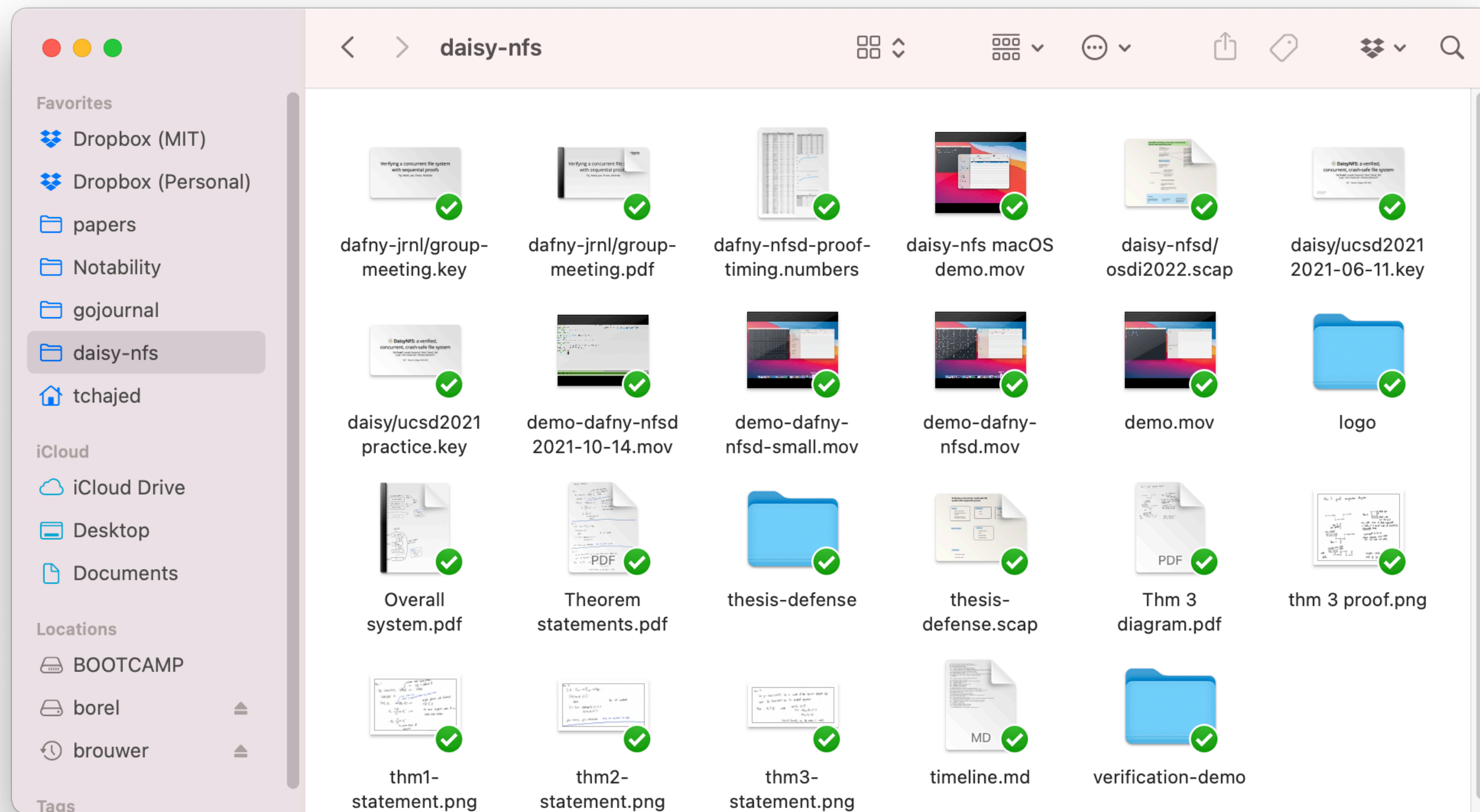


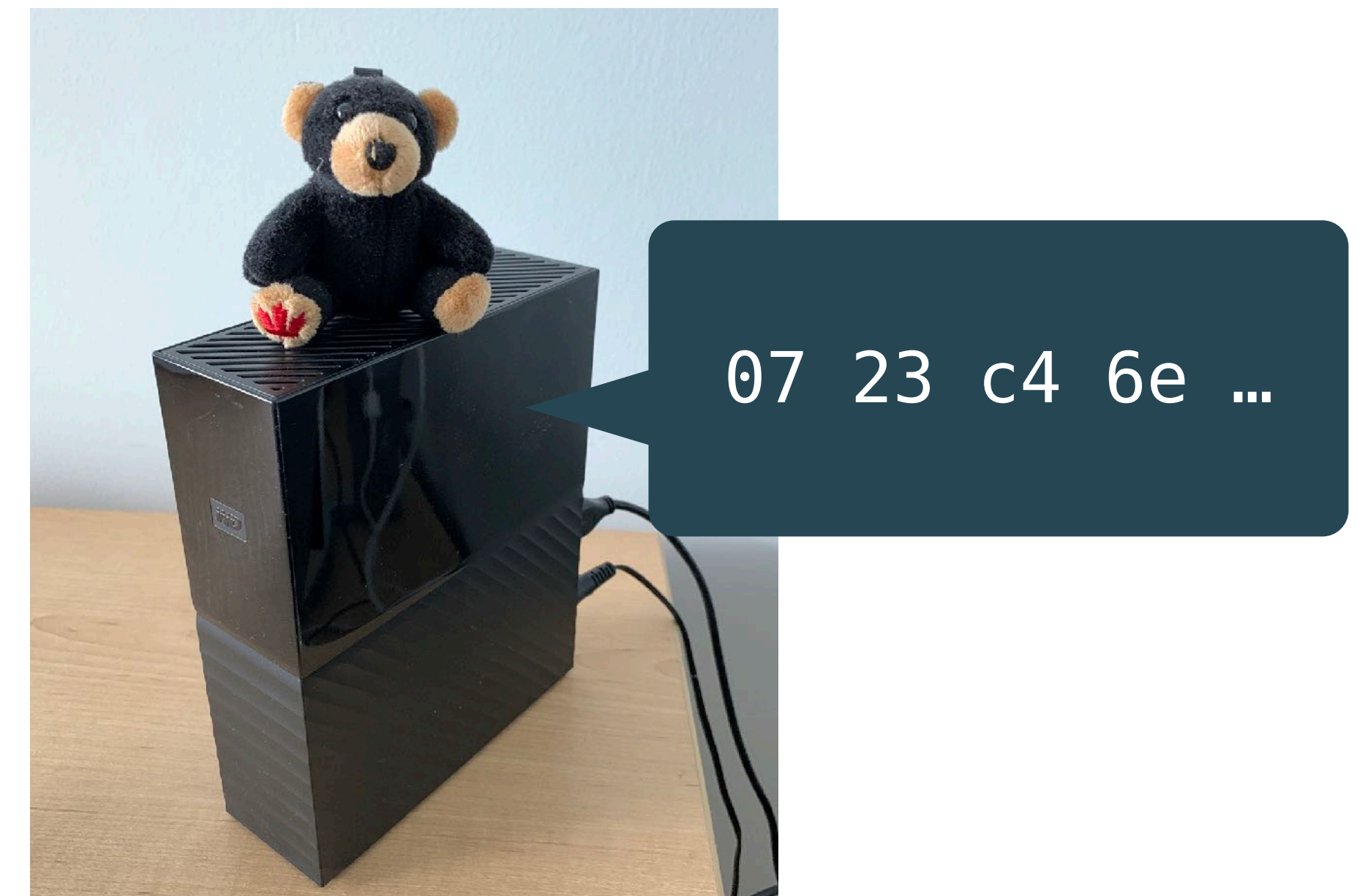
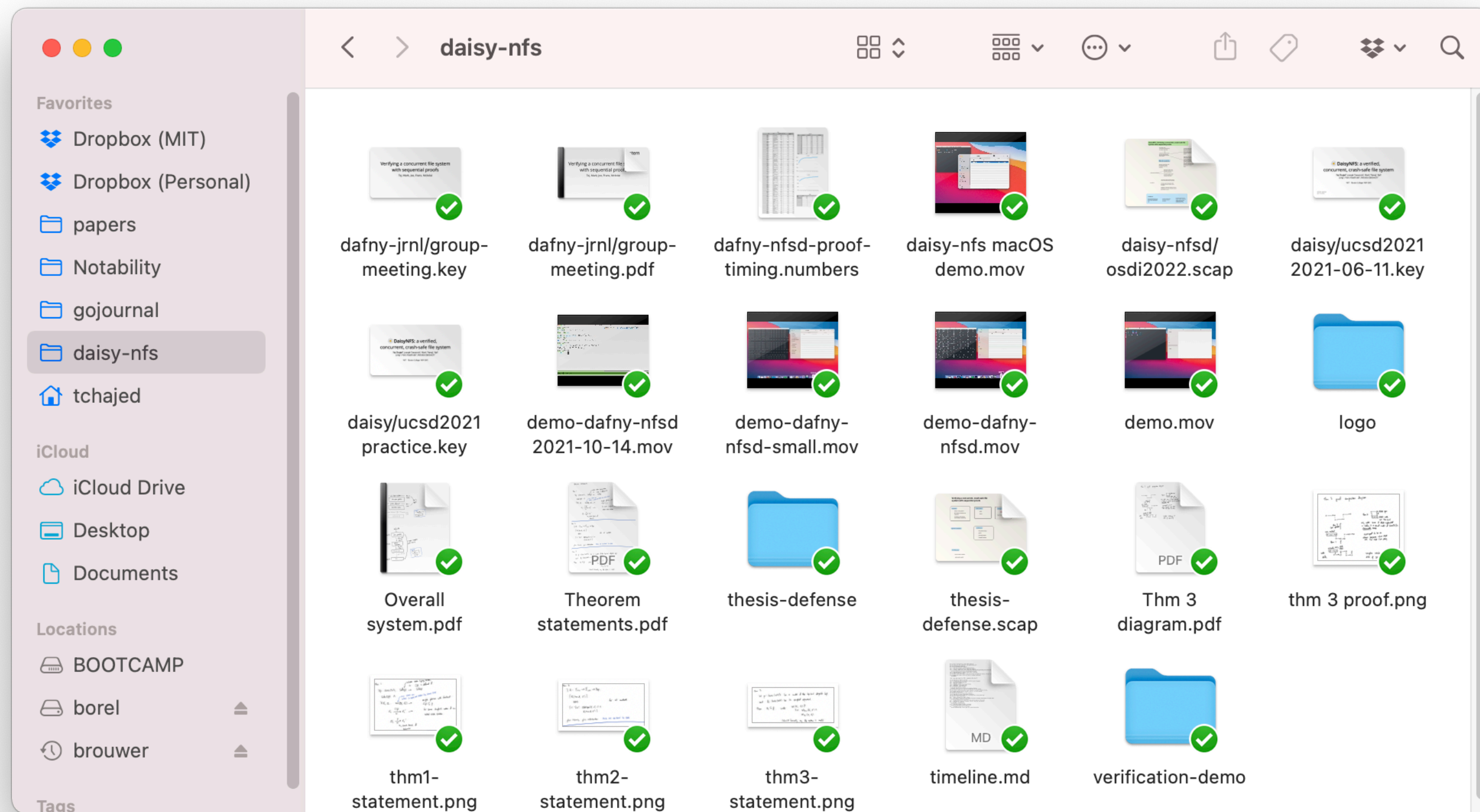


