

Formalising filesystems in the ACL2 theorem prover: an application to FAT32

Mihir Parang Mehta

Department of Computer Science
University of Texas at Austin
Austin, TX, USA
`mihir@cs.utexas.edu`

In this work, we present an approach towards constructing executable specifications of existing filesystems and verifying their functional properties in a theorem proving environment. We detail an application of this approach to the FAT32 filesystem.

We also detail the methodology used to build up this type of executable specification through a series of models which incrementally add features of the target filesystem. This methodology has the benefit of allowing the verification effort to start from simple models which encapsulate features common to many filesystems and which are thus suitable for re-use.

1 Introduction and overview

Filesystems are ubiquitous in computing, providing application programs a means to store data persistently, address data by a name instead of a numeric index, and communicate with other programs. Thus, the vast majority of application programs directly or indirectly rely upon filesystems, which makes filesystem verification critically important. Here, we present a formalisation effort in ACL2 for an implementation of the FAT32 filesystem, and a proof of the read-over-write properties for this filesystem. By starting with a high-level abstract model and adding more filesystem features in successive models, we are able to manage the complexity of this proof, which has not, to our knowledge, been previously attempted. Thus, this paper contributes an implementation of FAT32, a widely-used filesystem, formally verified against an abstract specification and tested for binary compatibility by means of co-simulation.

In the rest of this paper, we describe these filesystem models and the properties proved with examples; we proceed to a high-level explanation of these proofs and the co-simulation infrastructure; and further we offer some insights about the low-level issues encountered while working the proofs.

2 Related work

Filesystem verification research has largely followed a pattern of synthesising a new filesystem based on a specification chosen for its ease in proving properties of interest, rather than similarity to an existing filesystem. Our work, in contrast, follows the FAT32 specification closely. In spirit, our work is closer to previous work which uses interactive theorem provers and explores deep functional properties than to efforts which use non-interactive theorem provers such as Z3 to produce fully automated proofs of simpler properties.

2.1 Interactive theorem provers

An early effort in the filesystem verification domain was by Bevier and Cohen [3], who specified the Synergy filesystem and created an executable model of the same in ACL2 [8], down to the level of processes and file descriptors. On the proof front, they certified their model to preserve well-formedness of their data structures through their various file operations; however, they did not attempt to prove, for instance, read-over-write properties or crash consistency. Later, Klein et al with the SeL4 project [9] used Isabelle/HOL [12] to verify a microkernel; while their design abstracted away file operations in order to keep their trusted computing base small, it did serve as a precursor to their more recent COGENT project [1]. Here the authors built a "verified compiler" of sorts, generating C-language code from specifications in their domain-specific in a manner guaranteed to avoid many common filesystem bugs. Elsewhere, the SibylFS project [13], again using Isabelle/HOL, provided an executable specification for filesystems at a level of abstraction that could function across multiple operating systems including OSX and Unix. The Coq prover [2] has also been used, for instance, for FSCQ [4], a state-of-the-art filesystem which was built to have high performance and formally verified crash consistency properties.

2.2 Non-interactive theorem provers

Non-interactive theorem provers such as Z3 [5] have also been used; Hyperkernel [11] is a recent effort which focusses on simplifying the xv6 microkernel until the point that Z3 can verify it with its SMT solving techniques. However, towards this end, all system calls in Hyperkernel are replaced with analogs which can terminate in constant time; while this approach is theoretically sound, it increases the chances of discrepancies between the model and the implementation which may diminish the utility of the proofs or even render them moot. A stronger effort in the same domain is Yggdrasil [14], which focusses on verifying filesystems with the use of Z3. While the authors make substantial progress in terms of the number of filesystem calls they support and the crash consistency guarantees they provide, they are subject to the limits of SMT solving which prevent them from modelling essential filesystem features such as extents, which are central to many filesystems including FAT32.

3 Program architecture and performance considerations

We number our abstract filesystem models L1 through L6, and our concrete models M1 through M2 (concrete). These models are constructed incrementally to allow for reuse of features in general, and a refinement relation where possible.

