



Formal verification of a concurrent file system

Tej Chajed
MIT

Systems software is challenging to get right

Expose an interface to applications



Run in a complex environment

2

Systems software continues to have bugs

LWN.net News from the source

A tale of two data-corruption bugs

By Jonathan Corbet | May 24, 2015 | FOSDEM 2019 / Schedule / Events / Main tracks / Databases / PostgreSQL vs. fsync

There have been two bugs causing filesystem corruption in the news recently. One of them, a bug in ext4, has gotten the bulk of the attention, despite the fact that it is an old bug that is hard to trigger.

How is it possible that PostgreSQL used fsync incorrectly for 20 and what we'll do about it.

= threatpost

MoleRats APT Launches Spy Campaign on Bankers, Politicians, Journalists

BRATA Android Trojan Update

About a year ago the possibly disastrous cc will walk you through its misunderstanding like PostgreSQL deploys

Linux Servers at Risk of RCE Due to Critical CWP Bugs

INFOSEC Supply

WIRED BACKCHANNEL BUSINESS CULTURE GEAR IDEAS SCIENCE SECURITY

LILY HAN REED SECURITY DEC 18, 2021 2:54 PM

'The Internet Is on Fire'

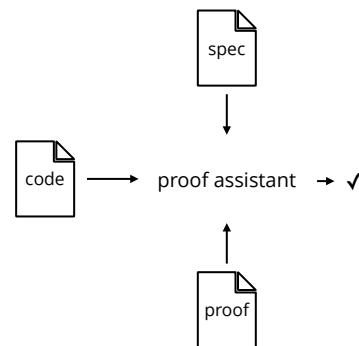
A vulnerability in the Log4j logging framework has security teams scrambling to put in a fix.

Vision: reliable systems software using verification



4

Approach: systems verification



5

Verification can eliminate whole classes of bugs

Memory safety

Concurrency bugs

Logic errors

Security flaws

6

Systems verification is practical

Microkernels (e.g., seL4, CertiKOS)

Cryptography libraries (e.g., Fiat Crypto, HACL*)

Distributed systems (e.g., Ironfleet, Verdi)

File systems (e.g., FSCQ)

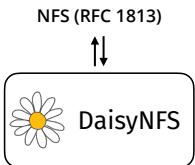
7

My research spans from foundations to systems

systems	draft '22	DaisyNFS	concurrent, crash-safe file system
	OSDI'21	GoTxn	high-performance transaction system
foundations	SOSP '15, SOSP '17	FSCQ	verified sequential file system
	CoqPL '20	Goose	reasoning about Go code
foundations	PLDI '19	Argosy	layered recovery procedures
	OSDI '18	CSPEC	concurrency using movers
foundations	SOSP '19, OSDI '21	Perennial	concurrency + crash safety
	SOSP '15	Crash Hoare Logic (CHL)	sequential crash safety

8

DaisyNFS is a verified NFS server



Theorem: Every NFS operation appears to execute atomically and correctly, despite crashes and concurrency.

9

System design



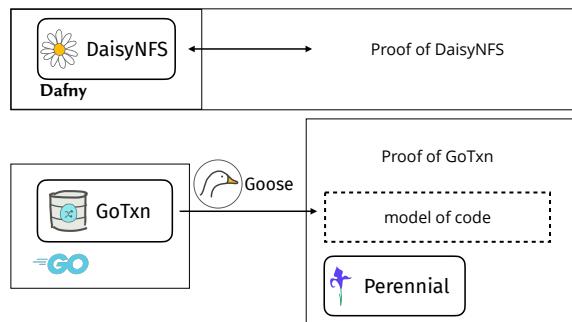
File-system code:
uses one transaction per operation,
implements all the features



Transaction system:
automatically gives atomicity,
internally concurrent

10

Verification overview



11

Main research contribution



First approach for verifying a *high-performance, concurrent, crash-safe file system*.

12

Related work

crash safety:

FSCQ [SOSP '15], Yggdrasil [OSDI '16], DFSCQ [SOSP '17],
Argosy [PLDI '19], VeriBetrFS [OSDI '20]

concurrency:

Concurrent GC [MSR Tech Report '15],
CertiKOS [OSDI '16], AtomFS [SOSP '19]

crash safety + concurrency:

Flashix [2021]

13

Outline



transaction system



file system



Evaluation



Future work

14



GoTxn

15

File system needs crash safety

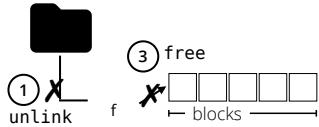
Crash could be due to power failure, kernel panic, etc.

