

# **Model-Checking Support for File System Development**

Yifei Liu  
Department of Computer Science  
Stony Brook University

*Research Proficiency Exam*

Technical report FSL-22-01

January 07, 2022

# Contents

<b>List of Figures</b>	<b>2</b>
<b>List of Algorithms</b>	<b>2</b>
<b>1 Introduction and Motivation</b>	<b>1</b>
<b>2 Background</b>	<b>3</b>
2.1 Model Checking . . . . .	3
2.1.1 Model Checking Concepts . . . . .	3
2.1.2 Spin Model Checker . . . . .	4
2.2 File System Bugs . . . . .	5
<b>3 MCFS Design</b>	<b>7</b>
3.1 Design Goals . . . . .	7
3.2 MCFS Architecture . . . . .	7
<b>4 Challenges</b>	<b>9</b>
4.1 Access to In-Memory States . . . . .	9
4.2 Cache Incoherency . . . . .	9
4.3 State Explosion . . . . .	10
4.4 False Positives . . . . .	11
<b>5 MCFS Prototype Implementation</b>	<b>12</b>
<b>6 Tracking File-System States</b>	<b>14</b>
<b>7 Evaluation</b>	<b>16</b>
7.1 Performance and Memory Demands . . . . .	16
7.2 Assisting File System Development . . . . .	18
<b>8 Related Work</b>	<b>19</b>
8.1 Regression Test Suites . . . . .	19
8.2 Verified or Machine-Checkable File Systems . . . . .	19
8.3 Model Checking and Verification . . . . .	20
8.4 Fuzzing . . . . .	20
<b>9 Conclusions</b>	<b>22</b>
<b>10 Future Work</b>	<b>23</b>

# List of Figures

2.1	An example of a file system state graph: vertices denote specific file system states, and edges indicate state transitions produced by file system operations. This graph is an abstraction for state-space exploration in model checking. . . . .	4
3.1	MCFS Model checking framework . . . . .	8
5.1	Model checking workflow for different file system types: from left to right, block-based, FUSE, and character-device-based. . . . .	13
6.1	The architecture of VeriFS, including the novel snapshot pool <code>stioctl_CHECKPOINT</code> and <code>ioctl_RESTORE</code> API. The interaction between VeriFS, libFUSE, and OS kernel. . . . .	15
7.1	Speed comparison for different experiments. Unless specified in parentheses, all experiments were run on RAM disks or entirely in memory. . . . .	16
7.2	Speed and memory consumption of the model checker. <code>fsops_rate</code> is the number of file system operations that MCFS performs per second; <code>vmem_size</code> is the size of the virtual memory; and <code>pmem_size</code> is the amount of physical memory consumed. The axes in both figures have the same meanings but use different scales. . . . .	17
7.3	File system operation rate and swap usage in a two-week MCFS experiment on VeriFS1. . .	18

# List of Algorithms

1	Abstraction Functions . . . . .	10
---	---------------------------------	----

## Abstract

Developing and maintaining a file system is time-consuming, typically requiring years of effort. Developers often test compliance with APIs such as POSIX with hand-written regression suites that, alas, examine only a fraction of a file system's state space. Conversely, formal model checking can explore vast state spaces efficiently, increasing confidence in the file system's implementation. Yet model checking is not currently part of file system development. *Our position is that file systems should be designed a priori to facilitate model checking.* To this end, we introduce MCFS, an architecture for efficient and comprehensive file-system model checking. MCFS relies on two new APIs that save and restore a file system's in-memory and on-disk state. We describe our earlier attempts at model-checking file systems, including unsuccessful or inefficient ones. Those attempts led us to develop VeriFS, which implements the new APIs. We illustrate MCFS's model-checking principles with VeriFS, a FUSE-based file system we were able to quickly develop with MCFS's help.

# Chapter 1

## Introduction and Motivation

File system development is time-consuming, complex, and error-prone, requiring precise logic, standards compliance (*e.g.*, POSIX), and careful implementation of data structures and concurrency [24, 40, 73]. Yet even mature file systems suffer repeated bugs. For example, Btrfs [54] was designed 14 years ago yet reported 110 bugs solely in 2020 [35]. Such bugs can cause data corruption or loss and system crashes [36, 40, 55].

For instance, Ext4 [18] encountered data loss upon power off [5], which caused a configuration file from KTimeTracker (a time tracking application) to be replaced by an empty zero-byte version. As another example, Ext4 corrupted a “qcow2” image file [42] that was running a Linux KVM [37] guest when Ext4 enabled encryption functions [33]. A “qcow2” file is a disk image for providing persistent storage for QEMU [4] virtual machines, which stands for the second version of the QEMU Copy On Write format. XFS [68], a well-tested file system created in 1993, had exposed severe bugs that led to kernel panic [32] and metadata corruption [34]. Likewise, Btrfs [50] suffered from data loss [65], kernel panics [38], and hanging [64] due to constantly emerging file system bugs.

Moreover, file systems need to incorporate new features and adapt the most recent kernel, which inevitably brings in defects and bugs [40]. Therefore, file system development is an everlasting effort. Although file system users could trigger some hidden bugs and report them to developers [5, 32, 38], this process is cumbersome because developers can only fix a few bugs at a time. The remaining undisclosed bugs can exist for years, endanger system reliability, and be susceptible to malware. Thus, the bugs should be detected in advance to prevent users from triggering them, which is ideal for file system development [40].

File system bug detection is an active research topic, which assists developers in identifying bugs and thereby supporting file system development [40]. Bug-detection techniques for file systems can be grouped into regression testing, verified or machine-verifiable file systems, model checking, and fuzzing. Generally, regression testing [44, 57, 62] only tests fixed known cases of a file system, so it is impractical to cover enough corner cases that are barely encountered and more likely to hide bugs. Machine-verifiable file systems are formally and rigorously verified for some aspects (*e.g.*, crash consistency [10–12, 58]). Model-checking validates if the model of a file system meets the specification. Still, existing model checking work is either restricted to specific bugs or requires an abstract file system model to perform model checking [71, 72]. Fuzzing is a recent technique to find file system and kernel bugs. However, existing file system fuzzing can only uncover a limited type of bugs [70] or require external checkers to identify bugs [36], making existing fuzzing less convenient.

To help address these issues, we leverage model checking to explore bounded file system state thoroughly and apply behavioral comparisons to recognize potential bugs. Thus, we present MCFS, a model-checking framework for file systems that: **(1)** checks every corner case in a bounded state space; **(2)** does not require an abstract file system model; **(3)** retains the original file system behavior; **(4)** applies to most POSIX-compliant file systems, whether in-kernel or user-space; and **(5)** has high performance with the help of fast devices like memory.

MCFS compares file systems to each other with concurrent processes that non-deterministically issue file-system calls; it compares their outcomes (*e.g.*, file content and return values) to find discrepancies. MCFS can exhaustively execute system-call permutations and thus uncover abnormal behavior.

We implemented MCFS using the Spin model checker [59] and applied it to a number of file systems. Our position is that thanks to its exhaustive state-space analysis capabilities, model checking should be an integral part of the file-system development process (alongside traditional hand-written regression suites [44, 57]). However, file systems need to facilitate model checking. We make the following contributions:

1. We created MCFS to support file system checking that can exhaustively search bounded state spaces.
2. We uncovered two inherent challenges to model-checking: cache incoherency and I/O inefficiencies. Ultimately we designed state checkpoint/restore APIs to ease integrating file systems with model checkers.
3. We developed (two versions of) a FUSE file system, *VeriFS*, that efficiently checkpoints and restores file-system states using our proposed APIs.
4. We empirically evaluated MCFS's performance. Our results suggest that MCFS is a viable approach: one can use it to find behavioral deviations and bugs while developing or maintaining a file system.

The rest of this report is organized as follows. Chapter 2 describes the background knowledge of model checking and file system bugs. Chapter 3 presents the design goals and architecture of the MCFS model checking framework. Chapter 4 describes four main challenges we encountered in the process of implementing and executing MCFS. Chapter 5 provides more implementation details regarding MCFS. Chapter 6 illustrates different approaches to track full file system states in MCFS. Chapter 7 describes the evaluation results of MCFS to model-check file systems and the efficiency of *VeriFS*. Finally, we conclude our work in Chapter 9 and discussed the future work for MCFS in Chapter 10.

# Chapter 2

## Background

This section presents the background on model checking and file system bugs, including generic model-checking concepts, the Spin model checker, and types and patterns of file system bugs.

### 2.1 Model Checking

#### 2.1.1 Model Checking Concepts

Model checking is an automatic method for verifying finite-state concurrent systems. It performs an exhaustive search of a bounded state space to verify whether a system’s model conforms to its specification. The model-checking process typically consists of three phases [3]: modeling phase, running phase, and analysis phase.

- The modeling phase generates a model description for the system under testing by a model checker and formalizes the properties from the system’s specification.
- The running phase checks whether properties from the modeling phase are valid by running the model checker.
- The analysis phase performs different operations based on the outcome to confirm if the system satisfied the property. The model checking process will check the next property when the property is satisfied. On the other hand, if the system violates a property, the model-checking process will analyze this instance and then refine the model, design, or property [14].

The model in model checking incorporates states and transitions between states, which constitutes a directed state graph [14]. Therefore, to model-check a system, the model checking process needs to conduct exhaustive state-space exploration in the state graph to verify every property belonging to the system under consideration. The state-space of a system can grow exponentially and become too large to explore. It causes the model checker to fail to finish state-space exploration within a certain time and potentially result in an “out of the memory” condition [3], called the “state explosion” in model checking. In that case, the model size should be reduced using various techniques [13, 14].

The above demonstrates the generic model checking workflow and concepts. In MCFS, we adopted the state-space exploration of the model checking to check file system states but did not need to build a model because the checking process relies on comparing file system behaviors among multiple file system implementations. Figure 2.1 shows an example of the state graph in MCFS.

