup TM/ Flash Memory Array. In Proceedings of the Symposium on VLSI Technology (VLSI'93), 1993.
- [60] Lluis Pamies-Juarez, Filip Blagojević, Robert Mateescu, Cyril Gyuot, Eyal En Gad, and Zvonimir Bandić. Opening the chrysalis: On the real repair performance of MSR codes. In Proceedings of the 14th USENIX Conference on File and Storage Technologies (FAST'16), pages 81–94, 2016.

- [61] Thanumalayan Sankaranarayana Pillai, 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 USENIX Symposium on Operating Systems Design and Implementation (OSDI'14), October 2014.
- [62] 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), pages 206-220, October 2005.
- [63] Abhishek Rajimwale, Vijay Chidambaram, Deepak Ramamurthi, Andrea C Arpaci-Dusseau, and Remzi H Arpaci-Dusseau. Coerced cache eviction and discreet mode journaling: Dealing with misbehaving disks. In Proceedings of the 2011 IEEE/IFIP 41st International Conference on Dependable Systems & Networks (DSN'11), pages 518-529. IEEE, 2011.
- [64] Bianca Schroeder and Garth A. Gibson. Disk failures in the real world: What does an MTTF of 1,000,000 hours mean to you? In Proceedings of the 5th USENIX Conference on File and Storage Technologies (FAST'07), 2007.
- [65] 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), pages 67-80, February 2016.
- [66] Konstantin Shvachko, Hairong Kuang, Sanjay Radia, and Robert Chansler. The hadoop distributed file system. In Proceedings of the IEEE 26th Symposium on Mass Storage Systems and Technologies (MSST'10), pages 1–10. IEEE, 2010.
- [67] Helgi Sigurbjarnarson, James Bornholt, Emina Torlak, and Xi Wang. Push-button verification of file systems via crash refinement. In Proceedings of 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI'16), 2016.
- [68] Keith A. Smith and Margo I. Seltzer. File system aging—increasing the relevance of file system benchmarks. In Proceedings of the 1997 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems (SIG-METRICS'97), pages 203-213, 1997.

- [69] V. Svanberg. Fsck takes too long on multiplyclaimed blocks. http://old.nabble.com/Fscktakes-too-long-on-multiply-claimed-blockstd21972943.html, 2009.
- [70] Adam Sweeney, Doug Doucette, Wei Hu, Curtis Anderson, Mike Nishimoto, and Geoff Peck. Scalability in the XFS File System. In Proceedings of the 1996 USENIX Annual Technical Conference (USENIX ATC'96), volume 15, 1996.
- [71] Huang-Wei Tseng, Laura M. Grupp, and Steven Swanson. Understanding the impact of power loss on flash memory. In Proceedings of the 48th Design Automation Conference (DAC'11), 2011.
- [72] Stephen C. Tweedie. Journaling the linux ext2fs filesystem. In Proceedings of the 4th Annual Linux Expo, 1998.
- [73] Simeng Wang, Jinrui Cao, Danny V Murillo, Yiliang Shi, and Mai Zheng. Emulating Realistic Flash Device Errors with High Fidelity. In *Proceedings of* the IEEE International Conference on Networking, Architecture and Storage (NAS'16). IEEE, 2016.
- [74] Sage A. Weil, Scott A. Brandt, Ethan L. Miller, Darrell D. E. Long, and Carlos Maltzahn. Ceph: A Scalable, High-performance Distributed File System. In Proceedings of the 7th Symposium on Operating Systems Design and Implementation (OSDI'06), pages 307-320, 2006.
- [75] Mingyuan Xia, Mohit Saxena, Mario Blaum, and David A. Pease. A tale of two erasure codes in HDFS. In Proceedings of the 13th USENIX Conference on File and Storage Technologies (FAST'15), pages 213-226, 2015.
- [76] Gala Yadgar, Eitan Yaakobi, and Assaf Schuster. Write once, get 50% free: Saving ssd erase costs using wom codes. In Proceedings of the 13th USENIX Conference on File and Storage Technologies (FAST'15), pages 257–271, 2015.
- [77] Junfeng Yang, Can Sar, and Dawson Engler. EX-PLODE: a lightweight, general system for finding serious storage system errors. In Proceedings of the Seventh Symposium on Operating Systems Design and Implementation (OSDI'06), pages 131-146, November 2006.
- [78] Yiying Zhang, Leo Prasath Arulraj, Andrea C Arpaci-Dusseau, and Remzi H Arpaci-Dusseau. De-indirection for flash-based ssds with nameless writes. In Proceedings of the 10th USENIX Conference on File and Storage Technologies (FAST'12), 2012.

- [79] 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 Conference on File and Storage Technologies (FAST'10), pages 29-42, 2010.
- [80] Mai Zheng, Joseph Tucek, Dachuan Huang, Feng Qin, Mark Lillibridge, Elizabeth S. Yang, Bill W Zhao, and Shashank Singh. Torturing Databases for Fun and Profit. In Proceedings of the 11th USENIX Symposium on Operating Systems Design and Implementation (OSDI'14), pages 449-464, 2014.
- [81] Mai Zheng, Joseph Tucek, Feng Qin, and Mark Lillibridge. Understanding the robustness of SSDs under power fault. In Proceedings of the 11th USENIX Conference on File and Storage Technologies (FAST'13), 2013.
- [82] Mai Zheng, Joseph Tucek, Feng Qin, Mark Lillibridge, Bill W Zhao, and Elizabeth S. Yang. Reliability Analysis of SSDs under Power Fault. In Proceedings of the ACM Transactions on Computer Systems (TOCS'16), 2016.