Starting with L1, an abstract model representing the directory structure as a tree, we add file metadata in L2. We branch off from L2 into L3, a model with file contents broken up into blocks and stored in an external disk-like data structure, and L4, a model which also breaks file contents up into disk blocks, additionally with garbage collection through reference counting, and a fixed disk size bounding the total size of file contents. To implement the reference counter in L4, we use an allocation vector as in the CP/M filesystem, which happens to be the previous technology target of this work prior to FAT32. We are able to refine L4 into L6, a filesystem with exactly the same properties but additionally featuring a FAT32-like file allocation table for allocation and garbage collection, since this data structure happens to refine the allocation vector.

M1 is another tree model, which hews somewhat closely to the FAT32 specification; it is also used as the in-memory format for M2, a real model which reads and writes FAT32 disk images. The operations currently supported are `pread(2)`, `pwrite(2)`, and `stat(2)`. With these operations, we are able to

model the file operations used in standard disk utilities such as `cat(1)` and `dd(1)`, with a view to writing formal specifications of these programs later.

There is not a refinement relation between the abstract models and the concrete models, but we re-use many of our theorems from L6 in M2 since the file allocation data structure in L6 follows the FAT32 specification.

In order to obtain reasonable execution performance, we implement M2 as a single-threaded object. Most importantly, this allows us to model the file allocation table and data region, which are long arrays, as arrays in Common Lisp rather than as lists, which have poor performance for update operations on their elements because of the number of cons cells which must be created and garbage collected for each update. Array operations in single threaded objects are subject to certain syntactic restrictions to prevent copies of arrays from being created as usually happens with Lisp objects; to simplify the task of reading and writing disk images under these restrictions, we also create a library of useful macros.

We choose to implement a subset of the POSIX filesystem application programming interface. This allows us to easily compare the results of running filesystem operations on M2 and the Linux kernel's implementation of FAT32, which in turn allows us to test our implementation's correctness through co-simulation in addition to theorem proving. One trade-off for this choice is the necessity of emulating certain functionality provided by the Linux kernel to all filesystems; thus, we implement process tables and file tables through a straightforward approach similar to that used in Synergy [3].

4 The FAT32 filesystem

FAT32 was initially developed at Microsoft [10] in order to address certain shortcomings of the DOS filesystem previously in use in their operating systems. While it is simple by today's standards, it does add some complexity compared to the filesystems which came before.

All files, including regular files and directory files, are divided into *clusters* (sometimes called *extents*) of a fixed size, which is decided at the time a FAT32 volume is formatted, and constrained to be a multiple of the disk sector size. Directory files differ only in a metadata attribute which indicates that their contents should be treated as a sequence of directory entries. Each such directory entry is 32 bytes wide and contains information including name, size, first cluster index, and access times for the corresponding file.

The file allocation table itself contains a number of linked lists. It maps each cluster index used by a file to either the next cluster index for that file or a special end-of-clusterchain value ¹ This allows the contents of a file to be reconstructed by reading just the first cluster index from the corresponding directory entry, and building the list of clusters using the table. Unused clusters are mapped to 0 in the table; this fact is used for counting and allocating free clusters.

We illustrate the file allocation table and data layout for a small directory tree in figure 1.

5 The models

For every read or write operation, FAT32 requires one or more lookups into the file allocation table, followed by the corresponding lookups into the data region. This makes proof efforts about these operations complex, which serves as the motivation for modelling the filesystem in a series of steps.

¹ There is actually a range of end-of-clusterchain values in the specification, not just one. We support all values in the range.