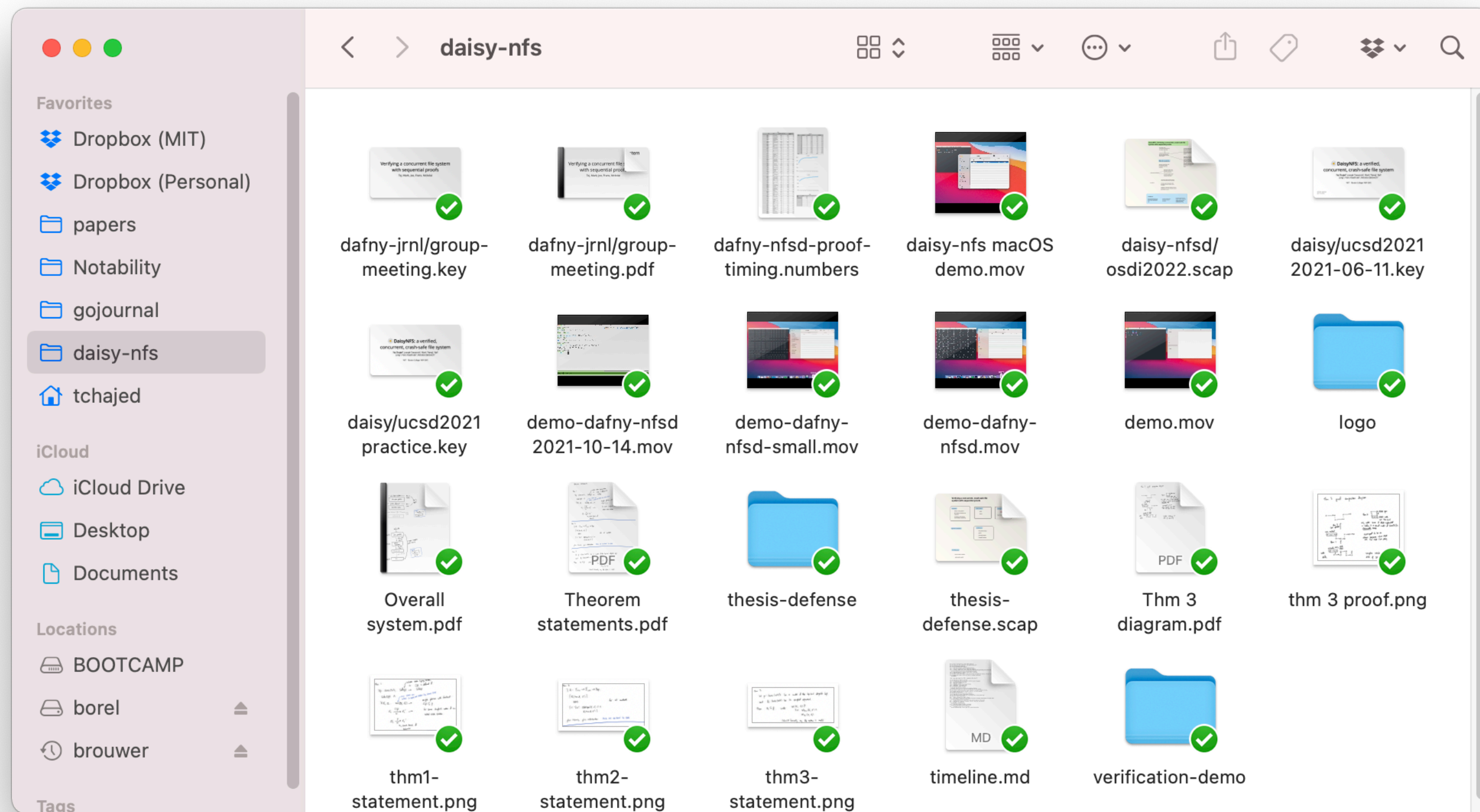
Verifying a concurrent, crash-safe file system with sequential reasoning