After reboot, file system should recover old data

16

REMOVE has several steps

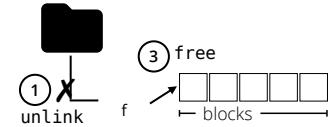
```
func REMOVE(d_ino: uint64,  
           name: []byte) {  
    ① f := unlink(d_ino, name)  
    ② blocks := getBlocks(f)  
    ③ free(blocks)  
}
```



17

Crashes create subtle bugs

```
func REMOVE(d_ino: uint64,  
           name: []byte) {  
    ① f := unlink(d_ino, name)  
    ② blocks := getBlocks(f) crash  
    ③ free(blocks)  
}
```

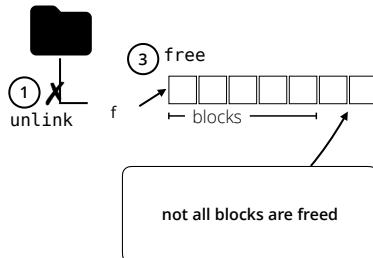


crash leaks f's blocks

18

Concurrency also create subtle bugs

```
func REMOVE(d_ino: uint64,  
           name: []byte) {  
    ① f := unlink(d_ino, name)  
    ② blocks := getBlocks(f) concurrent append  
    ③ free(blocks)  
}
```



19

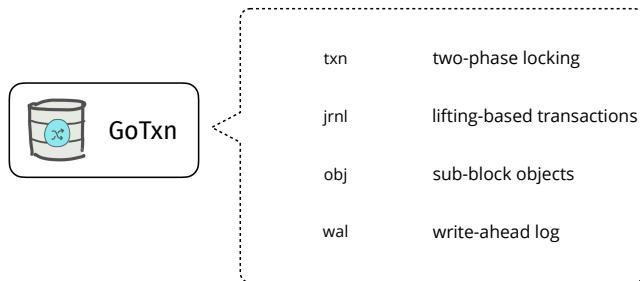
Transactions automatically give atomicity

```
func Begin() *Txn  
func (tx *Txn) Read(...)  
func (tx *Txn) Write(...)  
func (tx *Txn) Commit()  
GoTxn
```

Code between Begin() and Commit() is atomic both on crash and to other threads

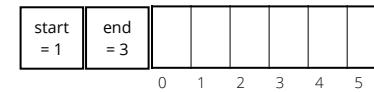
20

Crashes & concurrency are challenges throughout GoTxn



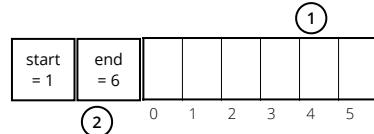
21

Multi-block atomic writes come from on-disk log



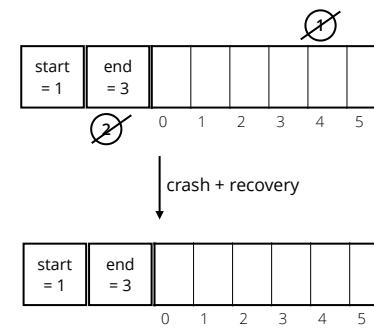
22

Multi-block atomic writes come from on-disk log



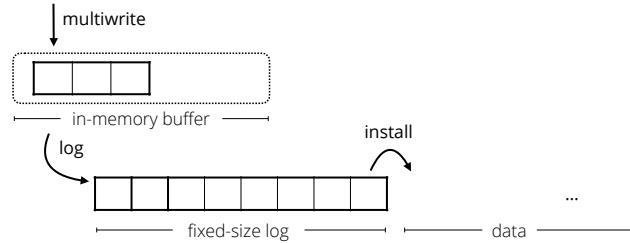
23

Crash never leaves a partial multiwrite



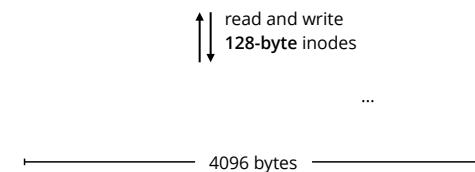
24

Writing, logging, and installation are concurrent



25

Object layer supports concurrent sub-block access



26

Two-phase locking handles concurrency

```
tx := Begin()          lock(0)
v := tx.Read(0)
tx.Write(1, v)
tx.Write(2, v)
tx.Commit()
obj.Multiwrite(...)
tx.releaseAll()
```

27

Verification challenges

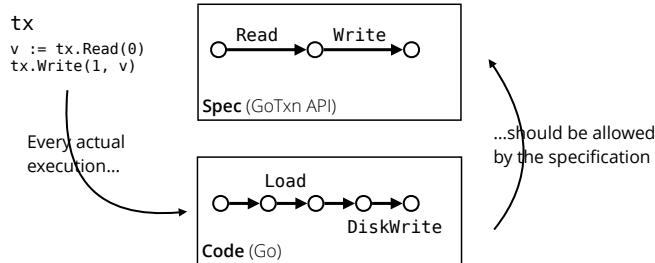
Concurrency and crash safety issues in every layer