Figure 1: A FAT32 directory tree

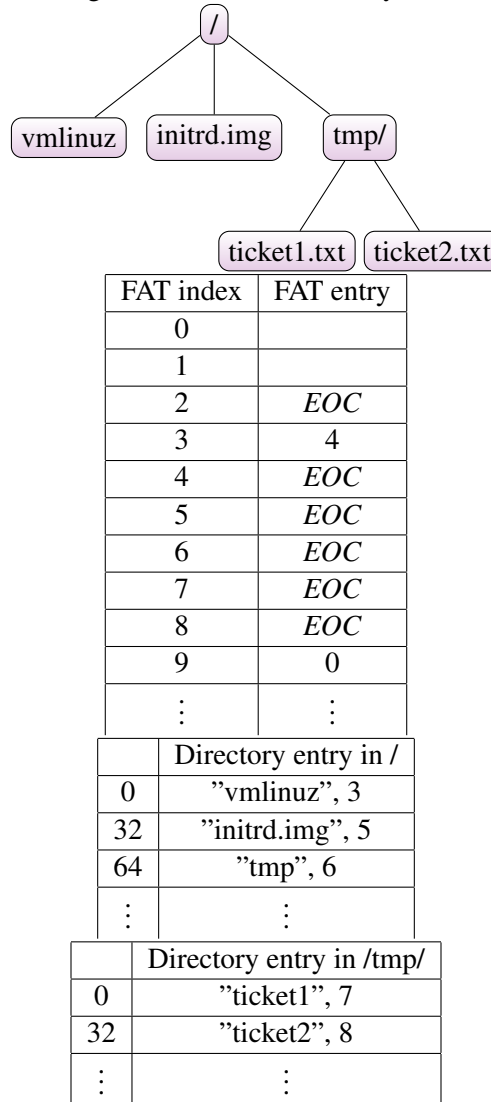


Figure 2: Refinement relationships between models

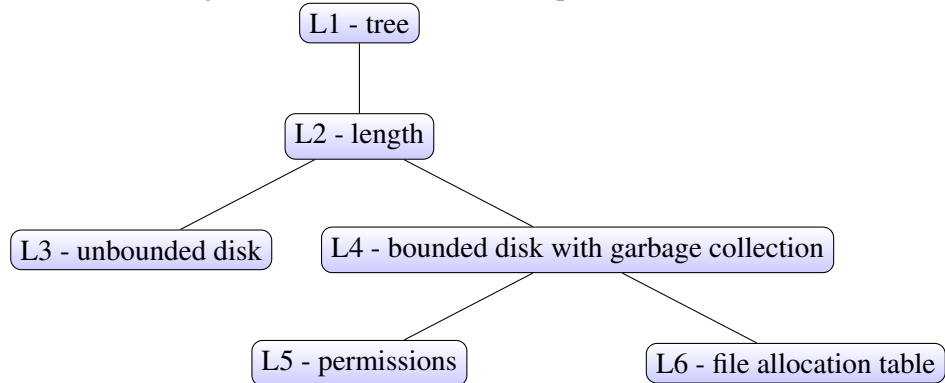


Table 1: Abstract models and their features

L1	The filesystem is represented as a tree, with leaf nodes for regular files and non-leaf nodes for directories. The contents of regular files are represented as strings stored in the nodes of the tree; the storage available for these is unbounded.
L2	A single element of metadata, <i>length</i> , is stored within each regular file.
L3	The contents of regular files are divided into blocks of fixed size. These blocks are stored in an external "disk" data structure; the storage for these blocks remains unbounded.
L4	The storage available for blocks is now bounded. An allocation vector data structure is introduced to help allocate and garbage collect blocks.
L5	Additional metadata for file ownership and access permissions is stored within each regular file.
L6	The allocation vector is replaced by a file allocation table, matching the official FAT specification.

Table 2: Concrete models

M1	The filesystem is represented as a tree, as in L1. Write operations show non-determinism in order to model disk capacity errors in a real filesystem.
M2	A single element of metadata, <i>length</i> , is stored within each regular file.

At this point in development, we have six models of the filesystem, here referred to as L1 through L6, described in table 1. Each model other than L1 refines a previous model, adding some features and complexity and thereby approaching closer to a model which is binary compatible with FAT32. These refinement relationships are shown in figure 2. L1 is the simplest of these, representing the filesystem as a literal directory tree; later models feature file metadata (including ownership and permissions), externalisation of file contents, and allocation/file allocation using an allocation vector after the fashion of the CP/M file system (this is a remnant of an earlier filesystem verification effort for CP/M, which we subsumed into the present work).