Tej Chajed
PhD defense

October 21st, 2021



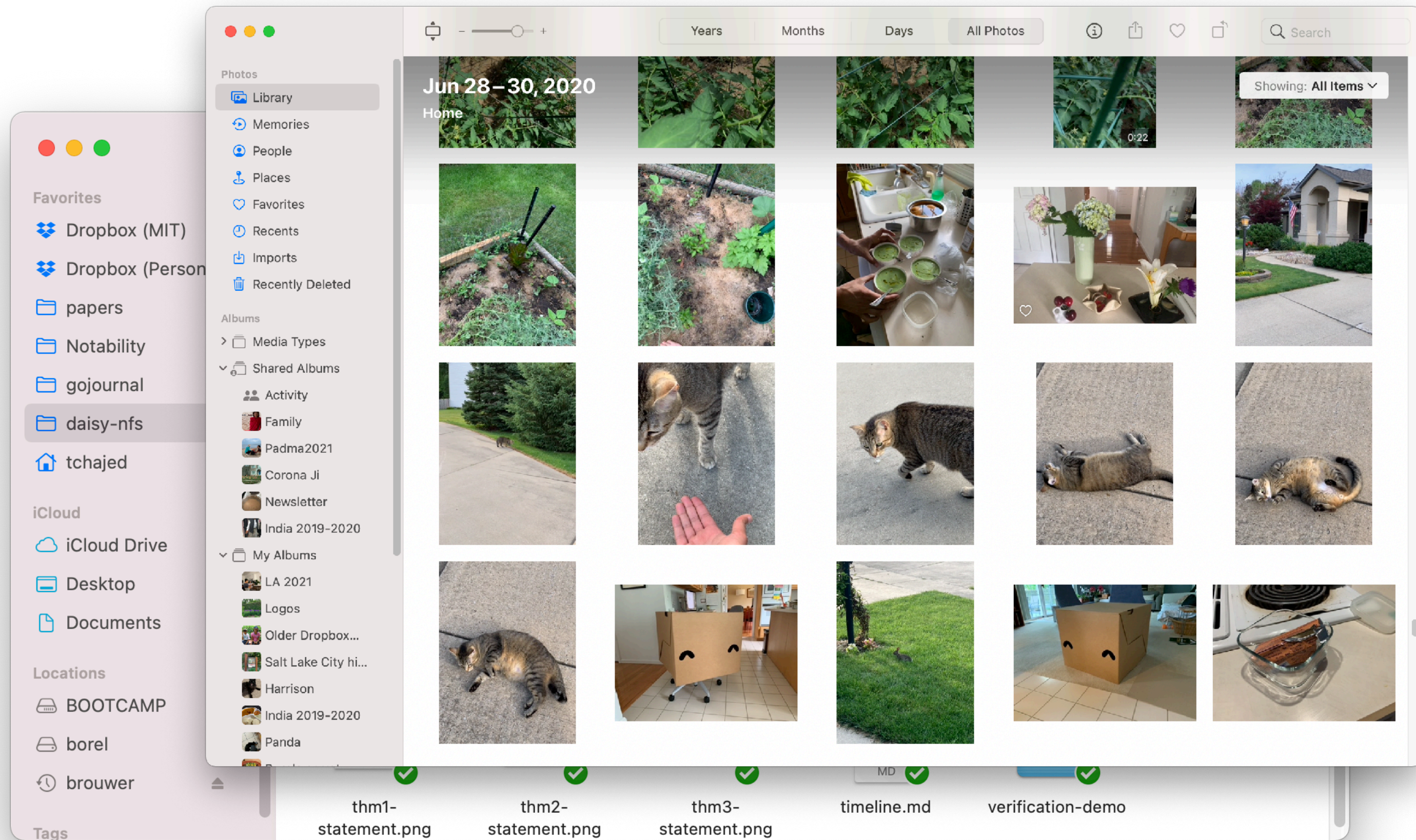




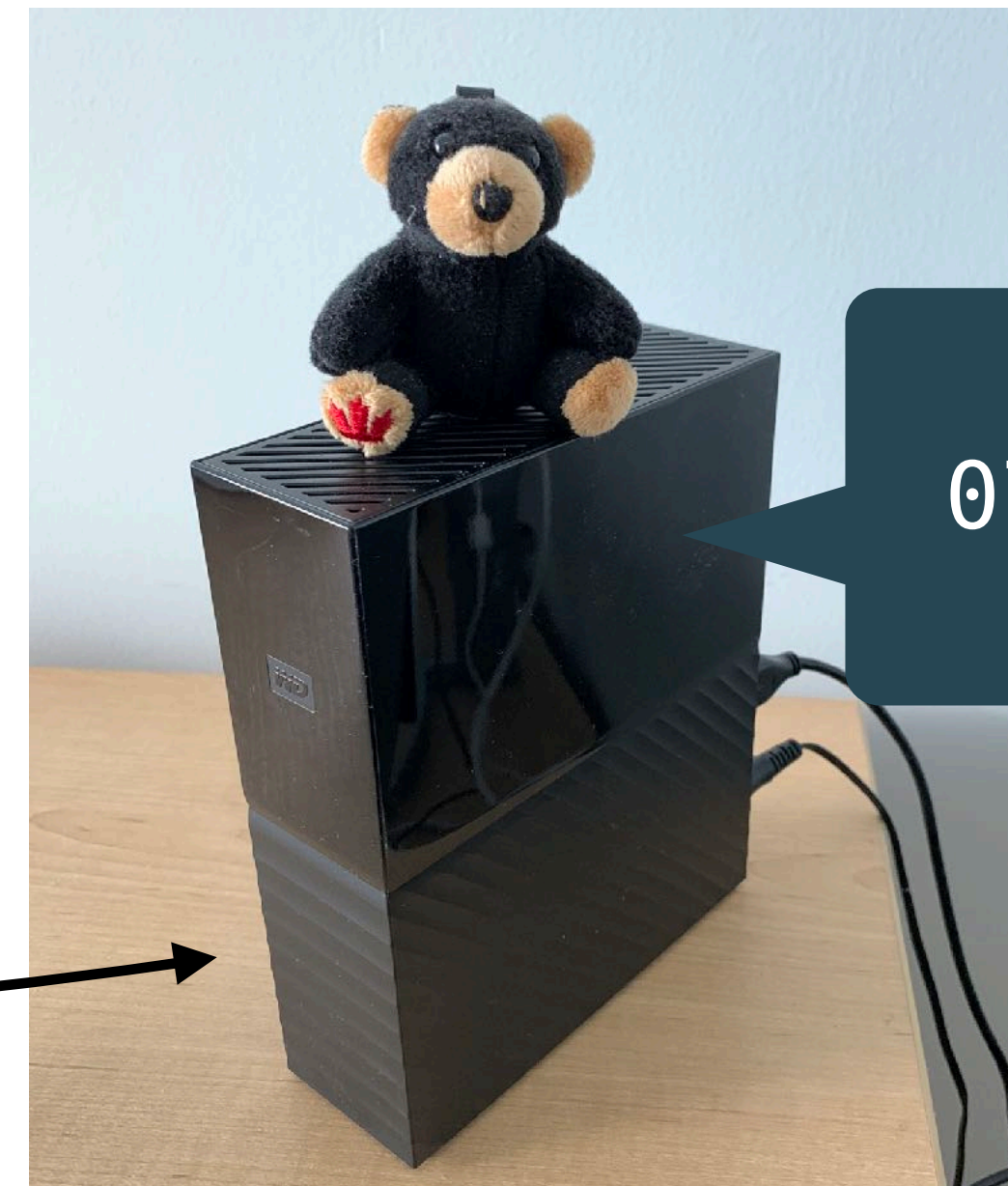
file system

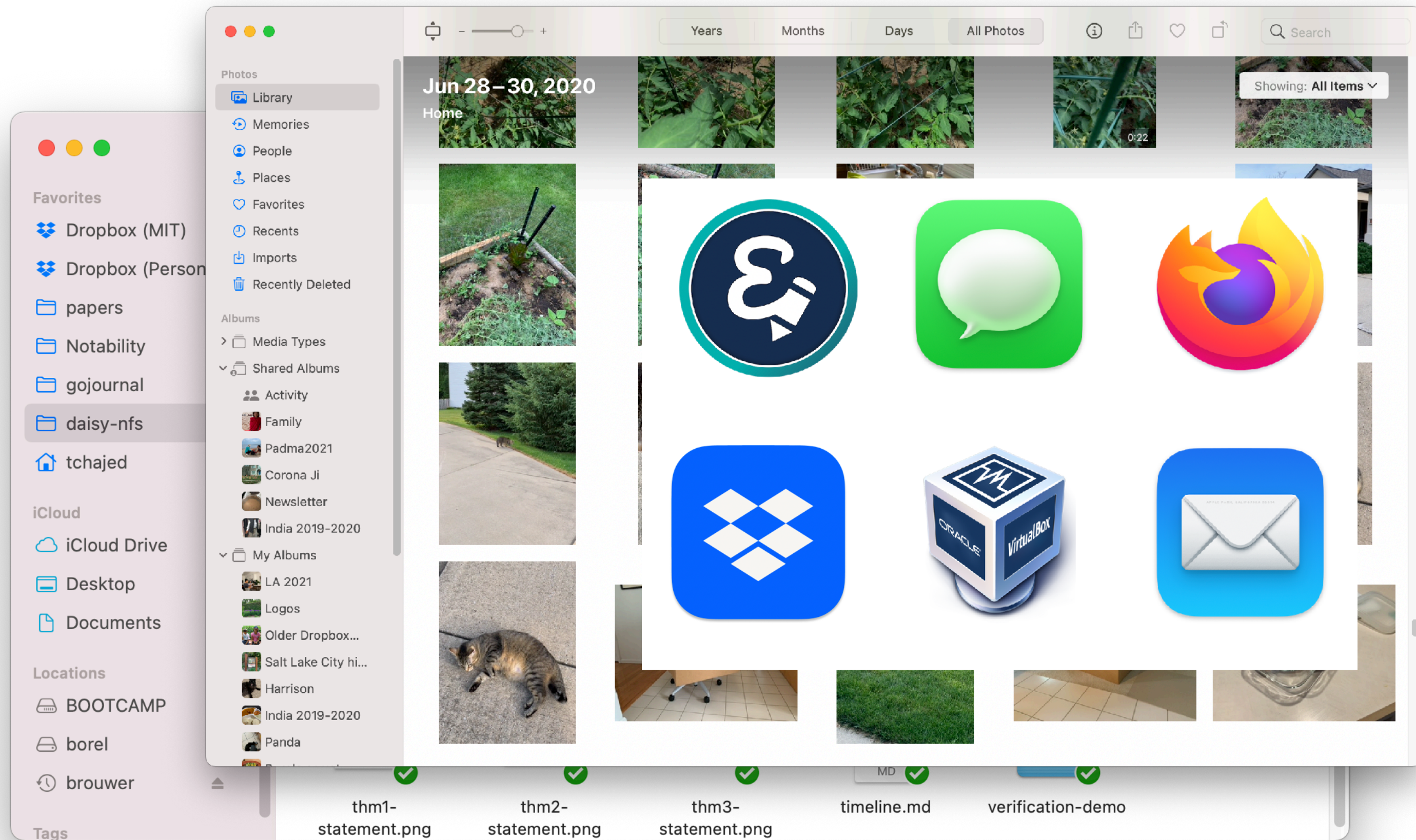


07 23 c4 6e ...

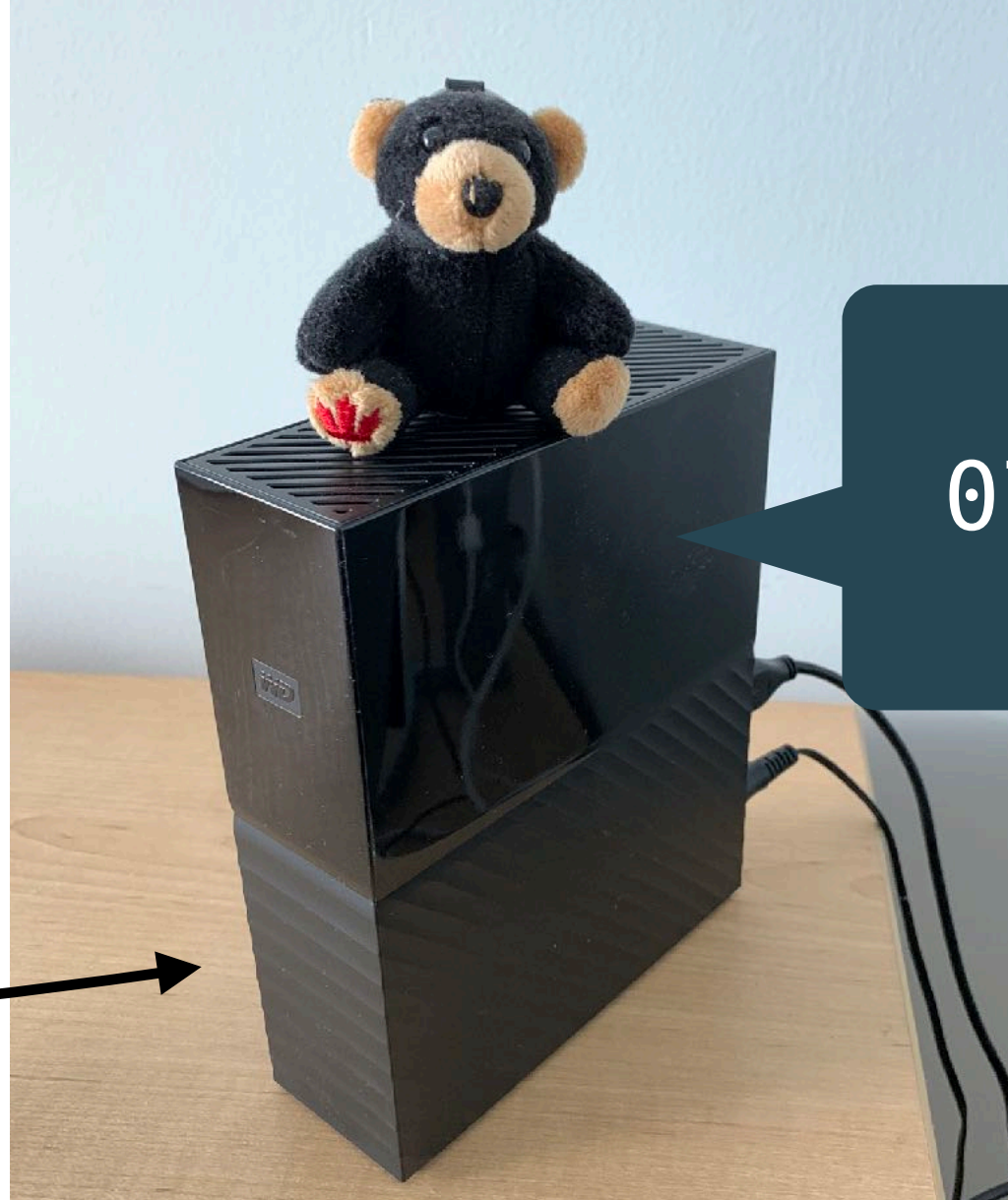


file system





file system



07 23 c4 6e ...