Composing proofs of layers

Specifying the transactional API

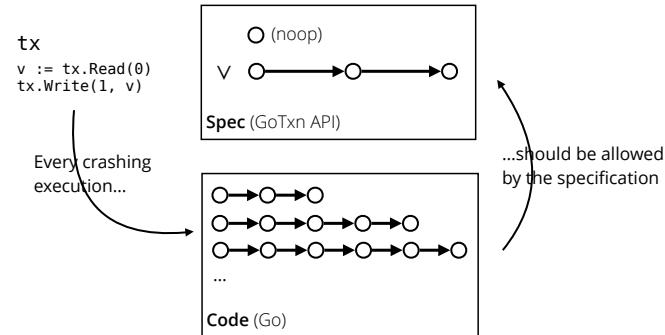
28

Specifying a transaction system



29

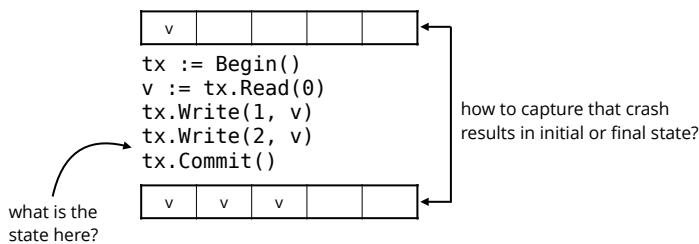
Crash atomicity for sequential transactions



30

Specifying sequential transactional API

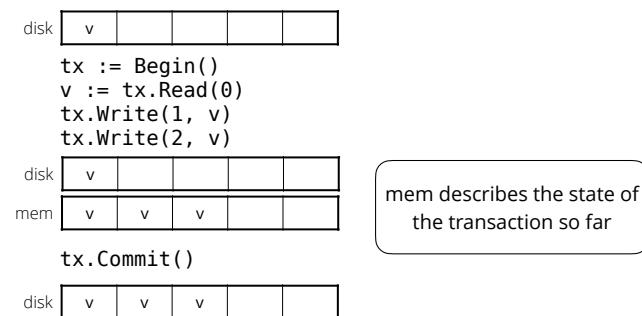
[SOSP '15]



31

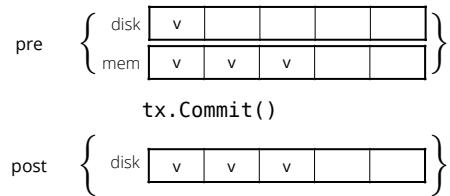
Specifying sequential transactional API

[SOSP '15]



32

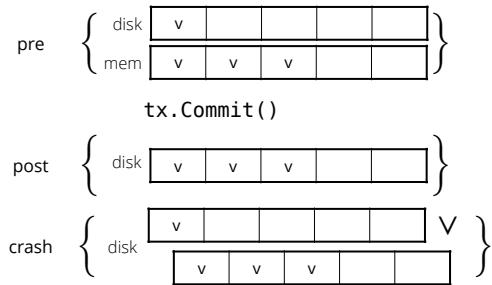
Hoare Logic to specify Commit without crashes



33

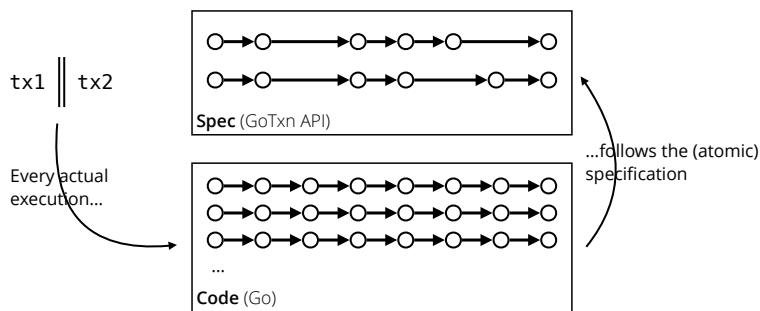
Crash Hoare Logic for sequential crash safety

[SOSP '15]



34

GoTxn also needs to handle concurrency



35

Challenge: specifying concurrent transactions

```
tx := Begin()
v := tx.Read(0)
tx.Write(1, v)
tx.Write(2, v)
tx.Commit()
```

```
||
```

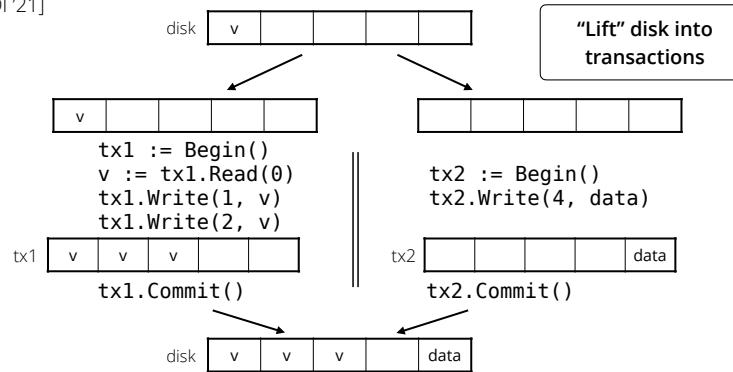
```
tx := Begin()
tx.Write(4, data)
tx.Commit()
```

How to specify effect of concurrent transactions?