In the state graph, a vertex represents a specific file system state that describes various file system properties, including files, directories, metadata, etc. The starting point of the state space is a clean file system



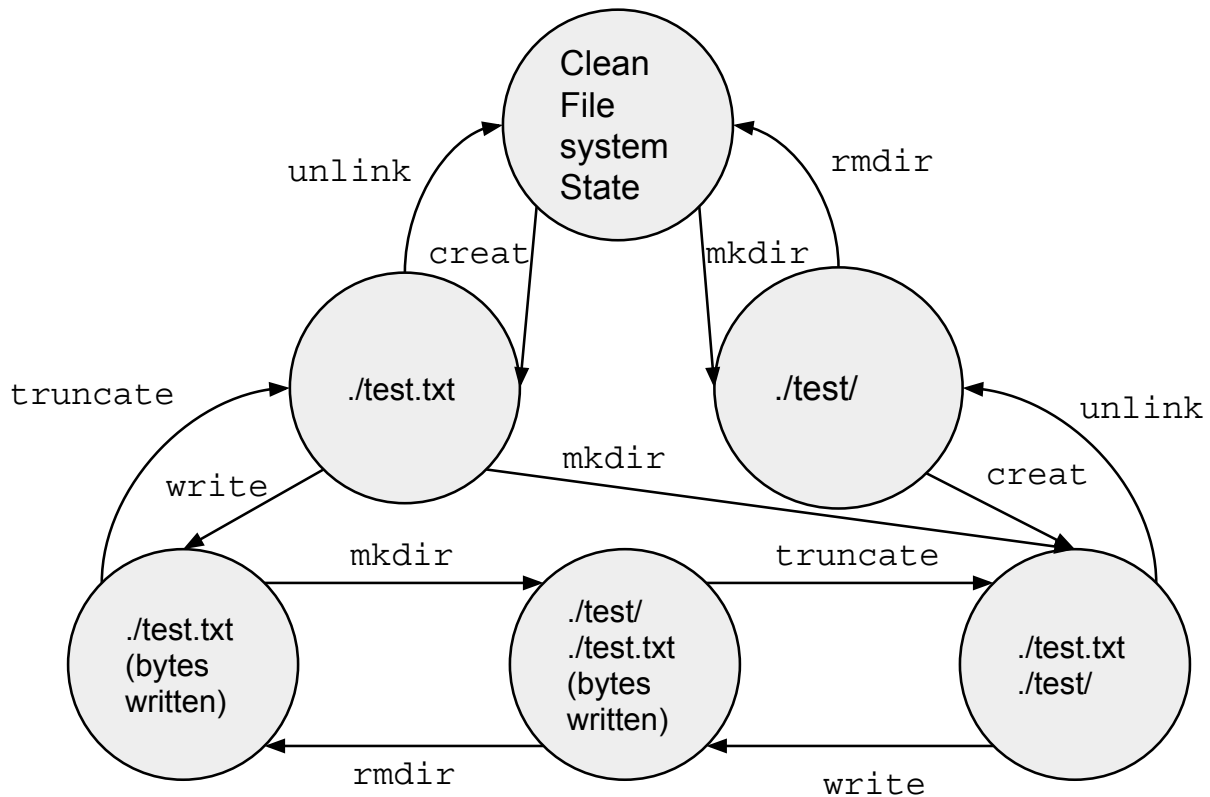


Figure 2.1: An example of a file system state graph: vertices denote specific file system states, and edges indicate state transitions produced by file system operations. This graph is an abstraction for state-space exploration in model checking.

state obtained from mounting an empty file system. The file system operations (*i.e.*, system calls) lead to the state transition (*i.e.*, edge in the graph) in the state-space exploration because the file system operation can potentially change the file system state. In practice, we also employed read-only file system operations that do not alter a file system state because we apply read-only operations to examine if the file system reaches a correct state and returns precise data and status.

### 2.1.2 Spin Model Checker

We require a model checker to drive the checking process for MCFS and identify which state to explore. The model checker must be efficient and able to traverse file system states thoroughly. We applied the well-known Spin model checker [25] to the MCFS framework as the model checking driver. Spin is a model checker that can verify multi-threaded software efficiently and exhaustively [17]. Spin supports the embedded C code as part of model specifications, which helps Spin integrate with the file system testing because most Linux file systems are implemented in C code [40]. Spin also works with a high-level description language called PROMELA to express the system model that Spin can analyze [25].

Spin supports on-the-fly verification, which prevents the need to construct a global state graph in advance. This feature is handy for file system model checking because the state graph is unpredictably large and cannot be predefined readily under hundreds of system calls and a vast parametric space [61]. In addition, Spin supports swarm verification [27, 28] to perform many parallel verification jobs using a multi-core CPU or distributed remote machines for solving one significant verification problem. Therefore, Spin provides a feasible method to tackle large state-space exploration in the model checking process, considering that a

modest number of file system operations can increase file system states exponentially.

## 2.2 File System Bugs

File systems are fundamental components in the operating systems (OSs), which generally have a massive codebase [40]. Typically, a sophisticated OS can work with diverse file systems to achieve various functionalities and features [52]. Therefore, file systems not only have different designs and characteristics but also continuously evolve to incorporate new features and improve performance and reliability [40]. In the process of file system development, it is inevitable to produce bugs and defects due to the complexity of a file system and unavoidable human mistakes. Thus, one of the critical tasks of file system development is to detect and fix emerging file system bugs and defects.

Based on bugs' effects, file system bugs can be classified into semantic bugs and memory-safety bugs [36]. Usually, semantic bugs are triggered silently and unlikely to make a strong signal such as kernel crash or panic that indicates a bug. Instead, memory safety bugs are more likely to have serious impact and thus can be observed more easily. File system bugs can also be classified according to their causes. In this way, file system bugs can be divided into the following categories that we describe next: concurrency bugs, crash inconsistency bugs, logic bugs, memory bugs, specification violation bugs, and other types of bugs [36].

**Concurrency Bugs.** Concurrency bugs in file systems are the bugs that are triggered under multiple interleaving threads, which come in diversified patterns, including atomicity violations, deadlocks, order violations, missed unlocks, double unlocks, and wrong locks [40]. Detecting concurrency bugs often requires multiple interleaving threads to access and test the same memory location [20, 21]. We have not applied multi-threading in MCFS yet for simplicity, but we believe MCFS can extend to multi-threading as Spin naturally supports the formal verification of multi-threaded applications [25, 29].

**Crash Consistency Bugs.** File system operations mainly modify in-memory state instead of persistent state for performance consideration. When a crash or power loss happens, crash consistency bugs make file systems unable to persist data and metadata in the file system correctly. For example, Btrfs has a field in the inode data structure to check if the inode needs to be persisted to the storage device through `fsync`. One crash-consistency bug discovered in 2015 [19] prevents Btrfs from updating this inode field when punching a file hole. Thus, the `fsync` call in Btrfs performed no operation when Btrfs crashed, resulting in data loss [45]. Unlike other types of bugs, crash consistency bugs only fail to persist data when a crash occurs [51]. MCFS can identify crash consistency bugs by simulating crash conditions or using corrupt file system images, which we leave as future work.

**Logic Bugs.** Logic bugs are specific to file system implementations because they often originate from misuse of designs and algorithms, inaccurate assumptions, and incorrect implementations [36, 40]. Developers generally add comprehensive assertions into the code to avoid logic bugs beforehand, but this technique is rarely enabled and used in real-world file system checkers [36]. MCFS can capture logic bugs with behavioral comparison. For example, Ext3 [63] used a wrong algorithm in function `find_group_other()` to find a block group for inode allocation, in which the linear search algorithm does not check all candidate groups and neglects the groups at the beginning [40]. MCFS can potentially detect this logic bug by running multiple file systems and performing operations that need to allocate inodes. The file system with this bug possibly returns an `ENOSPC` error even if the file system has free inodes, but other file systems with the correct algorithm can allocate an inode successfully without error returned. MCFS compares return value and error code for all the file systems under consideration after each file system operation. Therefore, MCFS can catch this discrepancy and report a potential bug.

**Memory Bugs.** Because file systems regularly involve memory operations (*e.g.*, initialize, allocate, free, access memory), the memory bugs are common in file system development and can induce severe issues like memory leaks, uninitialized reads, out-of-bound accesses, dangling pointers, double frees, and buffer overflow [36, 40]. The *integrity checks* in MCFS monitor the return status and file system state after each operation to test whether memory is correctly handled.

**Specification Violation Bugs.** Specifications (*e.g.*, POSIX [30]) tell developers how to implement a file system with specific behavior such as error codes, path resolution, invariants, etc. [53, 55]. Accordingly, file systems must conform to the specifications for robustness, reliability, and compatibility. Specification violation bugs expose serious security issues and make file systems more vulnerable to attackers and malware. MCFS can check for such specification violation bugs without file system knowledge. As an illustration, if one file system uses an incorrect error code not specified in POSIX, MCFS can capture this discrepancy by comparing it to a correctly-implemented file system. We assume that most file systems implement a particular file system operation accurately, so it is rare for both file systems to implement the same function incorrectly [40].

## Chapter 3

# MCFS Design

We exploit this efficient search technology to explore file system states exhaustively. The MCFS model-checking framework is designed to detect *discrepant behaviors* in file systems; i.e., situations where two file systems behave differently given the same inputs.

### 3.1 Design Goals

MCFS has five key design goals: **(1) Thorough coverage:** It should thoroughly explore the state spaces of the file systems we check, so that it can uncover as many corner cases as possible. This requires MCFS to check as many permutations of file system operations as possible, exploring all possible changes that the given set of operations can make. **(2) Eliminate need for an abstract model:** Traditional model checking requires one to define a model for the system under investigation. MCFS can verify file systems without a model because it directly executes system code to perform state-space exploration. **(3) Absence of observer effects:** To ensure accuracy, MCFS should avoid changing the behavior of the investigated file system, ideally treating it as a black box. If we have to modify it, we must be careful to minimize MCFS's footprint. **(4) Universality:** The model-checking framework should support a wide range of file systems (*e.g.*, in-kernel, user-level, networked, distributed). **(5) High performance:** Since a file system's state space is the product of many system calls and their parameters, the number of states can be exponential. The model checker must be able to enumerate new states as fast as possible, exploring large fractions of the state space within a reasonable amount of time.