Important that the file system be correct

Responsible for storing all persistent data

Nearly all applications rely on it

Bugs can cause permanent data loss

File systems are just programs and
therefore **they have bugs**



File systems struggle with
crash safety + concurrency + high performance

File systems struggle with
crash safety + concurrency + high performance

a crash is any sudden interruption, like a power failure

File systems struggle with **crash safety + concurrency + high performance**

a crash is any sudden interruption, like a power failure

concurrency comes from devices, simultaneous user requests

File systems struggle with **crash safety + concurrency + high performance**

a crash is any sudden interruption, like a power failure

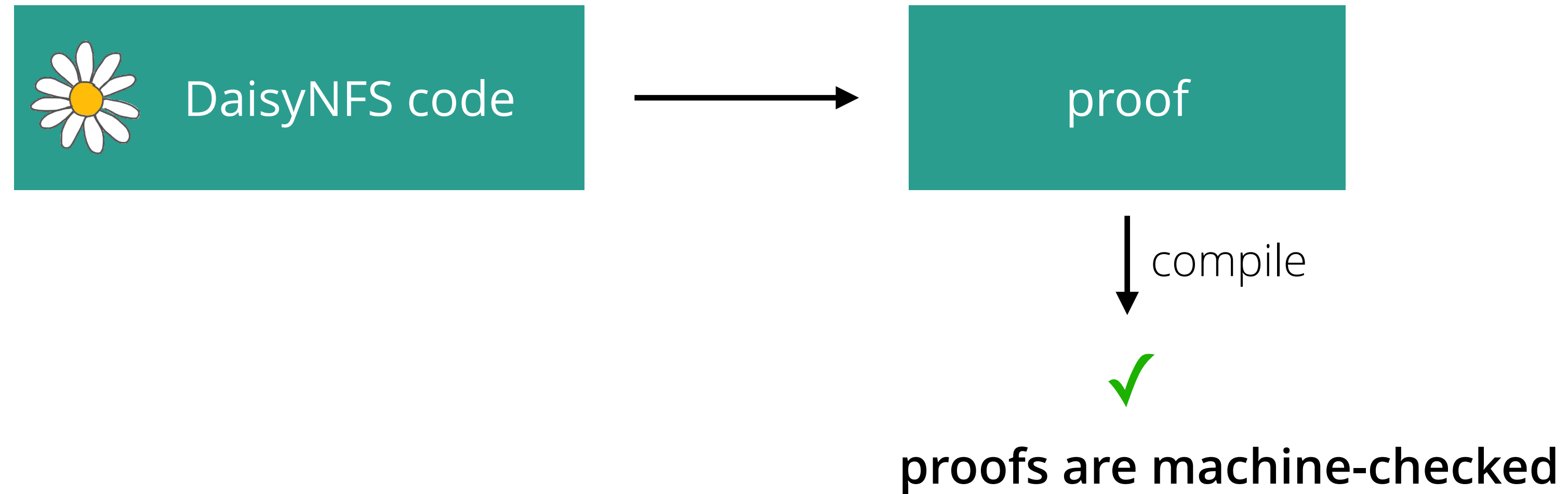
concurrency comes from devices, simultaneous user requests

high performance makes both of these hard

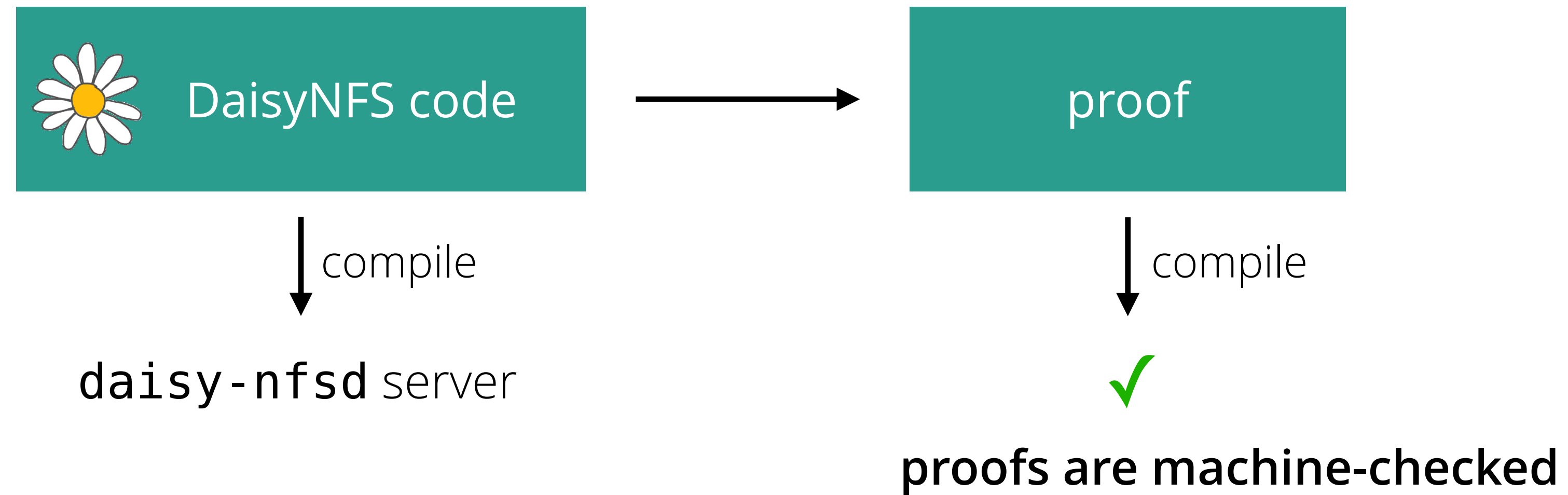
DaisyNFS is a new, verified file system



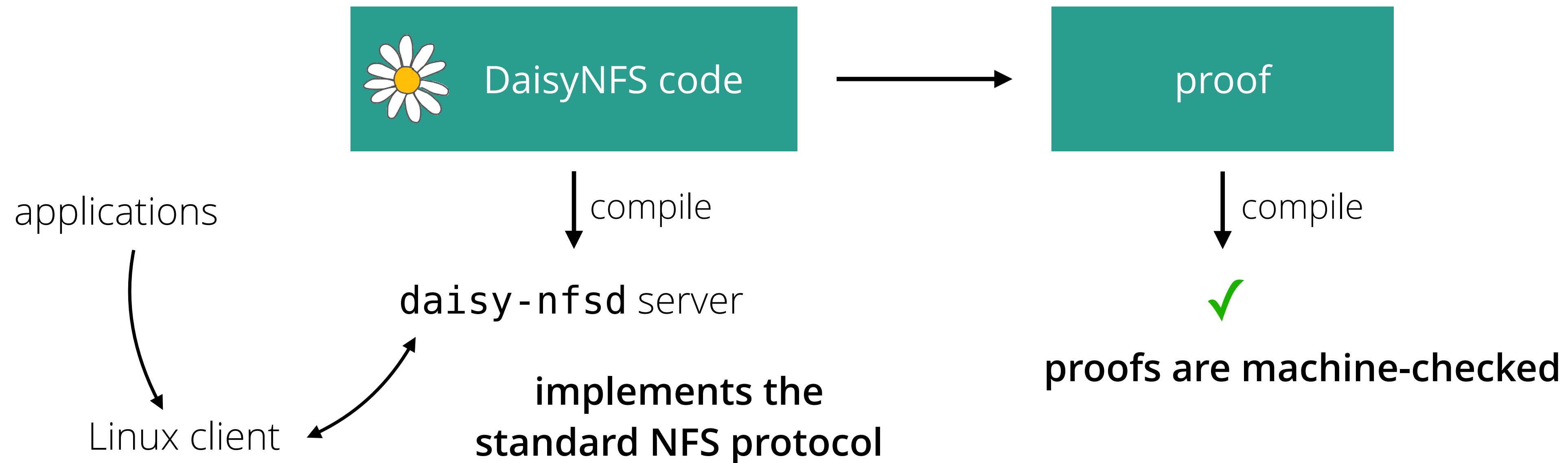
DaisyNFS is a new, verified file system



DaisyNFS is a new, verified file system



DaisyNFS is a new, verified file system



DaisyNFS is a real file system



```
demo:1:fish - "fish /home/tchajed"
@hamilton ~
16:31 > |
drwxr-xr-x - 8 Jun 12:44 /mnt/nfs

[demo] 1:fish* 16:31 Oct 20
```

The image shows a terminal window with a title bar that reads "demo:1:fish - 'fish /home/tchajed'". The terminal content is split into two columns. The left column shows the prompt "@hamilton ~" and the time "16:31" followed by a green arrow cursor and a black block cursor. The right column shows the output of a command: "drwxr-xr-x - 8 Jun 12:44 /mnt/nfs". At the bottom of the terminal, a green status bar displays "[demo] 1:fish*" on the left and "16:31 Oct 20" on the right.

DaisyNFS is a real file system



```
demo:1:fish - "fish /home/tchajed"
@hamilton ~
16:31 > |
drwxr-xr-x - 8 Jun 12:44 /mnt/nfs

[demo] 1:fish* 16:31 Oct 20
```

The image shows a terminal window with a light gray background and a dark green status bar at the bottom. The window title is "demo:1:fish - 'fish /home/tchajed'". The prompt is "@hamilton ~". The time is "16:31" and there is a cursor. The output shows the permissions "drwxr-xr-x", a dash "-", the size "8", the date "8 Jun", the time "12:44", and the path "/mnt/nfs". The status bar at the bottom shows "[demo] 1:fish*" on the left and "16:31 Oct 20" on the right.