36

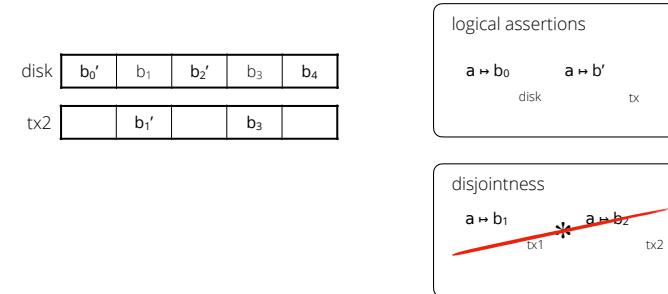
Idea: lifting-based specification

[OSDI '21]



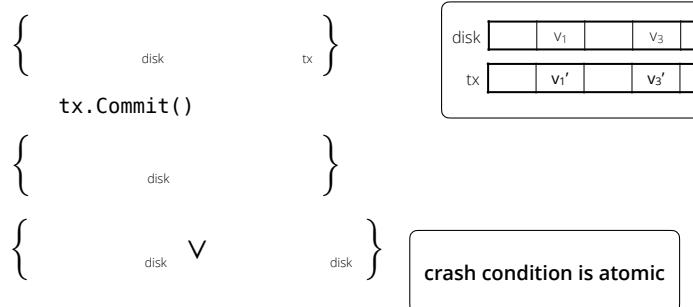
37

Separation logic describes lifting



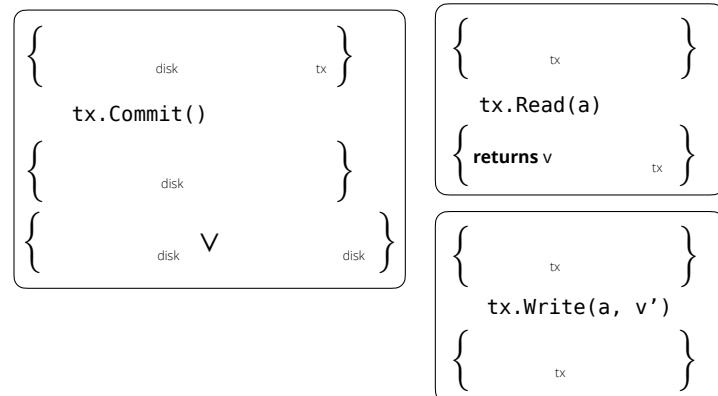
38

Commit spec captures atomicity



39

Lifting specification describes the GoTxn API



40

Perennial allows us to encode the spec



Builds on Iris [JFP 2018]

Extend with crash conditions

41

Perennial and Goose make the proof possible



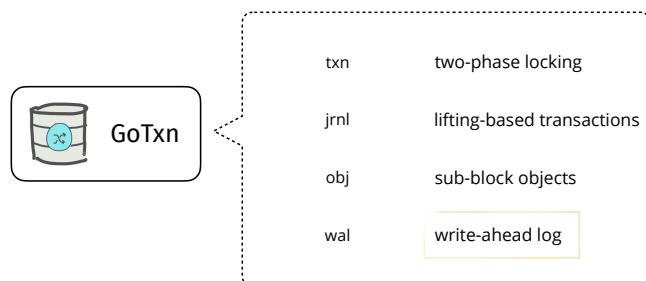
Reason about concurrency and crash safety in each layer