### 3.2 MCFS Architecture

Figure 3.1 depicts the MCFS model checking framework. MCFS's architecture has four main components: randomized test engines, optimized state-space exploration, integrity checks, and abstraction functions. Driven by the Spin model checker [25], *randomized test engines* nondeterministically issue operation sequences to each file system under consideration. Spin's efficient partial-order reduction algorithm [29] allows MCFS to execute all permutations of the given set of calls and their parameters without duplication. Its randomized driver processes generate both valid and invalid call sequences. Valid sequences should succeed on all file systems, while invalid ones (*e.g.*, `write()` before `open()`) should produce consistent error behavior. Invalid sequences are critical because they exercise error paths, where bugs often lurk.

We chose Spin to perform *optimized state-space exploration* because: (i) it is open-source and actively maintained; (ii) C code can be embedded in Spin's model-description language (Promela) to issue system calls, and `c_track` statements let Spin record C buffers as states [26]; and (iii) Spin can perform parallel

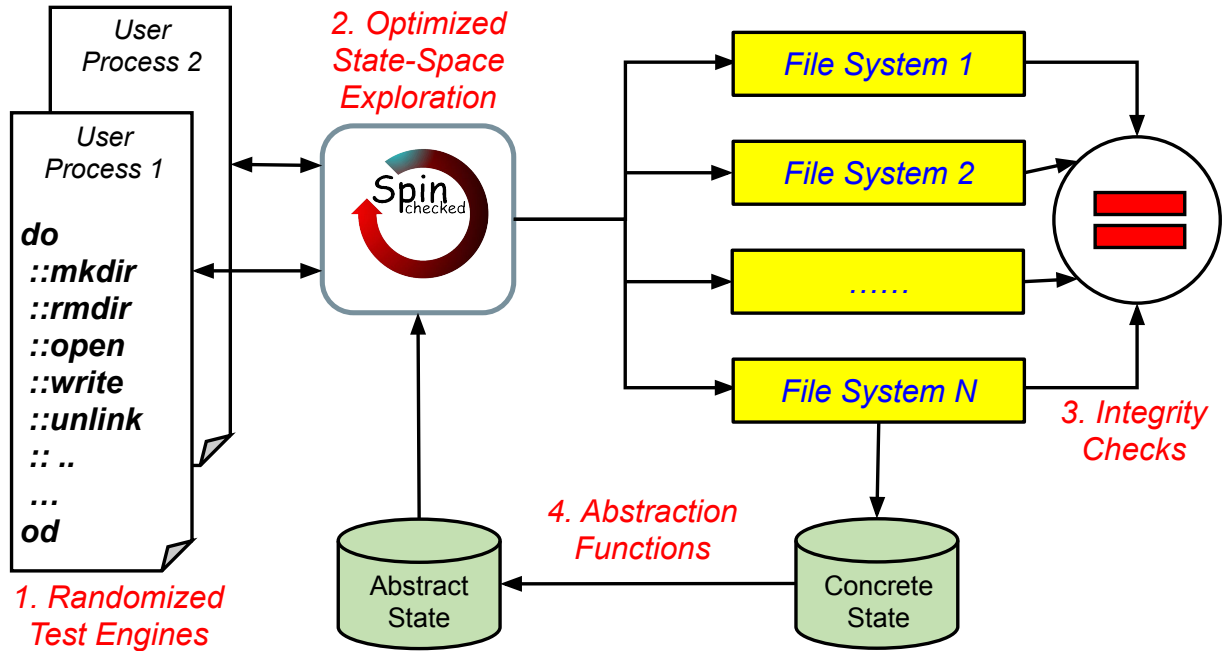


Figure 3.1: MCFS Model checking framework

model checking using Swarm verification techniques [27], substantially speeding up exploration of large state spaces.

After each system call, *integrity checks* verify that all tested file systems have identical states by asserting equality of return values, error codes, file data, and metadata. If a discrepancy is detected, the integrity checker reports a potential bug and halts. Spin logs the precise sequence of operations, parameters, and starting and ending states that led to a problem, simplifying reproducibility. Because file systems have implementation-specific features [53], not all discrepancies are bugs. We believe though that many of these non-bug behavior differences are also interesting, since they still affect application behavior [51].

Finally, *abstraction functions* convert concrete states into abstract ones. In MCFS, each concrete state contains all the information needed to describe the file systems, including file data and metadata. MCFS uses the abstract state to determine whether a state was previously visited. If MCFS reaches a state that is logically equivalent to a previous one, it will backtrack corresponding concrete states to restore the file systems to their earlier versions. The abstraction functions hash the important data in the concrete states (including file paths, data, and relevant metadata) to distinguish *logically* unique states but omit noisy attributes such as `atime` timestamps and the physical locations of data blocks. Hashing the entire on-disk state would fail because of those less interesting attributes.

Designing abstraction functions requires domain knowledge. We wrote the functions to conform to the POSIX specification, so they are generic and applicable to most POSIX-compliant file systems.

# Chapter 4

## Challenges

We now describe challenges we encountered while developing MCFS, and our attempts to work around them.

### 4.1 Access to In-Memory States

MCFS must save and restore all information related to the tested file systems, including their persistent (on-disk) and dynamic (in-memory) states. Spin could use the underlying block device to track persistent states, but there is no simple way to access in-memory states—in kernel or user space—because they are not part of Spin itself.

**In-kernel file systems.** Many file systems (*e.g.*, Ext4 [18], XFS [57]) run in the Linux kernel, with states in the kernel address space. In theory, one could use `/dev/kmem` to save and restore those states, but in reality they are so intertwined with other kernel data structures that doing so is impractical.

**User-space file systems.** Tracking kernel file systems is hard, so we turned to file systems built on lib-FUSE [60] (*e.g.*, fuse-ext2 [1]). Such file systems, however, are separate processes, so Spin cannot directly track their internal state. We tried several alternatives. First, we modified `malloc` to allocate memory from a shared-memory pool accessible to Spin, so that it could be saved and restored. This failed because important state was stored in static (non-heap) variables outside the shared-memory segment, leading to incorrect pointers and crashes. We concluded that it was impractical to modify file-system code to avoid these static variables.

### 4.2 Cache Incoherency

We next explored a compromise in which Spin tracked only the persistent (on-disk) state. Doing so allowed MCFS to run without crashing, but our experiments encountered corrupted file systems. A typical symptom was directory entries with corrupted or zeroed inodes, caused by Spin backtracking and restoring a persistent state. Since we were not restoring in-memory state to match, cached information in the kernel was no longer consistent with the disk content. For example, the `dcache` might contain a recently created directory, but the restored state might reflect a time before its creation.

We tried to resolve the inconsistency by calling `fsync` after each operation, and by mounting the file system with the `sync` option. Neither approach was effective: although they guaranteed that the caches were flushed to persistent storage, they did not implement the opposite operation—loading any Spin-initiated change in the persistent storage back into the in-memory caches.

Finally, we compromised again: we unmounted and remounted the file system between each pair of operations. An unmount is the *only* way to fully guarantee that no state remains in kernel memory. [Re]mounting always loads the latest state from disk, ensuring that the caches are coherent between each Spin state exploration. This compromise caused two problems: (1) it considerably slowed state exploration (see Section 7) and (2) it prevented us from identifying file-system bugs caused by incorrect in-memory states. These problems led us to consider adding support for file system state checkpointing and restoring (see Section 6).

### 4.3 State Explosion

---

#### Algorithm 1: Abstraction Functions

---

**Input** : Path to the file system mount point *path*  
**Output**: 128-bit MD5 Hash

```

1 files ← list ([])
2 md5ctx ← md5_init ()
  // Recursively walk the mount-point directory
3 foreach file in recursively_traverse_dir (path) do
4   | files.append (file)
5 sort (files) // Sort files by pathnames
6 foreach file in files do
7   | fd ← open (file.path)
8   | content ← read (fd) // Read all file content
9   | md5_update (md5ctx, content)
10  | close (fd)
11  | attrs ← stat (file)
  // Get important metadata
12  | attrs' ← important_attributes (attrs)
13  | md5_update (attrs')
14  | md5_update (file.path)
15 return get_md5_hash (md5ctx)

```

---

We use Spin’s `c_track` statement to declare memory buffers used by the C code. During state exploration, Spin detects an already-visited state by comparing these buffers against all previously visited states. This causes state explosion because *any* change in a buffer is considered a new state, yet some changes, such as access-time updates, are rarely relevant to bugs. Consequently, Spin could not fully explore file systems with even moderate parameter spaces. Fortunately, Spin’s `c_track` allows one to define abstract states used for matching, and concrete states used only in state restoration. We thus introduced an abstract state that contained an MD5 hash of file paths, data, and important metadata (*e.g.*, mode, size, nlink, UID, and GID) for all files and directories.

Algorithm 1 shows the procedure to obtain an abstract state of a file system. We first identify all files and directories in the file system by traversing from the mount point. We then sort them by their pathnames so that they appear in a consistent order. We then `read` each file’s contents and call `stat` to obtain its metadata. The `important_attributes` function extracts only the important metadata mentioned above. Finally, we calculate the MD5 hash for the files’ content, important metadata, and pathnames.

We then marked persistent states as concrete, and instructed Spin to track only the hashes that distinguish different abstract states. Doing so not only prevented visiting duplicate states, it also greatly reduced the amount of memory needed to track states, increasing Spin’s exploration capacity.