Approach: formal verification

Give **mathematical proof** that code does what it's supposed to

Formalize desired behavior as a **specification**

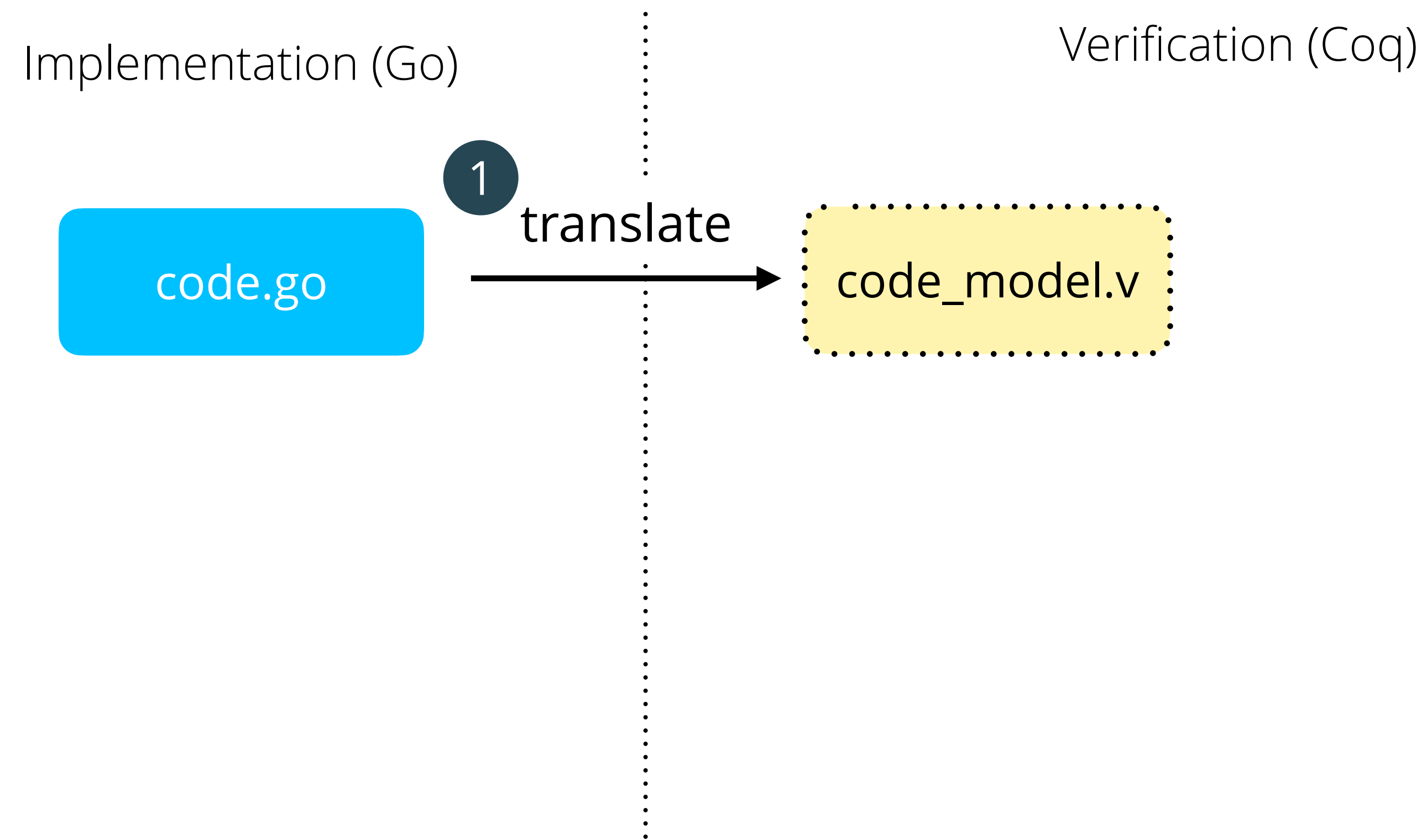
Approach: formal verification

Implementation (Go)

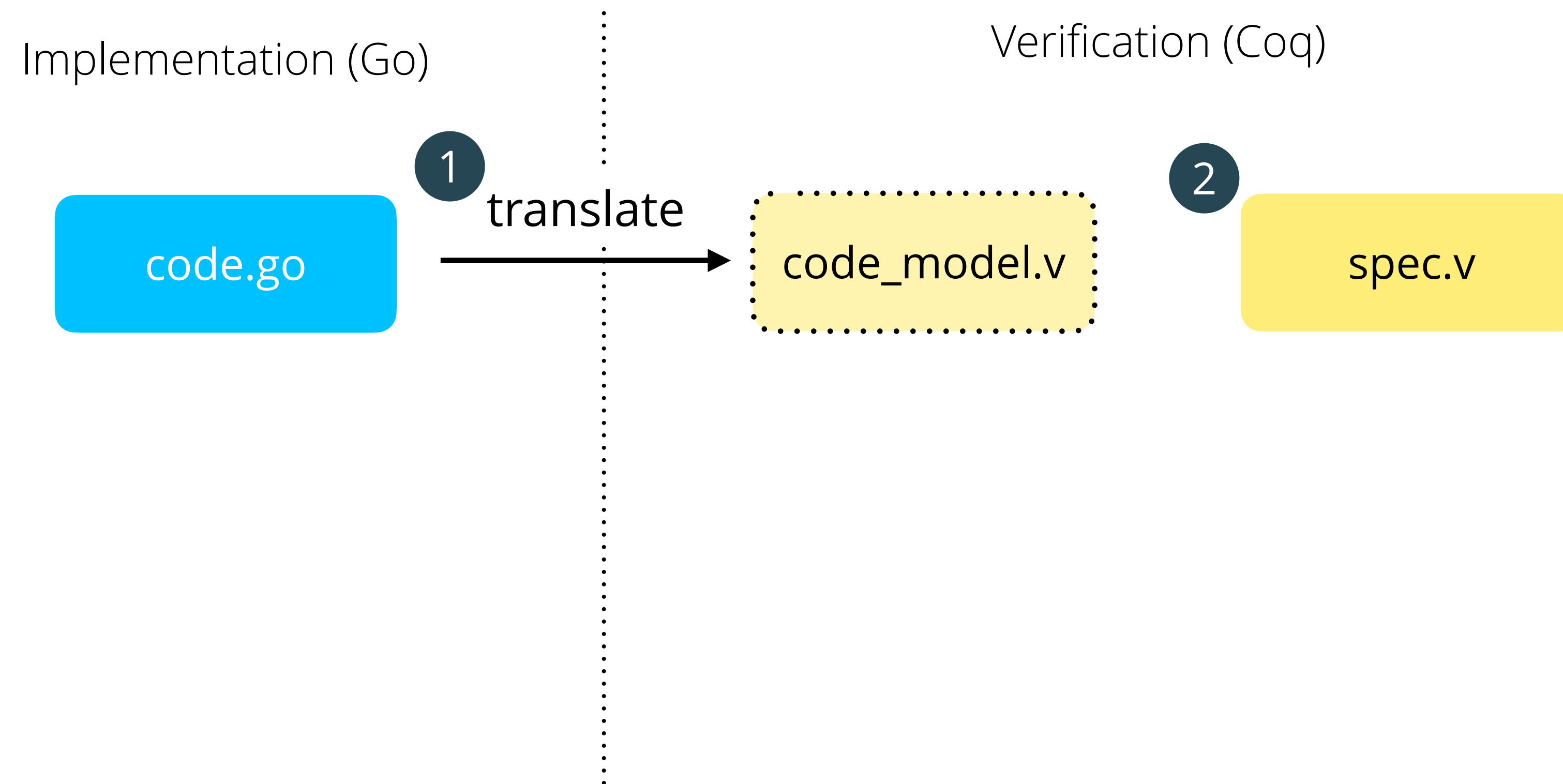


Verification (Coq)