Connect proofs to efficient code

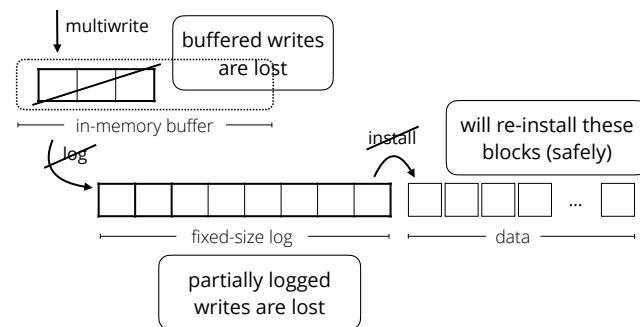
42

Excerpt from proof: focus on WAL



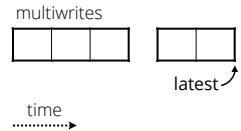
43

Crashes are complicated in the WAL



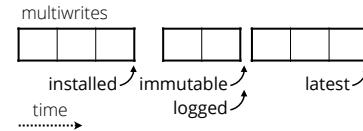
44

Idea: model WAL using history of multiwrites



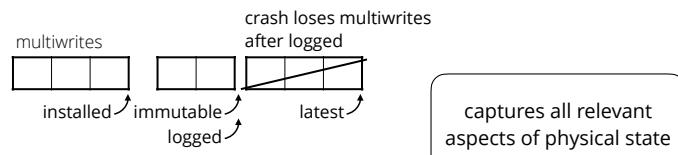
45

Idea: model WAL using history of multiwrites



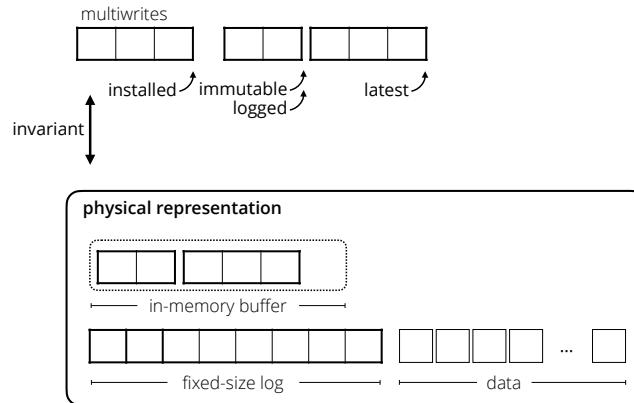
46

Crashes are easily expressed in this abstraction



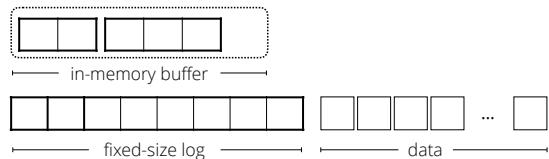
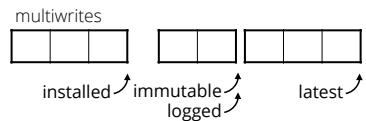
47

Invariant connects abstraction to physical representation



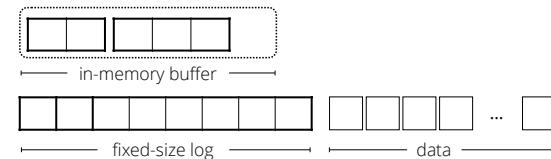
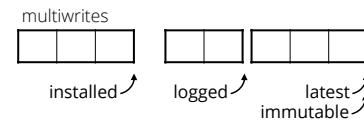
48

Pointers advance as write propagates



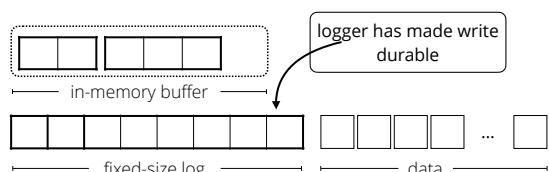
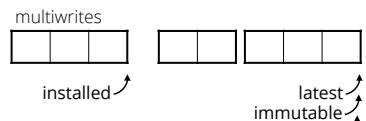
49

Write propagation: flush



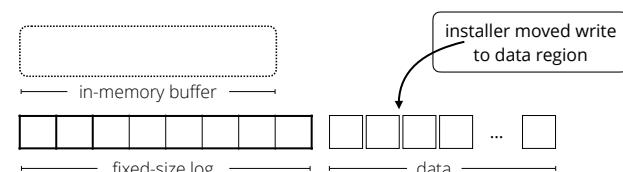
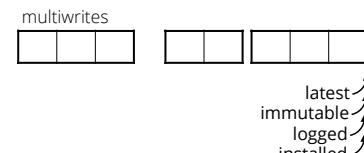
50