## 4.4 False Positives

MCFS performs integrity checks that assert state equality of the file systems under investigation. In case of any discrepancy, MCFS terminates and reports a bug, logging the operations it executed and their parameters. We found, however, that MCFS sometimes terminated on discrepancies that were not bugs. We took measures to prevent these false positives and describe several such cases below.

**Directory-size reporting and ordering.** In Ext4, directory sizes are a multiple of the block size; other file systems report sizes based on the number of active entries. Thus, directories on two different file systems might have the same contents but different sizes; we thus ignore directory sizes. Similarly, file systems return directory entries in different orders, so we sort the output of `getdents` before comparing them.

**Special folders.** Some file systems create special folders: For example, Ext4 has a `lost+found` folder to save lost and damaged files, while XFS does not. This caused namespace discrepancies between file systems. We added an exception list of special files and directories; MCFS ignores anything on this list when comparing abstract states.

**Differing data capacity.** Although we tested all file systems on block devices of the same size, they exposed different usable data capacities. This is a problem when the file systems are nearly full: calling `write` can succeed on one file system and fail on another, reporting a false bug. We thus equalize free space among file systems being checked: when MCFS starts, it queries all file systems and records the smallest free space as  $S_L$ ; then on each file system with free space  $S_n$ , it creates a dummy file and writes  $S_n - S_L$  bytes of zeros.

None of these workarounds introduce false negatives, because they are all dealing with unstandardized behavior. For example, an application should not expect sorted output from `getdents`, so if a given file system suddenly stops sorting, that is not a bug. (However, if the change introduces other misbehavior, we will detect the consequences.)



## Chapter 5

# MCFS Prototype Implementation

Figure 5.1 shows how MCFS’s prototype handles different types of file systems. The file system syscall engine, written in Promela with embedded C code, consists of a multi-entry `do ... od` nondeterministic loop; each entry contains code to issue file-system operations, perform integrity checks, and record logs. Using Spin, MCFS nondeterministically selects an operation and its parameters, then executes it on all tested file systems. MCFS `mmaps` the file systems’ backend storage devices into Spin’s address space so it can track their states.

To prevent cache-coherency problems, MCFS unmounts and remounts kernel file systems (*e.g.*, Ext4 and JFFS2) before and after each operation (see Section 4.2). However, syscalls such as `write` that depend on kernel state (*e.g.*, open file descriptors) cannot be used in isolation. We thus developed meta-operations comprising small sequences of syscalls that avoid tracking kernel state: `create_file` creates and then `closes` a file; `write_file` opens, writes some data to, and `closes` a file. Other operations that can execute alone are run directly by MCFS, *e.g.*, `truncate` and `mkdir`. Each operation’s parameters are selected nondeterministically from a pre-defined (bounded) parameter pool. Because we limited our exploration space to fixed syscalls and parameters, the entire exploration—while large—is guaranteed to be bounded.

For FUSE-based file systems (*e.g.*, `fuse-ext2`), the syscall is issued to the OS. FUSE’s normal behavior results in several user/kernel messages being passed, coordinated via `/dev/fuse` (see `fuse-ext2` in Figure 5.1).

To avoid being slowed by I/Os, we used RAM block devices as backend storage for block-based file systems (*e.g.*, Ext4 and XFS). Linux’s RAM block device driver (`brd`) requires all RAM disks to be the same size; we slightly modified it (renamed `brd2`), to allow different-sized RAM disks for file systems with different minimum-size requirements.

Some file systems must be mounted using special devices. For example, JFFS2 [67] requires an MTD character device [66] instead of a regular block device. MCFS sets up JFFS2 differently: it (1) loads the `mtddram` kernel module, which creates a virtual MTD device in RAM, and then (2) loads the `mtdblock` module to provide a block interface for the virtual MTD device. This approach allows Spin to `mmap` the MTD storage via the block device.

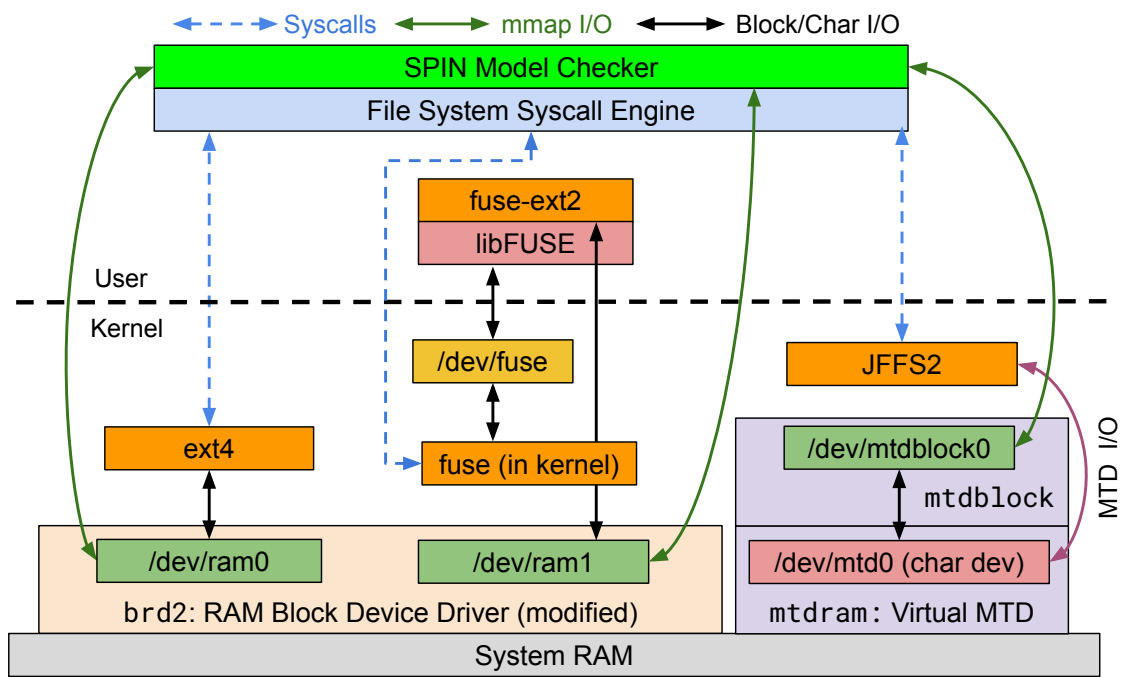


Figure 5.1: Model checking workflow for different file system types: from left to right, block-based, FUSE, and character-device-based.

## Chapter 6

# Tracking File-System States

Unmounting and remounting a file system after each operation solved cache-incoherency problems but slowed the model checking and deviated from a normal use case. Due to its backtracking search process, Spin saves and restores all information related to the tested file systems. Therefore, we must track all file-system states, including persistent and in-memory ones. We investigated three approaches: (1) process snapshotting and (2) VM snapshotting, which culminated in our (3) MCFS-enabled VeriFS.

**Process snapshotting.** User-space file systems run as independent processes. To keep their in-memory states coherent with their persistent states, we can snapshot them using existing tools. We explored a popular snapshotting tool called CRIU [15] and tried integrating it with MCFS. Unfortunately, CRIU refused to checkpoint processes that have opened or mapped any character or block device (with a few unhelpful exceptions). Since FUSE-based file systems use the character device `/dev/fuse` to communicate with the kernel, CRIU could not checkpoint them. However, CRIU was able to snapshot the user-space NFS server Ganesha [49]; we are investigating model-checking Ganesha with CRIU.

**Virtual-machine snapshotting.** Hypervisors can snapshot and restore an entire VM, including full file-system states enclosed therein. However, VM-level snapshotting is fairly slow and heavyweight. For example, LightVM claims that it takes 30ms to checkpoint a trivial unikernel VM and 20ms to restore it [41]. This may be fast enough for cloud platforms but is too slow for MCFS; LightVM’s latency limited our model-checking rate to only 20–30 operations/s.

**VeriFS.** It is difficult for MCFS to track the in-memory state of file systems (see Section 4.1). But if a file system *itself* could checkpoint and restore its state, MCFS could use that facility to easily perform state capture and restoration, avoiding cache incoherency. To demonstrate this idea, we developed a RAM-based FUSE file system, *VeriFS*. Figure 6.1 shows the architecture of VeriFS. Apart from standard POSIX operations such as `open`, `write`, and `close`, VeriFS provides checkpoint and restore APIs via `ioctl`s: `ioctl_CHECKPOINT` and `ioctl_RESTORE`.

When MCFS calls `ioctl_CHECKPOINT` with a 64-bit key, VeriFS locks itself, copies inode and file data into a *snapshot pool* under that key, and releases the lock.

Similarly, `ioctl_RESTORE` causes VeriFS to query the snapshot pool for the given key. If it is found, VeriFS locks the file system, restores its full state, notifies the kernel to invalidate caches, unlocks the file system, and discards the snapshot.