Broadly, we characterise the filesystem operations we offer as either *write* operations, which do modify the filesystem, or *read* operations, which do not. In each model, we have been able to prove *read-over-write* properties which show that write operations have their effects made available immediately for reads at the same location, but also that they do not affect reads at other locations.

The first read-after-write theorem states that immediately following a write of some text at some location, a read of the same length at the same location yields the same text. The second read-after-write theorem states that after a write of some text at some location, a read at any other location returns exactly what it would have returned before the write. As an example, listings for the L1 versions of these theorems follow.

```
(defthm l1-read-after-write-1
  (implies (and (l1-fs-p fs)
                (stringp text)
                (symbol-listp hns)
                (natp start)
                (equal n (length text)))
```

```

      (stringp (l1-stat hns fs)))
    (equal (l1-rdchs hns (l1-wrchs hns fs start text) start n) text)))

(defthm l1-read-after-write-2
  (implies (and (l1-fs-p fs)
                (stringp text2)
                (symbol-listp hns1)
                (symbol-listp hns2)
                (not (equal hns1 hns2))
                (natp start1)
                (natp start2)
                (natp n1)
                (stringp (l1-stat hns1 fs)))
            (equal (l1-rdchs hns1 (l1-wrchs hns2 fs start2 text2) start1 n1)
                  (l1-rdchs hns1 fs start1 n1))))

```

By composing these properties, we can reason about executions involving multiple reads and writes, as illustrated in the following throwaway proof.

```

(thm
  (implies (and (l1-fs-p fs)
                (stringp text1)
                (stringp text2)
                (symbol-listp hns1)
                (symbol-listp hns2)
                (not (equal hns1 hns2))
                (natp start1)
                (natp start2)
                (stringp (l1-stat hns1 fs))
                (equal n1 (length text1)))
            (equal (l1-rdchs hns1
                          (l1-wrchs hns2 (l1-wrchs hns1 fs start1 text1)
                                         start2 text2)
                          start1 n1)
                  (l1-rdchs hns1 (l1-wrchs hns1 fs start1 text1)
                          start1 n1))))

```

6 Proof methodology

In L1, our simplest model, the read-over-write properties were, of necessity, proven from scratch.

In each subsequent model, the read-over-write properties are proven as corollaries of equivalence proofs which establish the correctness of read and write operations in the respective model with respect to a previous model. A representation of such an equivalence proof can be seen in figures 3, 4 and 5, which respectively show the equivalence proof for `l2-wrchs`, the equivalence proof for `l2-rdchs` and the composition of these to obtain the first read-over-write theorem for model L2.

Figure 3: l2-wrchs-correctness-1

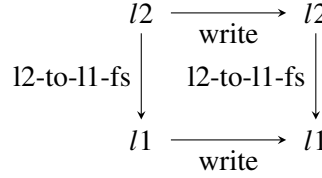
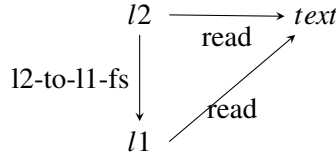


Figure 4: l2-rdchs-correctness-1



7 Some proof details

We have come to rely on certain principles for the proof effort for each new model. We summarise these below.

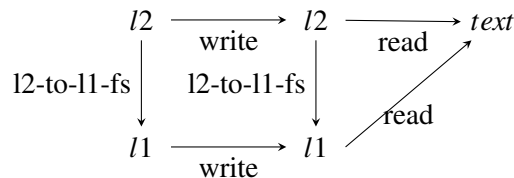
7.1 Invariants

As the models grow more complex, with the addition of more auxiliary data the "sanity" criteria for filesystem instances become more complex. For instance, in L4, the predicate $l4\text{-fs-p}$ is defined to be the same as $l3\text{-fs-p}$, which recursively defines the shape of a valid directory tree. However, we choose to require two more properties for a "sane" filesystem.

1. Each disk index assigned to a regular file should be marked as *used* in the allocation vector - this is essential to prevent filesystem errors.
2. Each disk index assigned to a regular file should be distinct from all other disk indices assigned to files - this does not hold true, for example, in filesystems with hardlinks, but makes our proofs easier.