Write propagation: logging



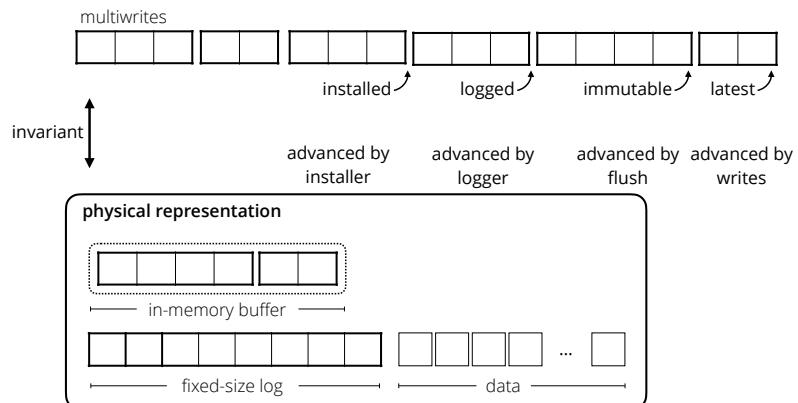
51

Write propagation: installation



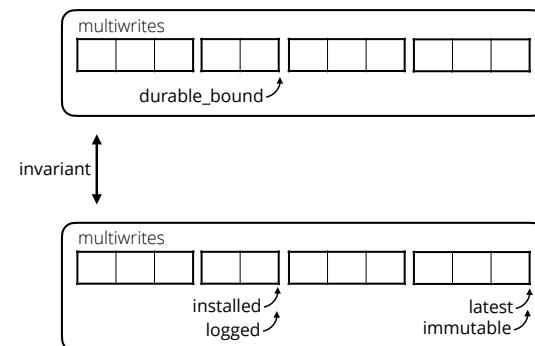
52

Pointers can all advance concurrently



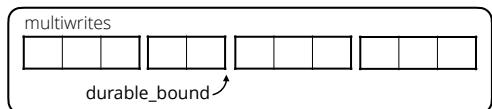
53

Further abstraction hides concurrency



54

Further abstraction hides concurrency

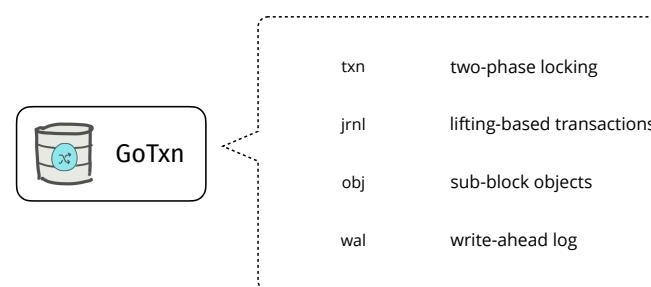


Flush() advances durable_bound

On crash keep writes before durable_bound

55

Summary



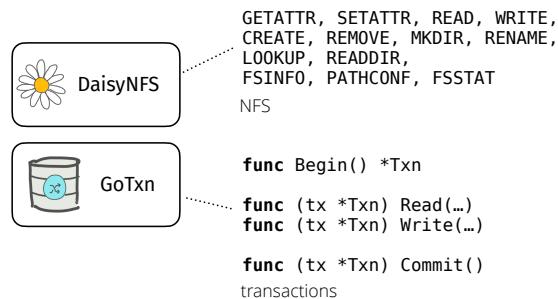
56



DaisyNFS

57

DaisyNFS is a verified file system on top of GoTxn



58

Challenges

Specification: formalizing NFS

Proof: leveraging atomicity from GoTxn

Implementation: fitting operations into transactions

59

Specification: how to formalize NFS (RFC 1813)?

INFORMATIONAL

Network Working Group
Request for Comments: 1813
Category: Informational

B. Callaghan
B. Pawlowski
P. Staubach
Sun Microsystems, Inc.
June 1995

NFS Version 3 Protocol Specification

Status of this Memo

This memo provides information for the Internet community.
This memo does not specify an Internet standard of any kind.
Distribution of this memo is unlimited.

IETF Note

Internet Engineering Steering Group comment: please note that
the IETF is not involved in creating or maintaining this
specification. This is the significance of the specification
not being on the standards track.

Abstract

This paper describes the NFS version 3 protocol. This paper is
provided so that people can write compatible implementations.

60

Define abstract state for NFS

```

type Ino = uint64
type Path = seq<byte>
datatype Attrs = Attrs(mode: uint32, ...)
datatype File =
| ByteFile(data: seq<byte>, attrs: Attrs)
| Directory(dir: map<Path, Ino>, attrs: Attrs)

type FilesysData = map<Ino, File>

```

Captures the state of the file system

61

Mathematical model of each operation

```

predicate REMOVE_spec(fs: FilesysData, fs': FilesysData,
                     args: RemoveArgs)
{
```

```

  && args.d_ino in fs
  && fs[args.d_ino].Directory?
  && var d := fs[args.d_ino].dir;
  && args.name in d

```

conditions needed for REMOVE to make sense

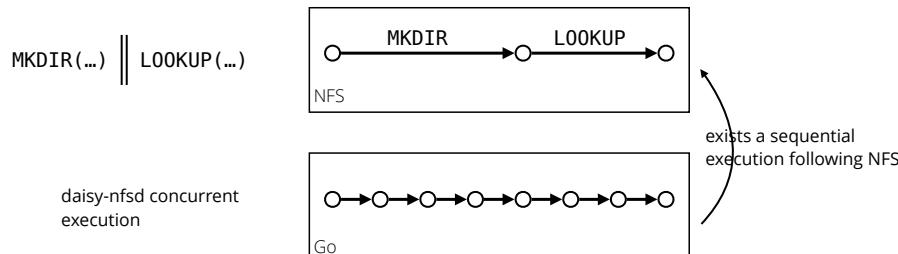
essence of REMOVE is this change to d_ino

```

  && var d' := delete(d, args.name);
  fs' == fs[args.d_ino] := Directory(d', d.attrs)
}
```