VeriFS is intended to demonstrate the idea of having file systems *themselves* support model checking by providing checkpoint and restore APIs. Therefore, the initial version, VeriFS1, was fairly simple. It used a fixed-length inode array with a contiguous memory buffer attached to each inode as the file data. It had

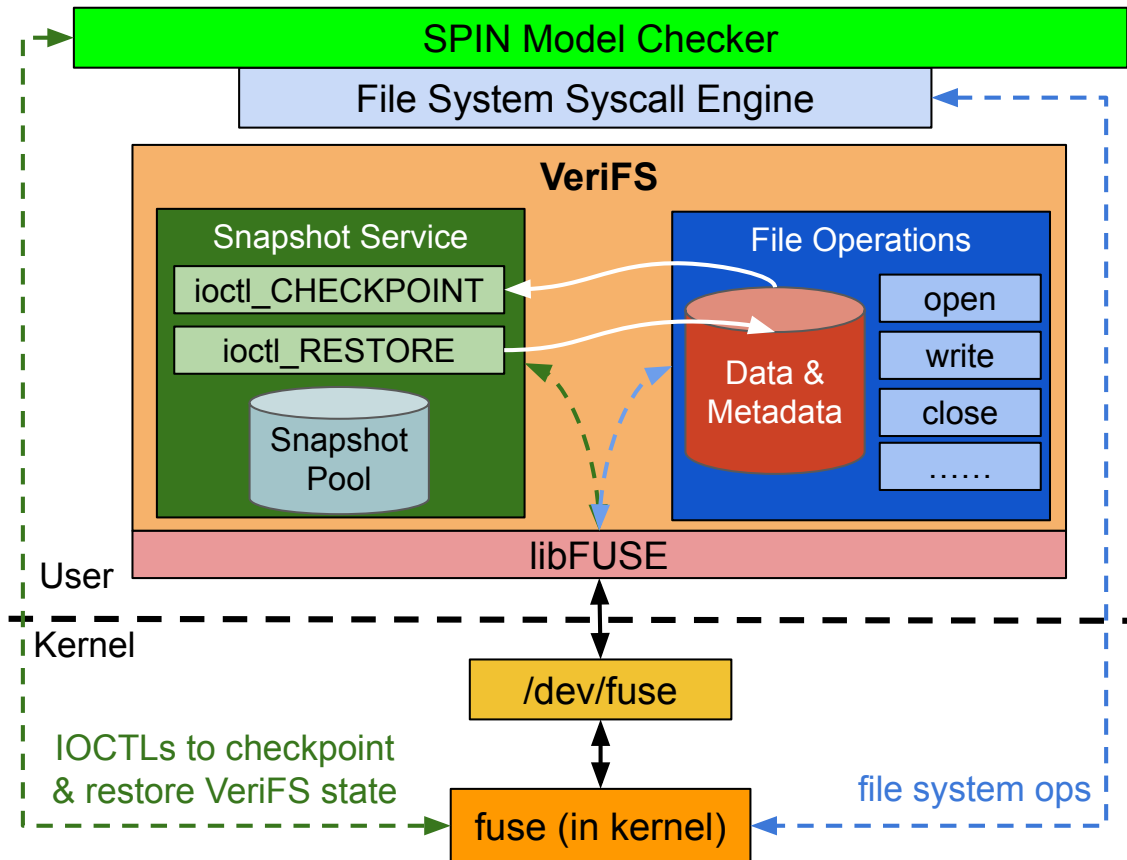


Figure 6.1: The architecture of VeriFS, including the novel snapshot pool `ioctl_CHECKPOINT` and `ioctl_RESTORE` API. The interaction between VeriFS, libFUSE, and OS kernel.

only a limited set of file system operations and lacked support for `access()`, `rename()`, symbolic and hard links, and extended attributes. It also did not limit the amount of data that could be stored.

We ran MCFS with Ext4 and VeriFS1 for over 5 days; MCFS executed over 159 million syscalls without any errors, behavioral discrepancies, or file system corruption. To demonstrate how MCFS supports file system development, we next developed VeriFS2 to add missing features. During development, we used MCFS to model-check VeriFS1 against VeriFS2 to find and fix bugs in VeriFS2. MCFS helped us find several bugs in VeriFS2, which we discuss further in Section 7.

**In sum,** model checking a file system can involve exploring a vast number of states. If state exploration takes too long (*e.g.*, is I/O-bound), then the entire model-checking process becomes impractical. Our position is that speedy and thorough file system model-checking requires the checkpoint/restore API we propose.

# Chapter 7

## Evaluation

We experimented with an MCFS prototype and a number of file systems running on a VM with 16 cores, 64GB RAM, and 128GB of swap space allocated on a local hypervisor SSD. We present preliminary performance results and discuss how MCFS helped our file system development.

### 7.1 Performance and Memory Demands

We ran MCFS and recorded key performance metrics for the following file system combinations: Ext2 vs. Ext4, Ext4 vs. XFS, Ext4 vs. JFFS2, and VeriFS1 vs. VeriFS2. We used 256KB RAM block devices for both Ext2 and Ext4, and 16MB for XFS, which allows a larger minimum file-system size. VeriFS is an in-memory file system and does not need a block device.

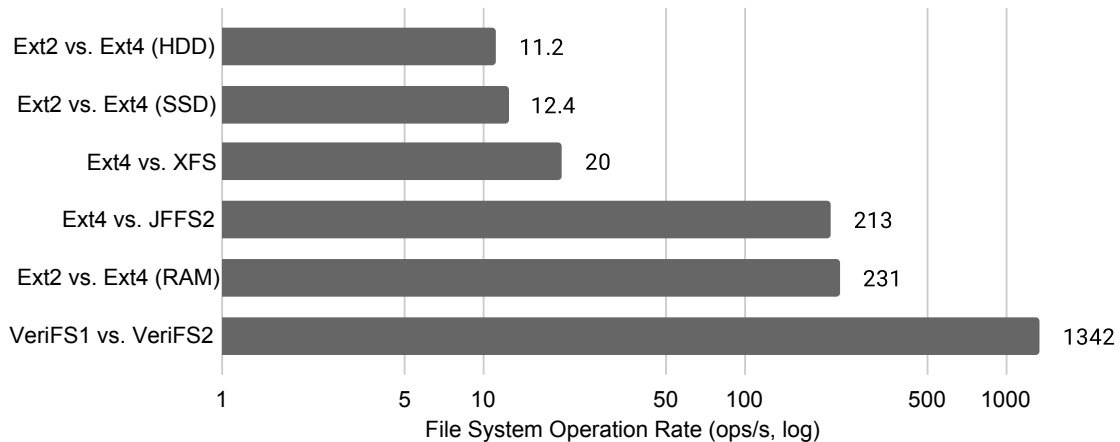


Figure 7.1: Speed comparison for different experiments. Unless specified in parentheses, all experiments were run on RAM disks or entirely in memory.

Figure 7.1 shows model-checking speeds observed in our experiments. Checking Ext4 vs. XFS on RAM disks was over  $11\times$  slower than Ext2 vs. Ext4 because MCFS consumed 105GB of swap space for the former, so swap time dominated. Checking Ext2 vs. Ext4 on HDD was  $20\times$  slower than on RAM disks; using SSD was still  $18\times$  slower. This illustrates the advantage of using RAM as backend storage. Checking VeriFS1 vs. VeriFS2 was  $5.8\times$  faster than Ext2 vs. Ext4 for two reasons: (i) MCFS used the checkpoint/restore APIs (see Section 6) and thus did not have to unmount and remount the VeriFS file systems; (ii) VeriFS runs entirely in-memory, so MCFS did not access block devices to checkpoint and restore persistent states.

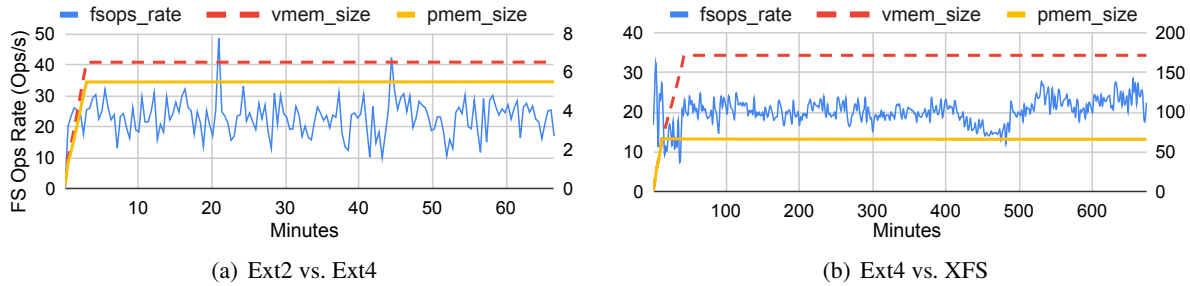


Figure 7.2: Speed and memory consumption of the model checker. `fsops_rate` is the number of file system operations that MCFS performs per second; `vmem_size` is the size of the virtual memory; and `pmem_size` is the amount of physical memory consumed. The axes in both figures have the same meanings but use different scales.

Figure 7.2 shows the speed of model checking and the memory consumption over time. Figure 7.2(a) shows an experiment that compared operations on Ext2 against those same operations performed on Ext4; Figure 7.2(b) shows experiments on Ext4 vs. XFS. We used 256KiB RAM block devices for both Ext2 and Ext4, and 16MiB for XFS, which imposes a larger minimum file system size.

When checking Ext2 vs. Ext4, MCFS consumed 5.5GB of physical memory (`pmem_size`) and its virtual memory (`vmem_size`) size was 6.5GB. When checking Ext4 vs. XFS, MCFS used 66GB of physical and 171GB of virtual memory; because this larger memory footprint exceeded physical RAM, MCFS used 105GB of swap space. The average model-checking speed for Ext2 vs. Ext4 was 231 operations per second (ops/s); for Ext4 vs. XFS it was 20 ops/s. Checking Ext4 vs. XFS was more than  $10\times$  slower than Ext2 vs. Ext4 because the time spent on swap I/O dominated the former experiment.

Note that by default MCFS remounts and unmounts the file systems before and after each operation. To evaluate the impact of that approach, we also measured its performance without the inter-operation remounts; in that case the average speed for Ext2 vs. Ext4 was 38% higher (316 ops/s), and that for Ext4 vs. XFS was 70% higher (34 ops/s).

We also tried MCFS with (1) VeriFS alone and (2) VeriFS vs. Ext4. MCFS does not need to unmount and remount VeriFS between every operation because it takes advantage of VeriFS’s `ioctl` APIs to checkpoint and restore states, but it must still do so when checking Ext4. The average speed of checking VeriFS alone was 713 ops/s, and that of checking VeriFS vs. Ext4 was 390 ops/s. Clearly, the added un/mounting steps slow the model checking rate considerably.