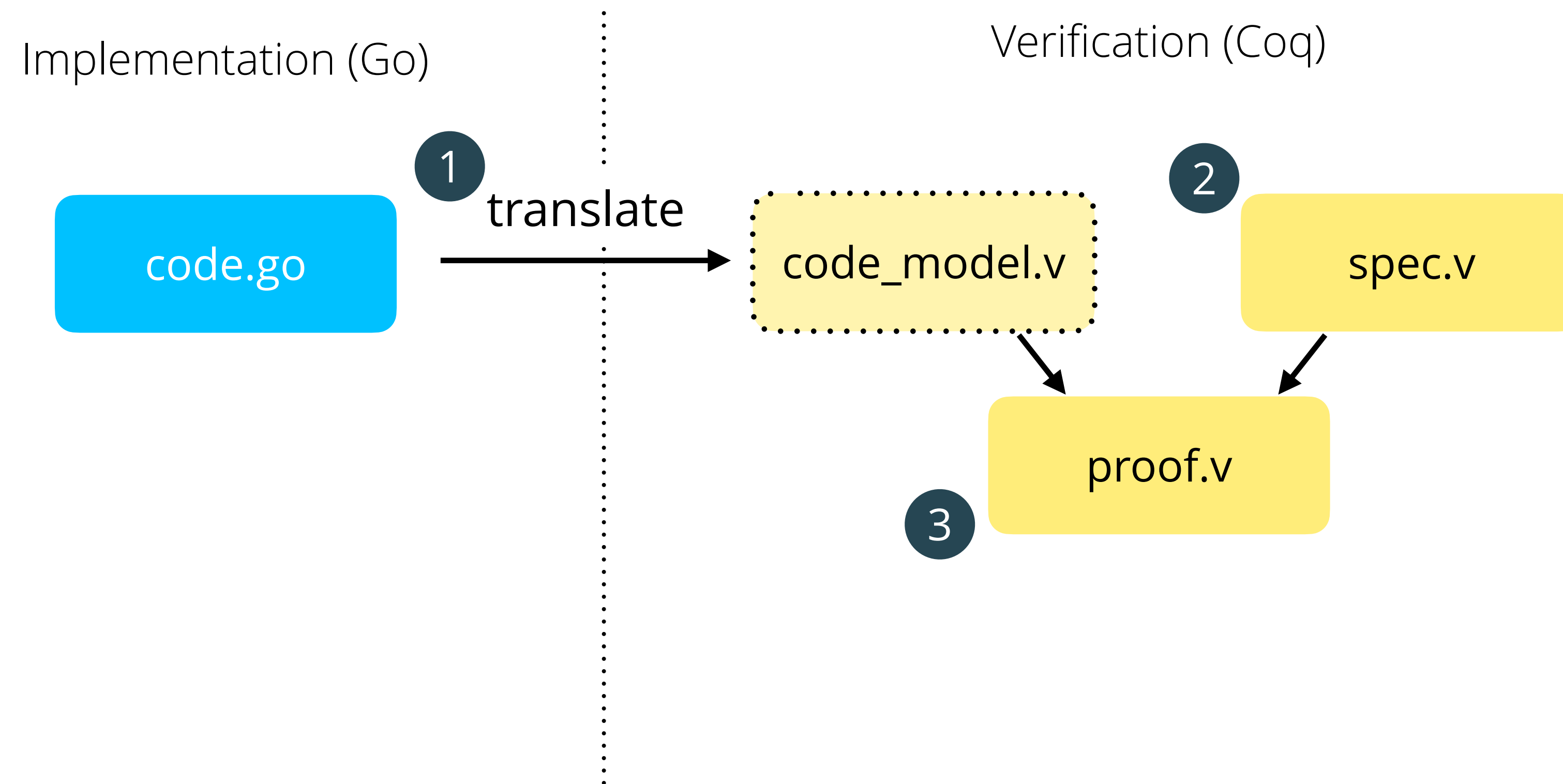
Approach: formal verification



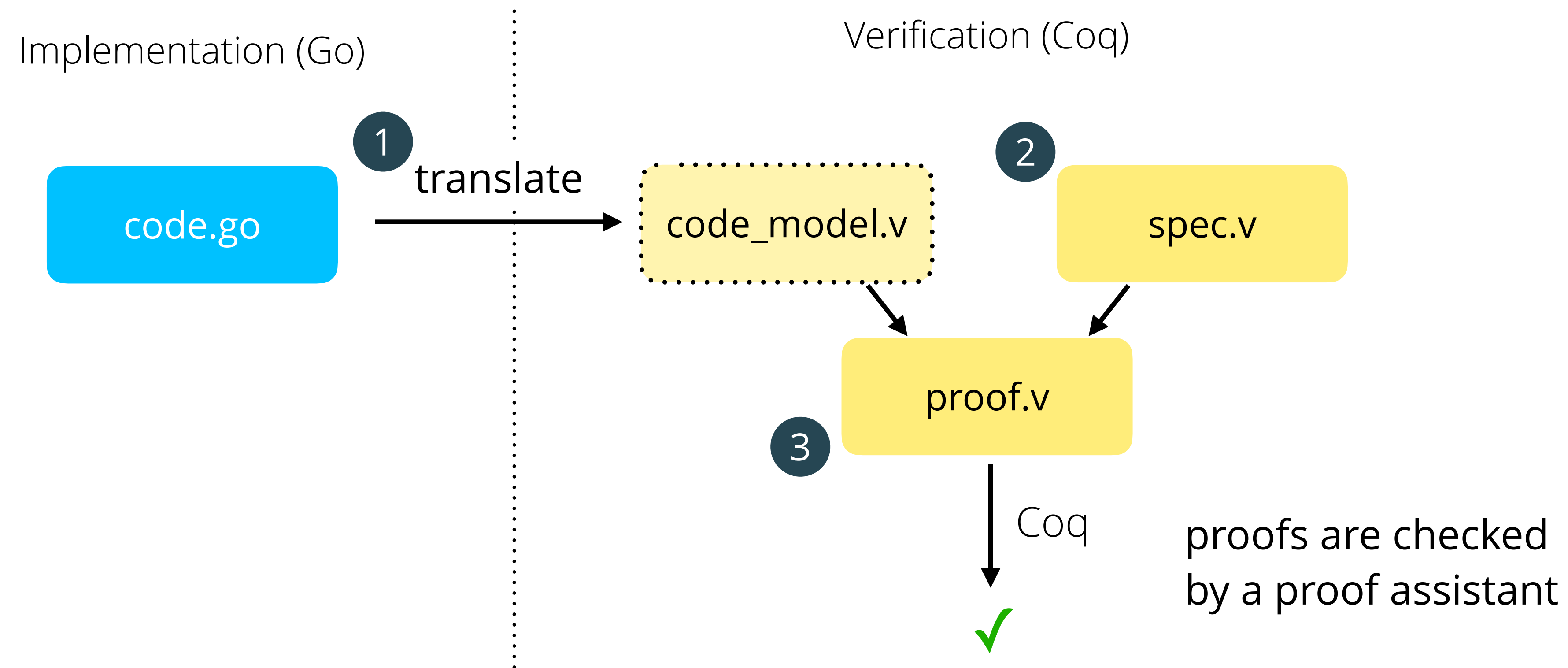
Approach: formal verification



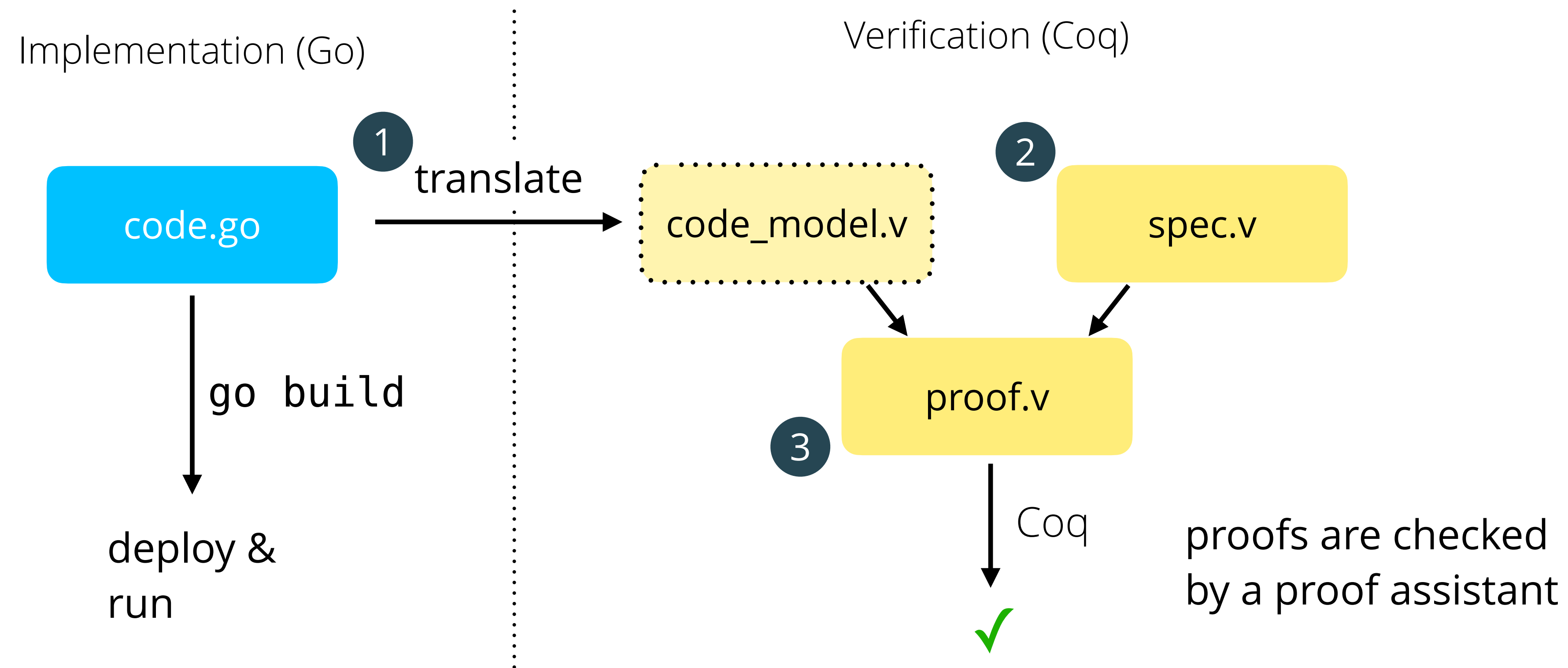
Approach: formal verification



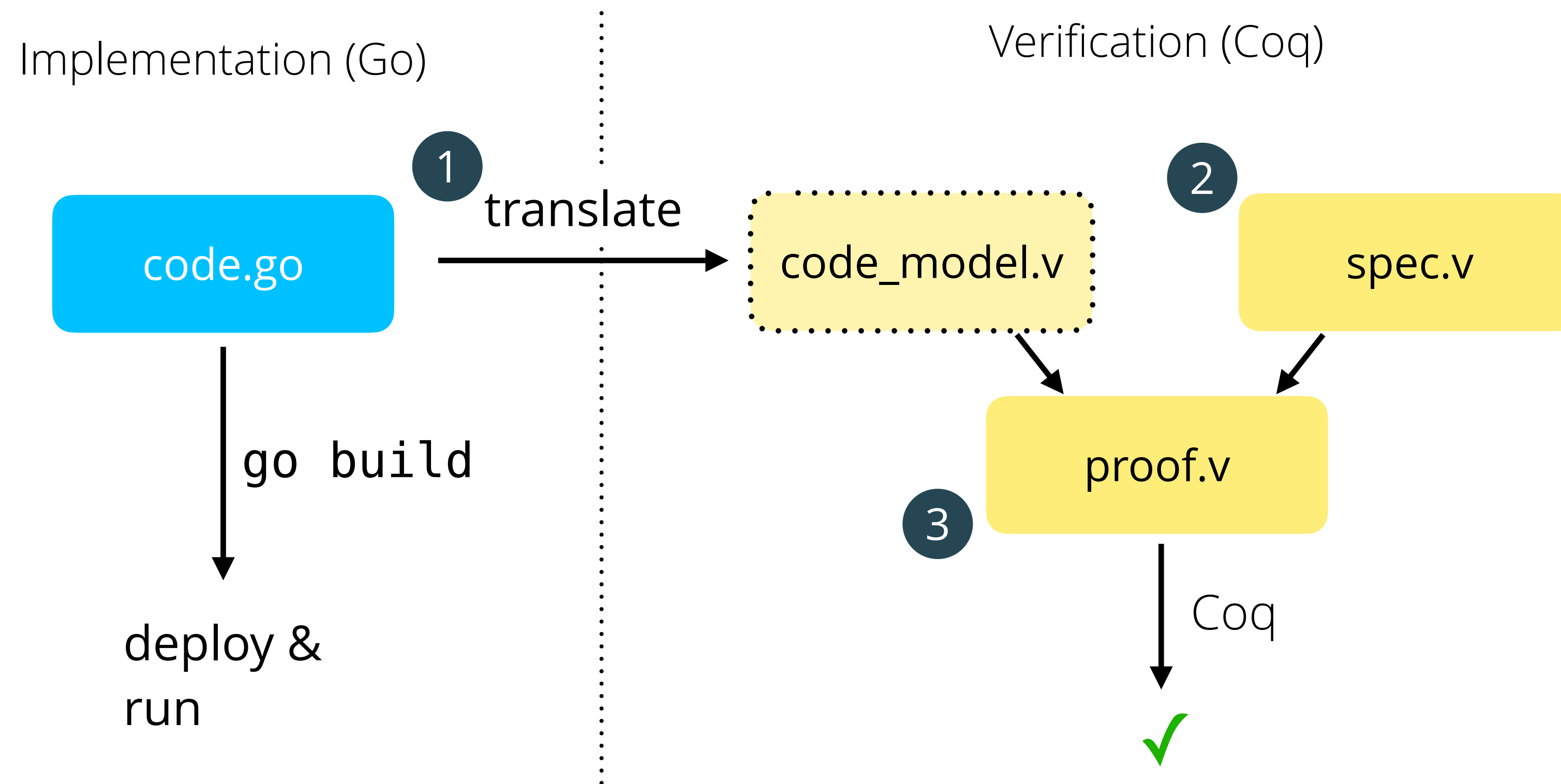