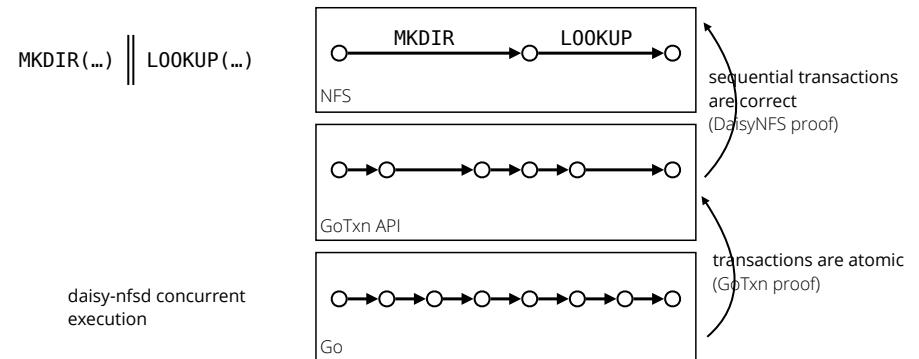
62

DaisyNFS's top-level specification



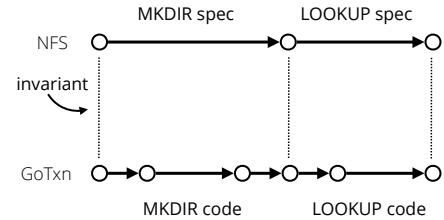
63

Proof: compose GoTxn and DaisyNFS proofs



64

Transactions are proven with sequential reasoning



65

Verify operations using Dafny

```

func REMOVE(args) RemoveReply {
    tx := Begin()
    r := fs.REMOVE(tx, args)
    tx.Commit()
    return r
}

method REMOVE(tx, args)
    returns (r: RemoveReply)
    requires Valid() ensures Valid()
    ensures
        REMOVE_spec(old(fs), fs, args, r)

```

Dafny (verified)

Go (unverified)

66

Implementing freeing using bounded transactions

```

func REMOVE(args) RemoveReply {
    tx := Begin()
    r := fs.REMOVE(tx, args)
    tx.Commit()
    go fs.ZeroFreeSpace(r.ino)
    return r
}

method REMOVE(tx, args)
    returns (r: RemoveReply)
    requires Valid() ensures Valid()
    ensures
        REMOVE_spec(old(fs), fs, args, r)

method ZeroFreeSpace(ino)
    requires Valid() ensures Valid()
    ensures fs == old(fs)

```

Dafny (verified)

Go (unverified)

Logical no-op but physically recovers space

67

Summary



Specified a mathematical version of RFC 1813



Proof strategy uses GoTxn atomicity

Implemented file system using transactions

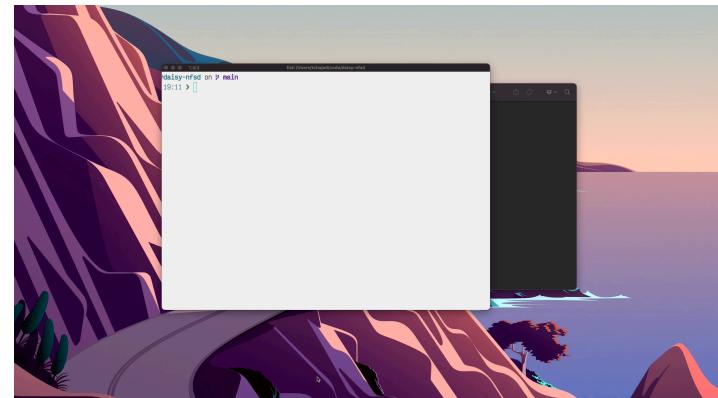
68

Implementation

	Total lines
DaisyNFS	9,600 (Dafny)
GoTxn	36,000 (Perennial)
Perennial	20,000 (Coq)
Goose	3,500 (Go)

69

DaisyNFS is a real file system



70

Evaluation

71

Evaluation questions

Does GoTxn reduce the proof burden?

What is assumed in the DaisyNFS proof?

Does DaisyNFS get acceptable performance?

72

GoTxn greatly reduces proof overhead

	Code	Proof	
DaisyNFS	3,200	6,400	2x proof:code
GoTxn	1,600 (Go)	35,000 (Perennial)	20x proof:code

73

Assumptions in the DaisyNFS proof



Theorem: Every NFS operation appears to execute atomically and correctly, despite crashes and concurrency.

74

Trusted computing base (TCB)

- Unverified code (e.g., network protocol)
- Dafny specification of NFS
- GoTxn spec in Dafny is faithful
- Goose accurately models Go

75

Testing methodology

- Unit tested Dafny support libraries
- Compare Goose model to Go executions
- Ran fsstress and fsx-linux test suites
- Ran CrashMonkey for crash testing [OSDI '18]
- Symbolic execution over NFS protocol

76

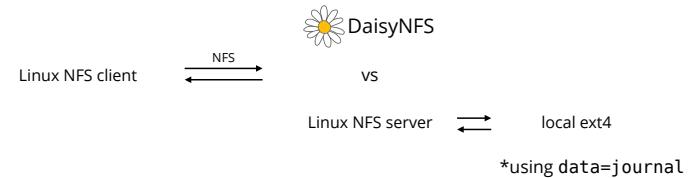
Bugs found in unverified code and spec

XDR decoder for strings can allocate 2^{32} bytes
File handle parser panics if wrong length
Panic on unexpected enum value
WRITE panics if not enough input bytes
Didn't find bugs in verified parts
Directory REMOVE panics in dynamic type cast
The names “.” and “..” are allowed
RENAME can create circular directories
CREATE/MKDIR allow empty name
Proof assumes caller provides bounded inode
RENAME allows overwrite where spec does not