MCFS keeps an in-memory copy of every concrete state that has a unique corresponding abstract state; this explains why the experiment with Ext4 vs. XFS used more memory than the one with Ext4 vs. Ext2: XFS’s concrete state alone was 16MiB, compared to 256–512KiB for Ext2/Ext4.

Figure 7.3 shows MCFS’s speed when checking VeriFS1 over two weeks. MCFS maintained a rate of around 1,500 ops/s in the first 3 days; this rate then dropped drastically and swap usage spiked because Spin was resizing its hash table of visited states. After rebounding, MCFS’s speed gradually decreased over time because the checkpointed states could not fit in memory and it began to consume swap space. Its speed increased again between days 13 and 14 because the RAM hit rate was high (states needed by MCFS happened to be in memory and did not need to be swapped in and out).

Note that by default MCFS remounts and unmounts the file systems before and after each operation. To evaluate the impact of that approach, we also measured MCFS’s performance without the inter-operation remounts. The average speed for Ext2 vs. Ext4 (in RAM disks) was 316 ops/s, 38% faster than that when remounts and unmounts were used; and for Ext4 vs. XFS it was 34 ops/s, which is 70% faster.

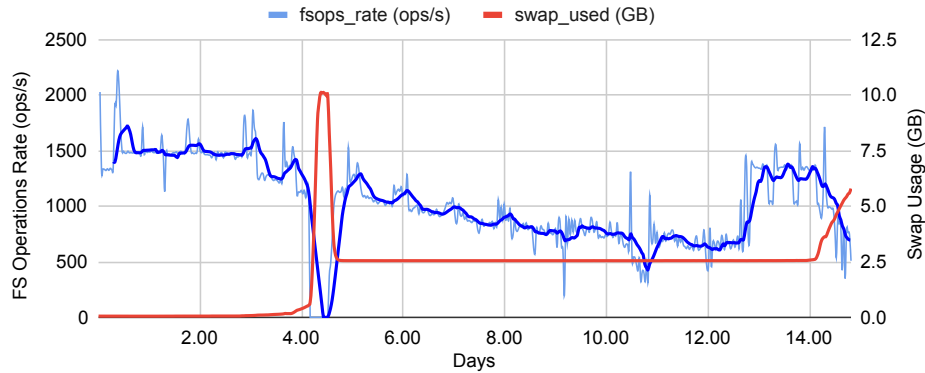


Figure 7.3: File system operation rate and swap usage in a two-week MCFS experiment on VeriFS1.

## 7.2 Assisting File System Development

While developing VeriFS1, we model-checked it vs. Ext4. MCFS found two bugs that we easily fixed, thanks to precise reports of operations and arguments. The first occurred after over 9K operations when test files on VeriFS and Ext4 had different content. The bug arose when `truncate` failed to clear newly allocated space when expanding a file. The second bug was detected after about 12K operations; MCFS created a test directory in Ext4 but VeriFS failed, claiming that the directory existed—but in fact it did not. This was due to cache incoherency between the kernel and VeriFS’s in-memory state. When VeriFS rolled back to an earlier state, the kernel’s inode and dentry caches did not keep up. The fix was to call FUSE’s cache-invalidation APIs (`fuse_lowlevel_notify_inval_entry` and `fuse_lowlevel_notify_inval_inode`).

We then developed VeriFS2 (see Section 6) and used MCFS to assist development by model-checking it vs. VeriFS1. MCFS found two more bugs during this phase. The first occurred after over 900K operations, when the test files in VeriFS1 and VeriFS2 had different data. VeriFS2 had failed to zero out the file’s buffer if `write` created a hole in the file. The second bug was detected after over 1.2M operations: the test file in VeriFS2 was shorter than that in VeriFS1. The reason was that the `write` method in VeriFS2 updated the file size only when the file was expanded beyond its buffer capacity, rather than whenever the file was appended to.

# Chapter 8

## Related Work

### 8.1 Regression Test Suites

Regression tools and suites (*e.g.*, `fsck` [62], `xfstests` [57] and `ltp` [44]) are useful in development and maintenance, but test only known defects and are unlikely to cover all corner cases. Even though these regression suites can incorporate new test cases, they cannot achieve systematic file system testing because of their hand-written nature [36]. Thus, there is no guarantee that regression suites would cover adequate test cases of a file system. Nevertheless, our MCFS framework can be a complementary tool to work with regression suites to cover sufficient state space and facilitate file system development.

### 8.2 Verified or Machine-Checkable File Systems

Prior work has demonstrated the benefits of building a machine-verifiable file system from scratch [2, 9–12, 31, 58, 74], but this approach does not apply to existing file systems. For instance, FSCQ [11, 12] applies extended Hoare logic to write the crash-safe specification. Then, based on the specification, developers can implement a new file system to meet this specification to avoid crash-consistency bugs.

Yggdrasil [58] is a toolkit for writing file systems with automated push-button verification and bug detection using a novel definition of correctness, called crash refinement. Yggdrasil first formulates file system verification as a satisfiability modulo theories (SMT) problem and then fully automates the proof process by solving the SMT problem via Z3 [16], a state-of-the-art SMT solver, relieving programmers from the proof burden.

DFSCQ [10] introduces metadata-prefix specification to specify the properties of `fsync` and `fdatasync`, which describes the possible states after a crash. The DFSCQ specification helps developers verify the crash safety of their applications related to file system synchronization and avoid application-level bugs. The advantage of DFSCQ is the high performance compared to other existing machine-checked file systems [11, 58].

SFSCQ [31] provides machine-checked security proof of confidentiality and utilizes data non-interference to capture discretionary access control. Consequently, SFSCQ can preclude confidentiality bugs in the file systems by implementing the confidentiality specification. In SFSCQ, the confidentiality specification is obtained from the file system code via the idea of sealed blocks.

However, this methodology runs into two significant problems. First, the specification of machine-checkable file systems can only be used to verify a particular dimension of a file system (*e.g.*, crash consistency [10, 11, 58]). Therefore, large parts of the file system are still unverified and thereby may contain bugs. Second, machine-checkable file systems are practical only when developing a new file system from scratch. It is unrealistic to adapt the specification to existing file systems such as Ext4 [18] and XFS [68] and detect real-world file system bugs, limiting the usage to support file system development. Worse, even formally



verified, the FSCQ file system still exposed crash-consistency bugs [36], which shows that even verified file systems are not ever bug-free.

### 8.3 Model Checking and Verification

File-system model checkers have verified test cases [6, 45, 71, 72] and located corner-case bugs, but are strictly focused on crash consistency.

FiSC [72] combines the CMC model checker [46–48] and a file system test driver to continually generate new file system states by executing systematic state transitions (*i.e.*, file system operations). In addition, FiSC exploits a permutation checker and `fscck` to check for failures during the model-checking process recursively. As the enhancement of FiSC, eXplode [71] performs implementation-level model checking, which does not require intrusive changes to the file systems and Linux kernel. However, both FiSC and eXplode tend only to detect crash consistency bugs.

JUXTA [43] automatically infers high-level semantics directly from file system source code by static analysis and exploits cross-checking to detect semantic bugs for multiple file systems in the Linux kernel. JUXTA can detect semantic bugs in existing file systems. However, the design of JUXTA requires substantial file system knowledge and expertise, not to mention different file systems have broadly contrasting semantics. Hence, JUXTA has to compromise on common file system semantics, which restricts the test coverage.

B3 [45], as a black-box crash testing tool, generates bounded workloads on corrupt file system images to simulate the crash circumstance and directly uses `read` and `write` system calls to acquire the file system state. Meantime, the same workloads are run on a clean file system images as an oracle to compare to the state on the corrupt image. Furthermore, B3 uses an automated checker to test whether file system operations upon crash exhibit aberrant behaviors and then correspondingly report a bug. However, B3 is designed to detect only crash-consistency bugs and cannot expose file system bugs that do not involve the crash state. Also, B3 has no guarantee to explore adequate state space of a file system.

Another approach is to build an abstract model and check whether a file system adheres to it [23, 39, 43, 52, 53]. Constructing a formally verifiable model, however, requires considerable domain knowledge and human effort. It is also impractical to build formal models for large and intricate file systems. Worse, the formal model requires updating when the file system’s code changes. CMC [46–48] inserts file system code directly into the model checker, for substantial speed advantages, but requires extensive changes to the *very code being verified*, making the results less trustworthy [71].

Likewise, symbolic-execution tools [7, 8] cannot guarantee thorough coverage because they focus on particular issues (*e.g.*, corrupt input [8]); they also require a behavioral model of each file system call [7].

### 8.4 Fuzzing

Fuzzing can find real-world bugs [22, 36, 56, 69, 70], but either is limited to specific types of bugs (*e.g.*, memory safety for Janus [70]), or needs human effort to create checkers [36].

kAFL [56] supports feedback-driven fuzzing for operating system kernels assisted by Intel’s Processor Trace (PT) technology that can obtain complete trace information of the OS under testing. The trace information provides coverage for ring-0 execution of OS code and gives feedback to kAFL for maximizing code coverage with only a little overhead. Still, kAFL tracks only strong indications like OS crashes and cannot detect hidden bugs triggered quietly. For example, Ext4 exposed a bug [36, 53] that returns `EISDIR` when calls `unlink` for a directory without proper permission. Based on the POSIX specification, the expected error code should be `EPERM` instead of `EISDIR`. However, this bug does not lead to apparent consequences such as kernel crash or panic. Thus, it is difficult for kAFL to identify this kind of bug.

JANUS [70], a feedback-driven fuzzer, explores the two-dimensional input space of a file system (*i.e.*, metadata in the file system image and image-directed system calls). The image mutation in JANUS operates only on metadata because the file system image is too large to mutate. The metadata constitutes barely 1% of the image size but comprises essential file system attributes and structures. The input system calls in JANUS are not randomly created. Instead, JANUS deduces the runtime status of every file object on the image after generating complete system calls. Therefore, this runtime status contributes context-aware file operations to the syscall fuzzer, which avoids runtime errors and gains high code coverage. Despite the impressive results, JANUS cannot detect silent semantic bugs because it relies on kernel crashes or panics to capture and recognize bugs.