Approach: formal verification



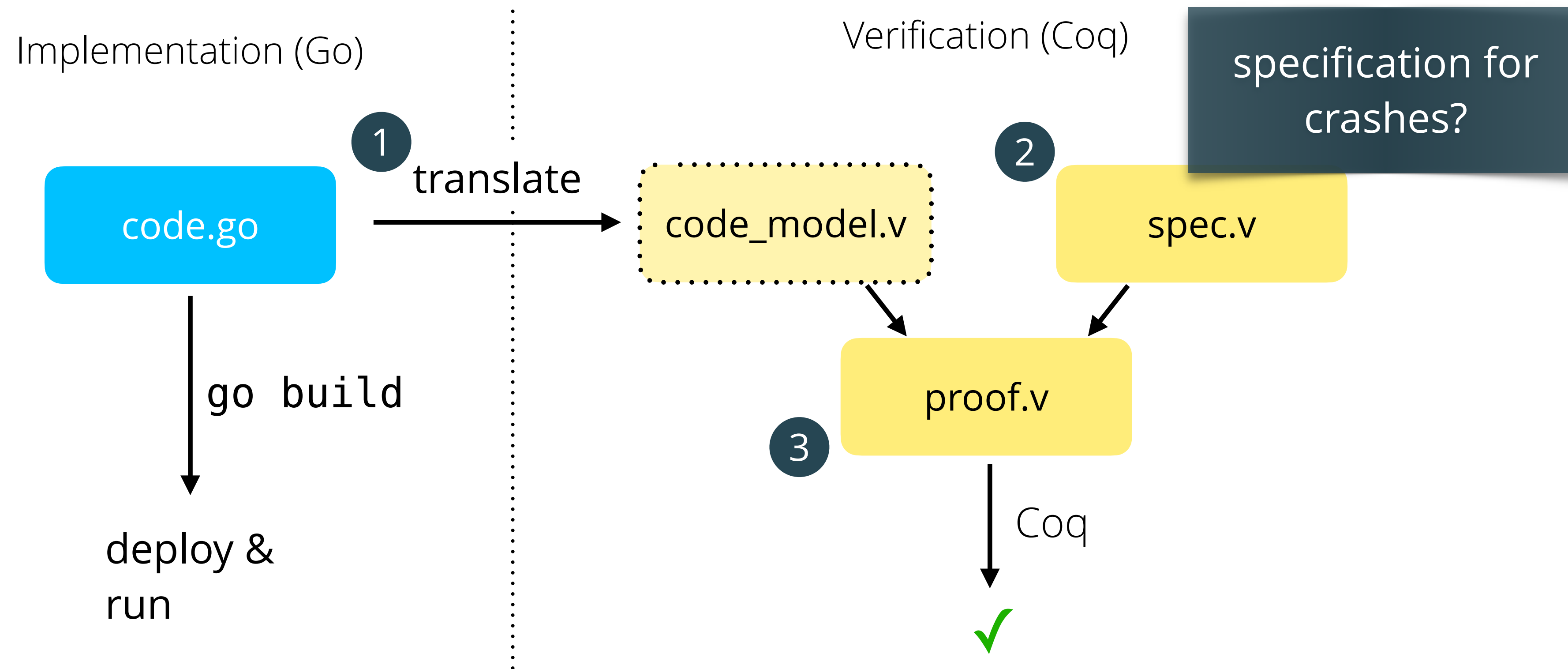
Approach: formal verification



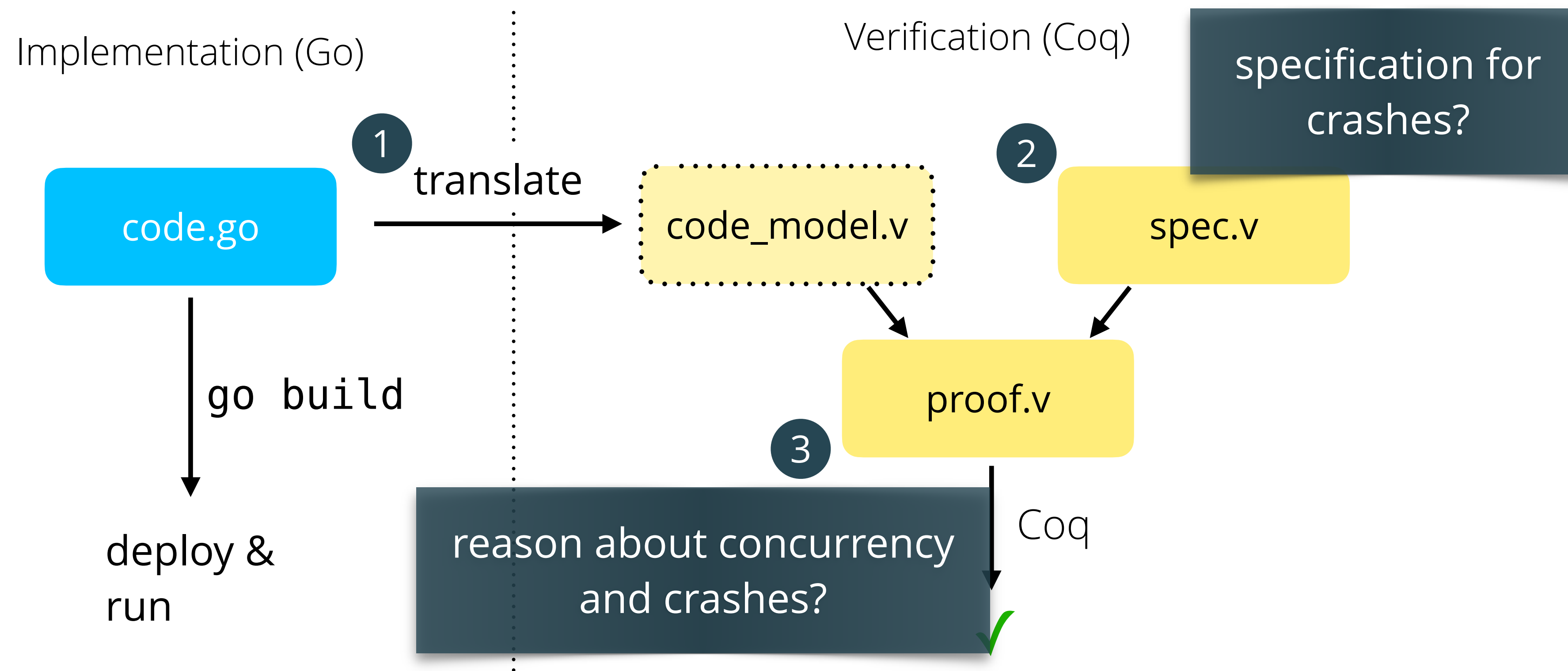
Verification of storage systems is challenging



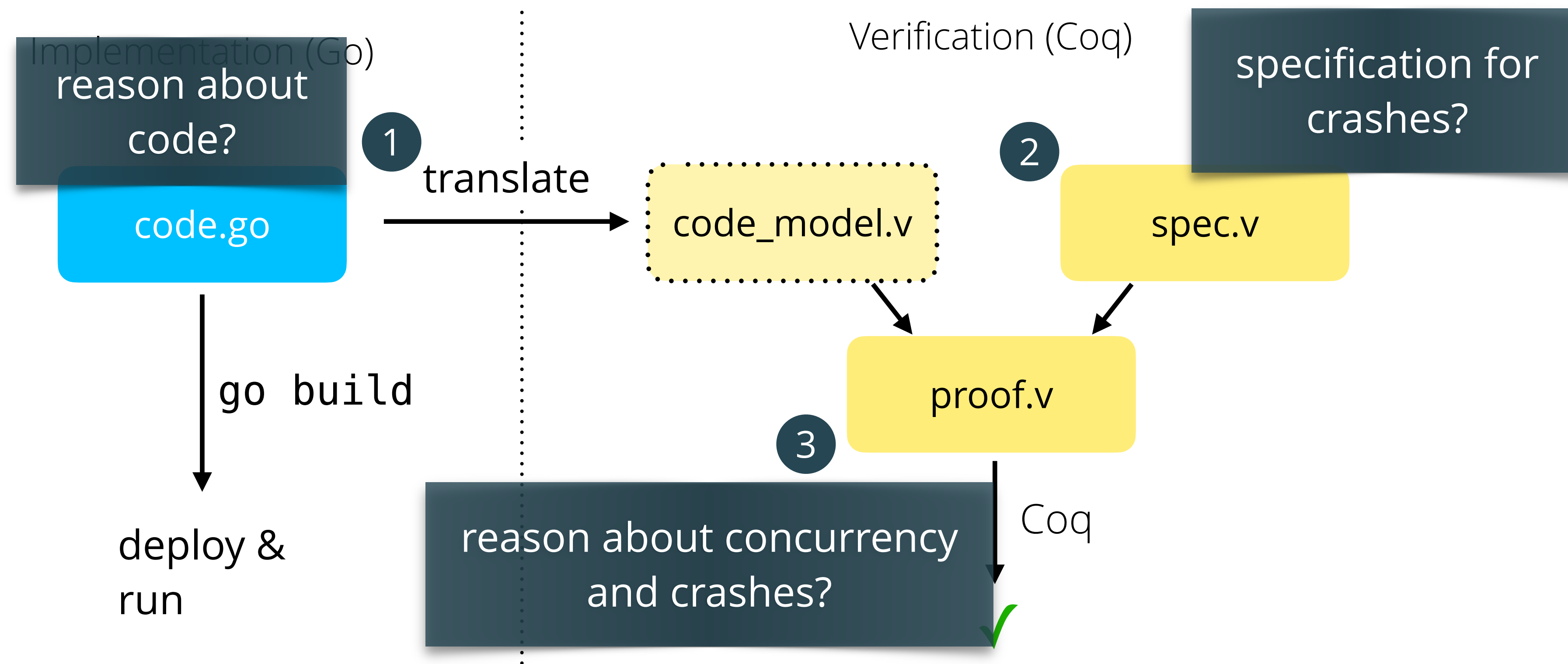
Verification of storage systems is challenging