Unverified glue code
Missing from specification

77

Compare against Linux NFS server with ext4



*using data=journal

78

Performance evaluation setup

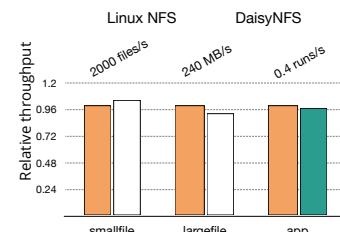
Hardware: i3.metal instance
36 cores at 2.3GHz, NVMe SSD

Benchmarks:

- smallfile: metadata heavy
- largefile: lots of data
- app: git clone + make

79

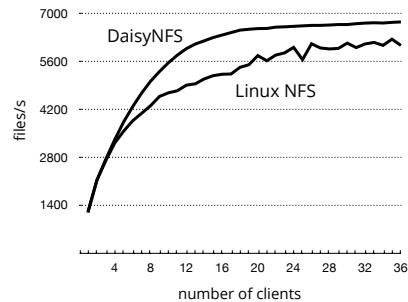
DaisyNFS gets good performance with a single client



Compare DaisyNFS throughput to Linux,
running on an in-memory disk

80

DaisyNFS can take advantage of multiple clients



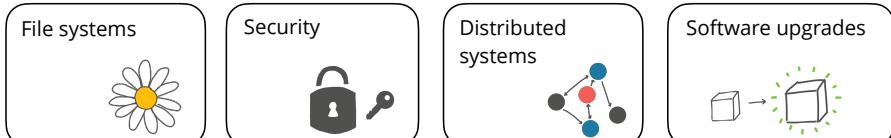
Run smallfile with many clients on an NVMe SSD

81

Future directions

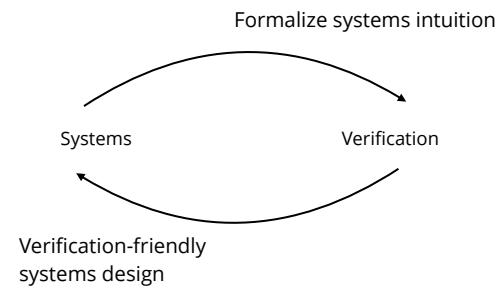
82

Imagine a future of reliable systems software



83

My research approach



84

Verification is promising for secure systems



System needs to cover all inputs but attacker only needs one vulnerability

Some verified libraries exist but not systems

85

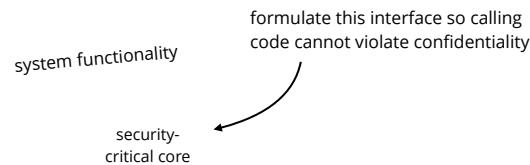
Goal: end-to-end security guarantees

Specifying application-level security

Prove confidentiality and integrity guarantees

86

Approach: specifications for confidentiality



87



Distributed systems are complex and multi-tiered

The Washington Post

Amazon Web Services' third outage in a month exposes a weak point in the Internet's backbone

The disruptions affect millions of people on an increasingly interconnected Web: "We are putting more eggs into fewer and fewer baskets. More eggs get broken that way."

Cloud services

replication + sharding

storage

GOOGLE SAYS YOUTUBE, GMAIL AND OTHERS WERE DOWN BECAUSE OF 'AN INTERNAL STORAGE QUOTA ISSUE'

CLOUDFLARE

Understanding How Facebook Disappeared from the Internet

88

Goal: prove a multi-service distributed system

Real systems are made from many smaller components

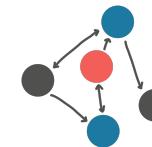
Proof burden should be manageable

89

Approach: alternate protocol proofs to RPC-style implementations

```
reply, err := RemoteCall(...)  
if err != nil {  
    if timeout(err) { ... }  
    // handle other errors  
}  
// normal processing  
  
Implementation uses  
remote procedures
```

Specification should
bridge the two



Protocol reasoning is over
whole system

90

Verifying safety of software upgrades

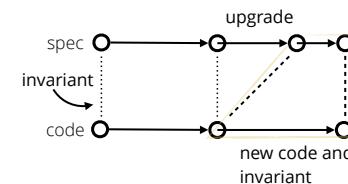


Upgrading a verified system deserves proof

File systems, distributed systems, and applications have
upgrades

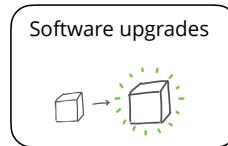
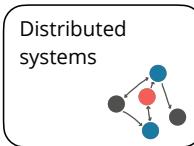
91

Approach: new specifications and proof patterns for software upgrades



92

Verification beyond these domains



93

Vision: reliable systems software with verification



DaisyNFS



GoTxn



Goose



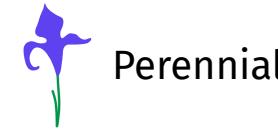
Perennial

Verified a concurrent, crash-safe file system

Expand the reach of verification in other domains

94

95



Perennial

96