Previous fuzzing frameworks failed to detect various types of file system bugs. As an extensible fuzzing framework, Hydra [36] was proposed to integrate multiple bug detectors to allow Hydra to discover more kinds of bugs. The Hydra framework provides mutated syscalls and file system images as fuzzing input. Hydra can combine existing bug checkers for file systems, including crash consistency, semantic, and memory bug checkers. However, the Hydra framework only supplies input for the fuzzing process. As a result, the file system developers have to implement the checkers by themselves, which is inconvenient and limits the usage of Hydra.

## Chapter 9

# Conclusions

We have proposed MCFS, a new model-checking framework for file system development that exhaustively explores a file system's (bounded) state space, without requiring manual modeling or significant changes to the kernel or the file system itself. We developed VeriFS (v1 and v2), prototypes that demonstrate state checkpointing and restoration functionality via `ioctl`s. These functions let MCFS track VeriFS's full state while avoiding cache incoherency. MCFS found real bugs while we were developing VeriFS. Because both versions of VeriFS are simple and fast to model-check, they serve as a useful baseline against which we can compare other file systems.

We discussed the challenges we encountered when implementing MCFS, convincing ourselves of the need for file-system-level support to allow MCFS to work correctly and efficiently. Finally, while MCFS is designed to check file systems, its underlying approach is applicable to other system software.

We are exploring methods to track code coverage while model-checking. We will also use Spin's swarm verification [27,28] to explore larger state spaces in parallel.

## Chapter 10

### Future Work

We plan to add checkpoint/restore API support to Linux kernel file systems (*e.g.*, Ext4) to eliminate the need for the current mount/remount workaround. We are implementing the checkpoint/restore API at the Linux VFS level, which we hope will apply to many Linux kernel file systems. We are also working on APIs that will checkpoint file system states to help us resume the model-checking process if an interruption occurs (*e.g.*, due to a kernel crash). We also plan to run more than two file systems concurrently with MCFS and use a majority-voting approach to recognize incorrect file-system behavior.

# Acknowledgments

I would like to express my sincere gratitude to my advisor Prof. Erez Zadok who has been immensely helpful and provided me with much valuable advice for this work. I also thank Wei Su for his outstanding contributions to this MCFS project. Wei always helps me catch and fix nasty system bugs, and I appreciate and benefit from every discussion with Wei. I would like to thank the mentors in the MCFS team, including Dr. Gerard Holzmann, Prof. Geoff Kuenning, and Prof. Scott Smolka, for their remarkable insights and constructive suggestions. They are profoundly knowledgeable and make this work on the right track throughout. I am also grateful to other students who have advanced this work, including Manish Adkar, Alex Bishka, Gomathi Ganesan, Tejeshwar Gurram, Pei Liu, Shushanth Madhubalan, Yatna Verma, and Haolin Yu.

I would like to thank Prof. Niranjana Balasubramanian and Prof. Scott Smolka for serving on my RPE committee and providing invaluable advice.

I also thank the ACM HotStorage anonymous reviewers and our shepherd, Ram Alagappan, for their valuable comments and helpful feedback. This work was made possible in part thanks to Dell-EMC, NetApp, and IBM support; and NSF awards CCF-1918225, CNS-1900706, CNS-1900589, CNS-1729939, CNS-1730726, DCL-2040599, and CPS-1446832.

# Bibliography

- [1] Alper Akcan. Fuse-ext2 GitHub repository, 2021. <https://github.com/alperakcan/fuse-ext2>.
- [2] Sidney Amani, Alex Hixon, Zilin Chen, Christine Rizkallah, Peter Chubb, Liam O’Connor, Joel Beeren, Yutaka Nagashima, Japheth Lim, Thomas Sewell, et al. Cogent: Verifying high-assurance file system implementations. In *Proceedings of the 21st International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, pages 175–188, Atlanta, GA, April 2016.
- [3] Christel Baier and Joost-Pieter Katoen. *Principles of model checking*. MIT press, 2008.
- [4] Fabrice Bellard. Qemu, a fast and portable dynamic translator. In *Proceedings of the FREENIX Track: 2005 USENIX Annual Technical Conference*, pages 41–46, Anaheim, CA, 2005.
- [5] Bogdan Gribincea. Ext4 data loss, 2009. <https://bugs.launchpad.net/ubuntu/+source/linux/+bug/317781?comments=all>.
- [6] James Bornholt, Antoine Kaufmann, Jialin Li, Arvind Krishnamurthy, Emina Torlak, and Xi Wang. Specifying and checking file system crash-consistency models. In *Proceedings of the 21st International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, pages 83–98, April 2016.
- [7] Cristian Cadar, Daniel Dunbar, and Dawson R. Engler. KLEE: Unassisted and automatic generation of high-coverage tests for complex systems programs. In *Proceedings of the 8th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, pages 209–224, San Diego, CA, December 2008.
- [8] Cristian Cadar, Vijay Ganesh, Peter M. Pawlowski, David L. Dill, and Dawson R. Engler. EXE: Automatically generating inputs of death. In *Proceedings of the 13th ACM Conference on Computer and Communications Security (CCS)*, pages 322–335, Alexandria, VA, November 2006.
- [9] Tej Chajed. Verifying an I/O-concurrent file system. Master’s thesis, Massachusetts Institute of Technology, 2017.
- [10] Haogang Chen, Tej Chajed, Alex Konradi, Stephanie Wang, Atalay Mert Ileri, Adam Chlipala, M. Frans Kaashoek, and Nickolai Zeldovich. Verifying a high-performance crash-safe file system using a tree specification. In *Proceedings of the 26th Symposium on Operating Systems Principles (SOSP)*, pages 270–286, Shanghai, China, October 2017.
- [11] 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 Symposium on Operating Systems Principles (SOSP)*, pages 18–37, Monterey, CA, October 2015.

- [12] Haogang Chen, Daniel Ziegler, Adam Chlipala, M. Frans Kaashoek, Eddie Kohler, and Nickolai Zeldovich. Specifying crash safety for storage systems. In *Proceedings of the 15th Workshop on Hot Topics in Operating Systems (HotOS)*, Kartause Ittingen, Switzerland, May 2015.
- [13] Edmund M Clarke, William Klieber, Miloš Nováček, and Paolo Zuliani. Model checking and the state explosion problem. In *LASER Summer School on Software Engineering*, pages 1–30. Springer, 2011.
- [14] Edmund M Clarke Jr, Orna Grumberg, Daniel Kroening, Doron Peled, and Helmut Veith. *Model Checking*. MIT press, 2018.
- [15] CRIU Community. Checkpoint/restore in userspace (CRIU), 2021. <https://criu.org/>.
- [16] Leonardo Mendonça de Moura and Nikolaj Bjørner. Z3: an efficient SMT solver. In *International conference on Tools and Algorithms for the Construction and Analysis of Systems*, pages 337–340, Budapest, Hungary, 2008.
- [17] C. Demartini, R. Iosif, and R. Sisto. Modeling and validation of Java multithreading applications using SPIN. In *Proceedings of the 4th SPIN Workshop*, Paris, France, 1998.
- [18] Ext4. <http://ext4.wiki.kernel.org/>.
- [19] Filipe Manana. Btrfs: add missing inode update when punching hole, 2015. <https://patchwork.kernel.org/patch/5830801/>.
- [20] Pedro Fonseca, Rodrigo Rodrigues, and Björn B. Brandenburg. SKI: exposing kernel concurrency bugs through systematic schedule exploration. In *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, pages 415–431, Broomfield, CO, October 2014.
- [21] Sishuai Gong, Deniz Altinbüken, Pedro Fonseca, and Petros Maniatis. Snowboard: Finding kernel concurrency bugs through systematic inter-thread communication analysis. In *ACM SIGOPS 28th Symposium on Operating Systems Principles (SOSP)*, pages 66–83, Virtual Event, 2021.
- [22] Google. syzkaller: Linux syscall fuzzer, 2021. <https://github.com/google/syzkaller>.
- [23] Ronghui Gu, Jérémie Koenig, Tahina Ramananandro, Zhong Shao, Xiongnan Wu, Shu-Chun Weng, Haozhong Zhang, and Yu Guo. Deep specifications and certified abstraction layers. In *Proceedings of the 42nd ACM Symposium on Principles of Programming Languages (POPL)*, Mumbai, India, January 2015.
- [24] John S. Heidemann and Gerald J. Popek. File-system development with stackable layers. *ACM Transactions on Computer Systems*, 12(1):58–89, 1994.
- [25] Gerard J. Holzmann. The model checker SPIN. *IEEE Transactions on software engineering*, 23(5):279–295, 1997.
- [26] Gerard J. Holzmann. Logic verification of ANSI-C code with SPIN. In *Proceedings of the 7th International SPIN Workshop on SPIN Model Checking and Software Verification*, pages 131–147, Berlin, Heidelberg, 2000. Springer-Verlag.
- [27] Gerard J. Holzmann, Rajeev Joshi, and Alex Groce. Swarm verification. In *2008 23rd IEEE/ACM International Conference on Automated Software Engineering*, pages 1–6. IEEE, 2008.
- [28] Gerard J. Holzmann, Rajeev Joshi, and Alex Groce. Swarm verification techniques. *IEEE Transactions on Software Engineering*, 37(6):845–857, 2010.