These properties are invariants to be maintained across write operations; while not all of them are strictly necessary for a filesystem instance to be valid, they do simplify the verification of read-after-write properties by helping us ensure that write operations do not create an "aliasing" situation in which a regular file's contents can be modified through a write to a different regular file.

Figure 5: l2-read-over-write-1



These properties, in the form of the predicates `indices-marked-listp` and `no-duplicatessp`, are packaged together into the `l4-stricter-fs-p` predicate, for which a listing follows.

```
(defun l4-stricter-fs-p (fs alv)
  (declare (xargs :guard t))
  (and (l4-fs-p fs)
        (boolean-listp alv)
        (let ((all-indices (l4-list-all-indices fs)))
          (and (no-duplicatessp all-indices)
                (indices-marked-p all-indices alv))))))
```

7.2 Reuse

As noted earlier, using a refinement methodology allows us to derive our read-over-write properties essentially "for free"; more precisely, we are able to prove read-over-write properties simply with `:use` hints after having done the work of proving refinement through induction.

At a lower level, we are also able to benefit from refinement relationships between components of our different models. For example, such a relationship exists between the allocation vector used in L4 and the file allocation table used in L6. More precisely, by taking a file allocation table and mapping each non-zero entry to `true` and each zero entry to `false`, we obtain a corresponding allocation vector with exactly the same amount of available space. This is a refinement mapping which makes it a lot easier to prove that L4, which uses an allocation vector, is an abstraction of L6, which uses a file allocation table. This, in turn, means that the effort spent on proving the invariants described above for L4 need not be replicated for L6.

8 Co-simulation

Previous work on executable specifications [] has shown the importance of testing these on real examples, in order to validate that the behaviour shown matches that of the system being specified. In our case, this means we must validate our filesystem by testing it in execution against a canonical implementation of FAT32; in this case, we choose the implementation which ships with Linux kernel 3.10.

We use `mkfs.fat(8)`, a program which produces FAT32 disk images, for our tests. When run with the `-v` flag, this program emits an English-language summary of the fields of the newly created disk image; we make use of this by writing an ACL2 program based on our model which reads the image and reproduces this summary. This validates our code for reading the various fields of the disk image and gives us a regression test to use while we modify our model to support proofs and filesystem calls.

Our first co-simulation test is for `cat(1)`, a simple program which reads its input and copies it to its output. Its functionality is reproduced in an ACL2 program which uses our implementations of the system calls `lstat(2)`, `open(2)`, `pread(2)`, `pwrite(2)`, and `close(2)`. This allows us to validate our code for reading and writing regular files and directories which span multiple clusters.

We further test the `dd(1)` program, which provides similar functionality to `cat` but with more options to customise the data transfer, for instance, by allowing the data transfer to be split into blocks of a specified size.

9 Conclusion

This work formalises a FAT32-like filesystem and proves read-over-write properties through refinement of a series of models. Further, it proves the correctness of FAT32’s allocation and garbage collection mechanisms, and provides artefacts to be used in a subsequent realistic model of FAT32.

10 Future work

Our primary goal is to dispense with the tree representation and implement filesystem traversal by looking up entries in directory files. This will also involve addressing a subtle issue where reads affect the state of a filesystem by means of updating the access time, which has been analysed earlier in the context of microprocessors [6]. This will yield a model which is entirely contained in a disk data structure, without auxiliary data structures such as the tree, and which can further be validated by co-simulation with a FAT32 implementation such as that of Linux.

Next, we hope to re-use some artefacts of verifying FAT32 in order to verify a more complex filesystem, such as ext4. Choosing a filesystem with journalling will allow us to model crash consistency.

Finally, we hope to support “code proofs”, by providing a basis for reasoning about filesystem operations in filesystem-specific utilities such as `fsck`, as well as other application programs. This is a large part of the motivation for pursuing binary compatibility.

Acknowledgments.