- [29] Gerard J. Holzmann and Doron Peled. An improvement in formal verification. In *Formal Description Techniques VII*, pages 197–211. Springer, 1995.
- [30] IEEE/ANSI. Information technology–test methods for measuring conformance to POSIX–part 1: System interfaces. Technical Report STD-2003.1, ISO/IEC, 1992.
- [31] Atalay Mert Ileri, Tej Chajed, Adam Chlipala, M. Frans Kaashoek, and Nickolai Zeldovich. Proving confidentiality in a file system using DiskSec. In *Proceedings of the 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, pages 323–338, Carlsbad, CA, October 2018.
- [32] Kernel.org Bugzilla. Ubuntu 10.04.2 random kernel panic on xfs write, 2011. [https://bugzilla.kernel.org/show\\_bug.cgi?id=41572](https://bugzilla.kernel.org/show_bug.cgi?id=41572).
- [33] Kernel.org Bugzilla. Corruption of vm qcow2 image file on ext4 with crypto enabled, 2016. [https://bugzilla.kernel.org/show\\_bug.cgi?id=118511](https://bugzilla.kernel.org/show_bug.cgi?id=118511).
- [34] Kernel.org Bugzilla. Metadata corruption after xfs\_fsr and reboot, 2018. [https://bugzilla.kernel.org/show\\_bug.cgi?id=198749](https://bugzilla.kernel.org/show_bug.cgi?id=198749).
- [35] Kernel.org Bugzilla. Btrfs bug entries, 2021. <https://bugzilla.kernel.org/buglist.cgi?component=btrfs>.
- [36] Seulbae Kim, Meng Xu, Sanidhya Kashyap, Jungyeon Yoon, Wen Xu, and Taesoo Kim. Finding semantic bugs in file systems with an extensible fuzzing framework. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles (SOSP)*, pages 147–161, Toronto, Ontario, Canada, October 2019.
- [37] A. Kivity, Y. Kamay, D. Laor, U. Lublin, and A. Liguori. kvm: the Linux virtual machine monitor. In *Proceedings of the 2007 Ottawa Linux Symposium (OLS 2007)*, volume 1, pages 225–230, Ottawa, Canada, June 2007.
- [38] Lenz Grimmer. Frequent kernel panics when doing heavy i/o in lxc containers on btrfs, 2015. <https://bugs.launchpad.net/ubuntu/+source/linux/+bug/1415510>.
- [39] David Lie, Andy Chou, Dawson Engler, and David L. Dill. A simple method for extracting models from protocol code. In *Proceedings of the 28th Annual International Symposium on Computer Architecture (ISCA)*, pages 192–203, Gothenburg, Sweden, July 2001.
- [40] Lanyue Lu, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau, and Shan Lu. A study of Linux file system evolution. In *Proceedings of the USENIX Conference on File and Storage Technologies (FAST '13)*, San Jose, CA, February 2013. USENIX Association.
- [41] Filipe Manco, Costin Lupu, Florian Schmidt, Jose Mendes, Simon Kuenzer, Sumit Sati, Kenichi Yasukata, Costin Raiciu, and Felipe Huici. My VM is lighter (and safer) than your container. In *Proceedings of the 26th Symposium on Operating Systems Principles (SOSP)*, pages 218–233, Shanghai, China, October 2017.
- [42] Mark McLoughlin. The qcow2 image format, 2008.
- [43] Changwoo Min, Sanidhya Kashyap, Byoungyoung Lee, Chengyu Song, and Taesoo Kim. Cross-checking semantic correctness: The case of finding file system bugs. In *Proceedings of the 25th Symposium on Operating Systems Principles (SOSP)*, pages 361–377, Monterey, CA, October 2015.
- [44] Subrata Modak. Linux test project (LTP), 2009. <http://ltp.sourceforge.net/>.



- [45] Jayashree Mohan, Ashlie Martinez, Soujanya Ponnappalli, Pandian Raju, and Vijay Chidambaram. Finding crash-consistency bugs with bounded black-box crash testing. In *Proceedings of the 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, pages 33–50, Carlsbad, CA, October 2018.
- [46] Madanlal Musuvathi, Andy Chou, David L. Dill, and Dawson R. Engler. Model checking system software with CMC. In *Proceedings of the 10th ACM SIGOPS European Workshop*, pages 219–222, Saint-Emilion, France, July 2002.
- [47] Madanlal Musuvathi and Dawson R. Engler. Model checking large network protocol implementations. In *Proceedings of the 1st Symposium on Networked Systems Design and Implementation (NSDI)*, pages 155–168, San Francisco, CA, March 2004.
- [48] Madanlal Musuvathi, David Y. W. Park, Andy Chou, Dawson R. Engler, and David L. Dill. CMC: A pragmatic approach to model checking real code. In *Proceedings of the 5th Symposium on Operating System Design and Implementation (OSDI)*, Boston, MA, December 2002.
- [49] NFS-Ganesha, 2016. <http://nfs-ganesha.github.io/>.
- [50] Oracle. Btrfs, 2008. <http://oss.oracle.com/projects/btrfs/>.
- [51] 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 USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, pages 433–448, Broomfield, CO, October 2014.
- [52] 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, Brighton, UK, October 2005. ACM Press.
- [53] Tom Ridge, David Sheets, Thomas Tuerk, Andrea Giugliano, Anil Madhavapeddy, and Peter Sewell. SibiYFS: formal specification and oracle-based testing for POSIX and real-world file systems. In *Proceedings of the 25th Symposium on Operating Systems Principles (SOSP)*, pages 38–53, Monterey, CA, October 2015.
- [54] Ohad Rodeh, Josef Bacik, and Chris Mason. BTRFS: The Linux B-tree filesystem. *ACM Transactions on Storage (TOS)*, 9(3):1–32, 2013.
- [55] Cindy Rubio-González, Haryadi S. Gunawi, Ben Liblit, Remzi H. Arpaci-Dusseau, and Andrea C. Arpaci-Dusseau. Error propagation analysis for file systems. In *Proceedings of the 30th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, pages 270–280, Dublin, Ireland, June 2009.
- [56] Sergej Schumilo, Cornelius Aschermann, Robert Gawlik, Sebastian Schinzel, and Thorsten Holz. kAFL: Hardware-assisted feedback fuzzing for OS kernels. In *26th USENIX Security Symposium (USENIX Security)*, pages 167–182, Vancouver, BC, Canada, August 2017.
- [57] SGI XFS. xfstests, 2016. [http://xfs.org/index.php/Getting\\_the\\_latest\\_source\\_code](http://xfs.org/index.php/Getting_the_latest_source_code).

- [58] Helgi Sigurbjarnarson, James Bornholt, Emina Torlak, and Xi Wang. Push-button verification of file systems via crash refinement. In *Proceedings of the 12th Symposium on Operating Systems Design and Implementation (OSDI)*, pages 1–16, Savannah, GA, November 2016.
- [59] Spin Project. Static analysis tools for C code. [www.spinroot.com/static/](http://www.spinroot.com/static/).
- [60] Miklos Szeredi. Filesystem in Userspace. <http://fuse.sourceforge.net>, February 2005.
- [61] Chia-che Tsai, Bhushan Jain, Nafees Ahmed Abdul, and Donald E. Porter. A study of modern Linux API usage and compatibility: What to support when youre supporting. In *Proceedings of the Eleventh European Conference on Computer Systems (EuroSys)*, pages 1–16, London, United Kingdom, April 2016.
- [62] T. Ts'o. E2fsprogs: Ext2/3/4 filesystem utilities, 2008. <http://e2fsprogs.sourceforge.net>.
- [63] Stephen Tweedie. Ext3, journaling filesystem. In *Ottawa Linux Symposium*, July 2000. <http://olstrans.sourceforge.net/release/OLS2000-ext3/OLS2000-ext3.html>.
- [64] Ubuntu bugs launchpad. kernel hangs on boot with btrfs, 2012. <https://bugs.launchpad.net/ubuntu/+source/linux/+bug/1088185>.
- [65] Vijay Chidambaram. Strange behavior (possible bugs) in btrfs, 2018. <https://www.spinics.net/lists/linux-btrfs/msg77415.html>.
- [66] David Woodhouse, Joern Engel, Jarkko Lavinien, and Artem Bityutskiy. General MTD documentation, 2009.
- [67] David Woodhouse, Joern Engel, Jarkko Lavinien, and Artem Bityutskiy. JFFS2, 2009.
- [68] XFS – high-performance 64-bit journaling file system. <https://www.linuxlinks.com/xfs/>. Visited February, 2021.
- [69] Wen Xu, Sanidhya Kashyap, Changwoo Min, and Taesoo Kim. Designing new operating primitives to improve fuzzing performance. In *Proceedings of the ACM SIGSAC Conference on Computer and Communications Security (CCS)*, pages 2313–2328, Dallas, TX, October 2017.
- [70] Wen Xu, Hyungon Moon, Sanidhya Kashyap, Po-Ning Tseng, and Taesoo Kim. Fuzzing file systems via two-dimensional input space exploration. In *IEEE Symposium on Security and Privacy (Oakland)*, pages 818–834, San Francisco, CA, May 2019.
- [71] Junfeng Yang, Can Sar, and Dawson Engler. eXplode: a lightweight, general system for finding serious storage system errors. In *Proceedings of the 7th symposium on Operating systems design and implementation (OSDI)*, pages 131–146, Seattle, WA, November 2006.
- [72] Junfeng Yang, Paul Twohey, Dawson Engler, and Madanlal Musuvathi. Using model checking to find serious file system errors. In *Proceedings of the 6th Symposium on Operating Systems Design and Implementation (OSDI)*, pages 273–288, San Francisco, CA, December 2004.
- [73] Erez Zadok, Rakesh Iyer, Nikolai Joukov, Gopalan Sivathanu, and Charles P. Wright. On incremental file system development. *ACM Transactions on Storage (TOS)*, 2(2):161–196, 2006.
- [74] Mo Zou, Haoran Ding, Dong Du, Ming Fu, Ronghui Gu, and Haibo Chen. Using concurrent relational logic with helpers for verifying the AtomFS file system. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles (SOSP)*, pages 259–274, Toronto, Ontario, Canada, October 2019.