This material is based upon work supported by the National Science Foundation SaTC program under contract number CNS-1525472. Thanks are also due to Warren A. Hunt Jr. and Matthew J. Kaufmann for their guidance.

11 Bibliography

References

- [1] Sidney Amani, Alex Hixon, Zilin Chen, Christine Rizkallah, Peter Chubb, Liam O’Connor, Joel Beeren, Yutaka Nagashima, Japheth Lim, Thomas Sewell et al. (2016): *Cogent: Verifying high-assurance file system implementations*. In: *ACM SIGPLAN Notices*, 51, ACM, pp. 175–188.
- [2] Yves Bertot & Pierre Castéran (2013): *Interactive theorem proving and program development: CoqArt: the calculus of inductive constructions*. Springer Science & Business Media.
- [3] William R Bevier & Richard M Cohen (1996): *An executable model of the Synergy file system*. Technical Report, Technical Report 121, Computational Logic, Inc.
- [4] Haogang Chen, Daniel Ziegler, Tej Chajed, Adam Chlipala, M. Frans Kaashoek & Nikolai Zeldovich (2016): *Using Crash Hoare Logic for Certifying the FSCQ File System*. In Gulati & Weatherspoon [7], doi:10.1145/2815400.2815402. Available at <https://www.usenix.org/conference/atc16>.
- [5] Leonardo De Moura & Nikolaj Bjørner (2008): *Z3: An efficient SMT solver*. In: *International conference on Tools and Algorithms for the Construction and Analysis of Systems*, Springer, pp. 337–340.
- [6] Shilpi Goel, Warren A Hunt & Matt Kaufmann (2017): *Engineering a formal, executable x86 ISA simulator for software verification*. In: *Provably Correct Systems*, Springer, pp. 173–209.

- [7] Ajay Gulati & Hakim Weatherspoon, editors (2016): *2016 USENIX Annual Technical Conference, USENIX ATC 2016, Denver, CO, USA, June 22-24, 2016*. USENIX Association. Available at <https://www.usenix.org/conference/atc16>.
- [8] Matt Kaufmann, Panagiotis Manolios & J. Strother Moore (2000): *Computer-aided reasoning: an approach*. Kluwer Academic Publishers.
- [9] Gerwin Klein, Kevin Elphinstone, Gernot Heiser, June Andronick, David Cock, Philip Derrin, Dhammika Elkaduwe, Kai Engelhardt, Rafal Kolanski, Michael Norrish et al. (2009): *seL4: Formal verification of an OS kernel*. In: *Proceedings of the ACM SIGOPS 22nd symposium on Operating systems principles*, ACM, pp. 207–220.
- [10] Microsoft (2000): *Microsoft Extensible Firmware Initiative FAT32 File System Specification*. Available at www.microsoft.com/hwdev/download/hardware/fatgen103.pdf.
- [11] Luke Nelson, Helgi Sigurbjarnarson, Kaiyuan Zhang, Dylan Johnson, James Bornholt, Emina Torlak & Xi Wang (2017): *Hyperkernel: Push-Button Verification of an OS Kernel*. In: *Proceedings of the 26th Symposium on Operating Systems Principles, SOSP '17*, ACM, New York, NY, USA, pp. 252–269, doi:10.1145/3132747.3132748. Available at <http://doi.acm.org/10.1145/3132747.3132748>.
- [12] Tobias Nipkow, Lawrence C Paulson & Markus Wenzel (2002): *Isabelle/HOL: a proof assistant for higher-order logic*. 2283, Springer Science & Business Media.
- [13] Tom Ridge, David Sheets, Thomas Tuerk, Andrea Giugliano, Anil Madhavapeddy & Peter Sewell (2015): *SibylFS: formal specification and oracle-based testing for POSIX and real-world file systems*. In: *Proceedings of the 25th Symposium on Operating Systems Principles*, ACM, pp. 38–53.
- [14] Helgi Sigurbjarnarson, James Bornholt, Emina Torlak & Xi Wang (2016): *Push-Button Verification of File Systems via Crash Refinement*. In: *OSDI*, 16, pp. 1–16.