Verification of storage systems is challenging



Verification of storage systems is challenging



Verifying a file system is a daunting task

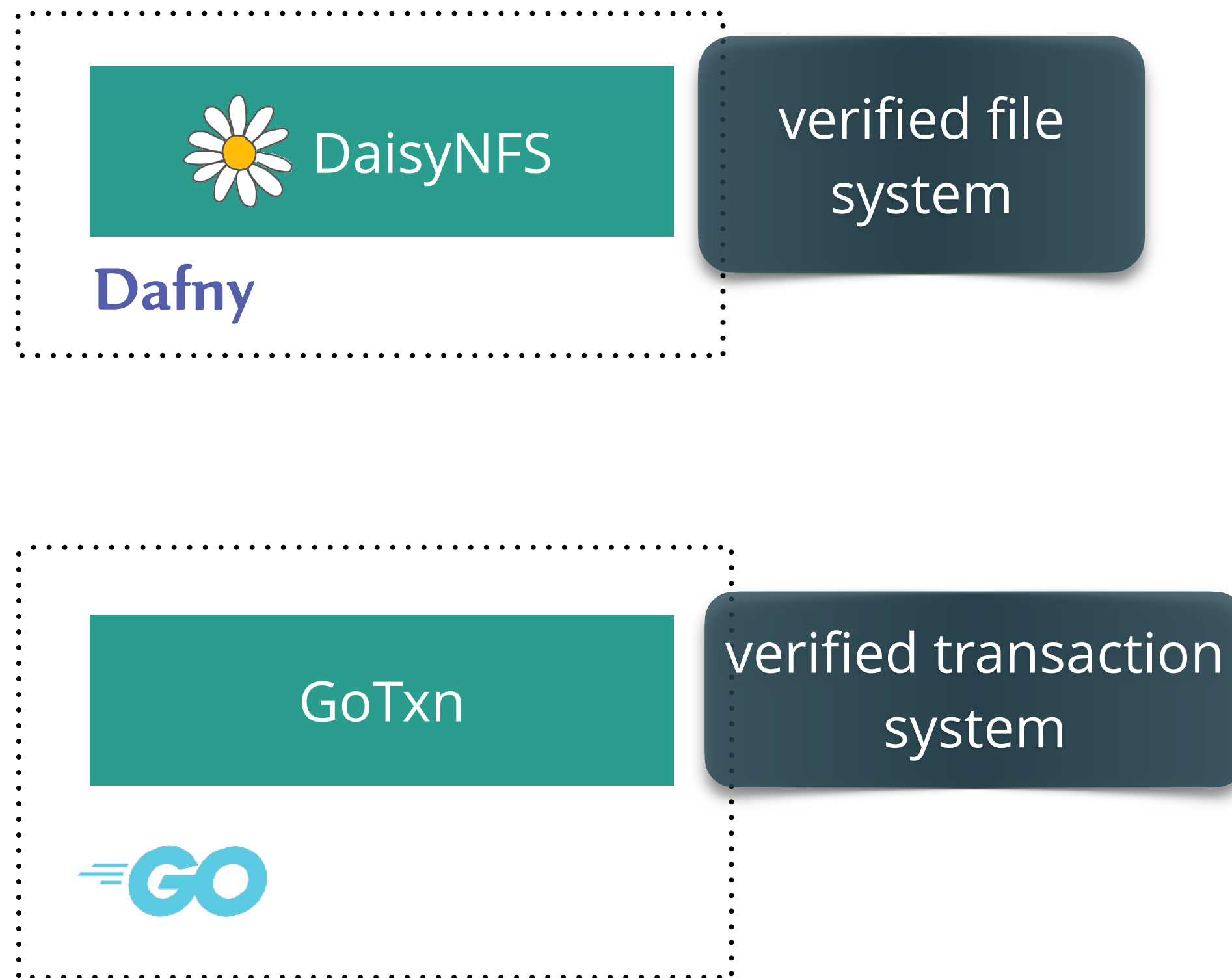
Still need to reason about **crash safety + concurrency** for a **high performance** implementation

DaisyNFS organizes the system and proof to make this manageable

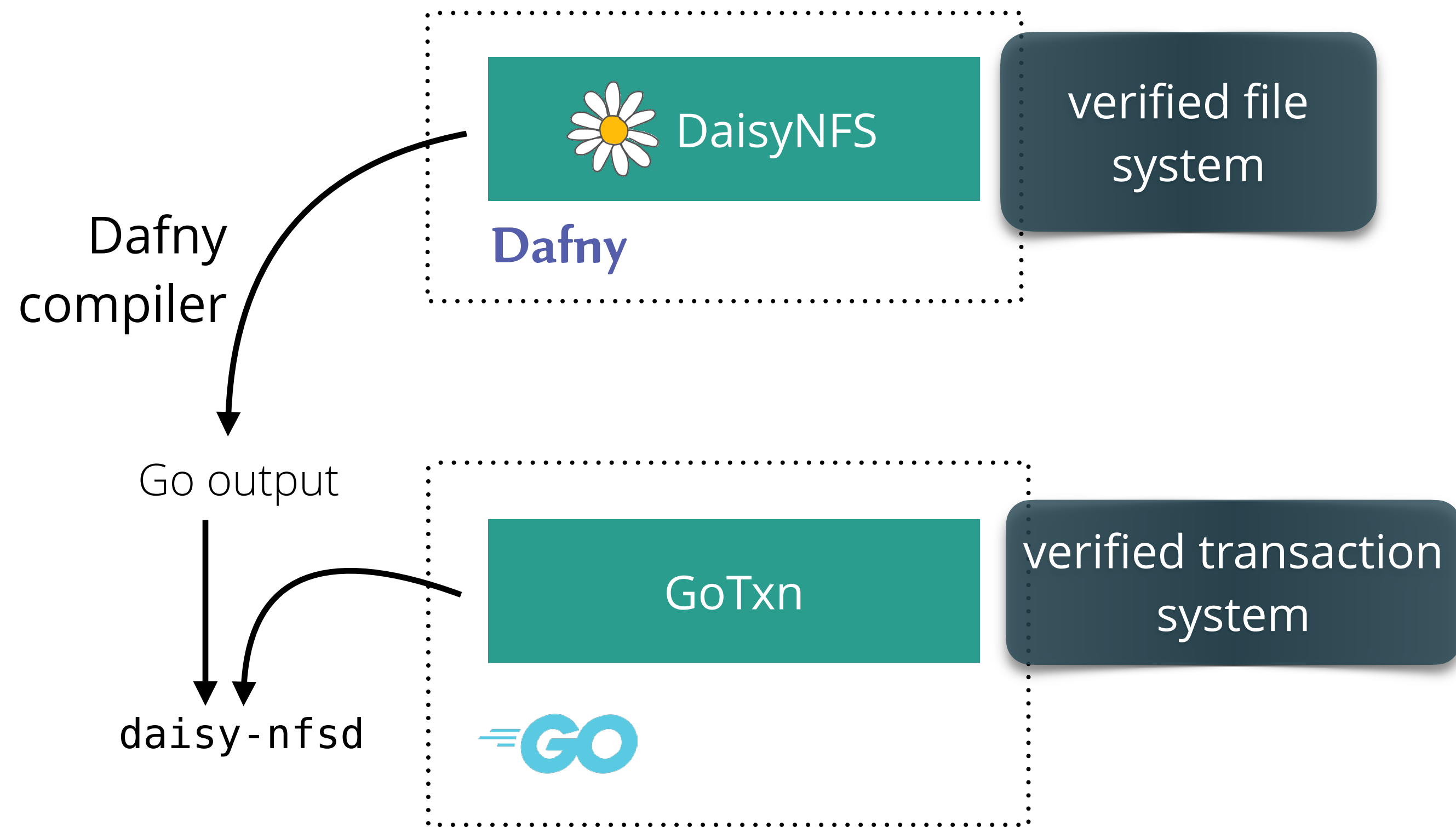
DaisyNFS architecture



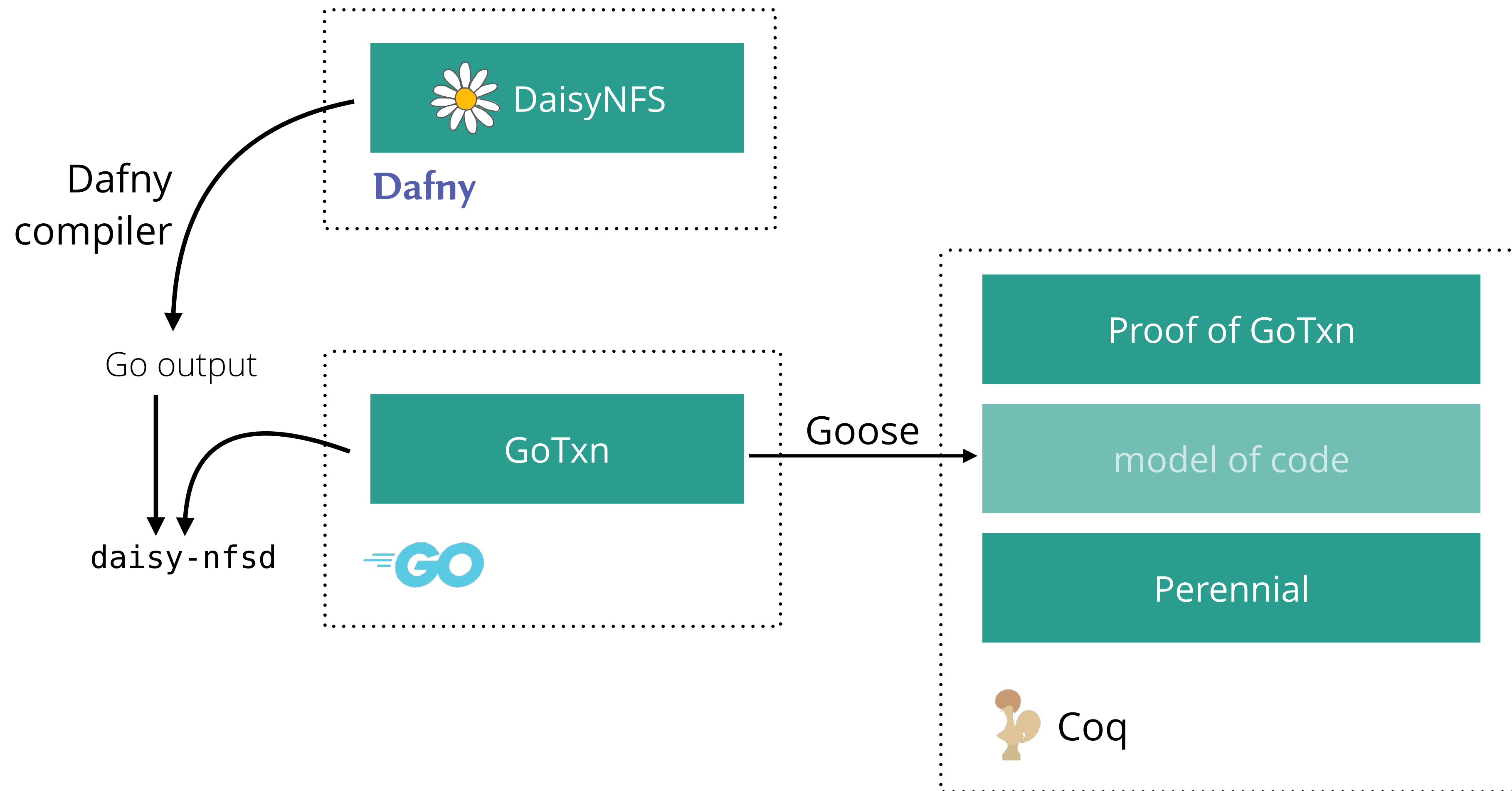
DaisyNFS architecture



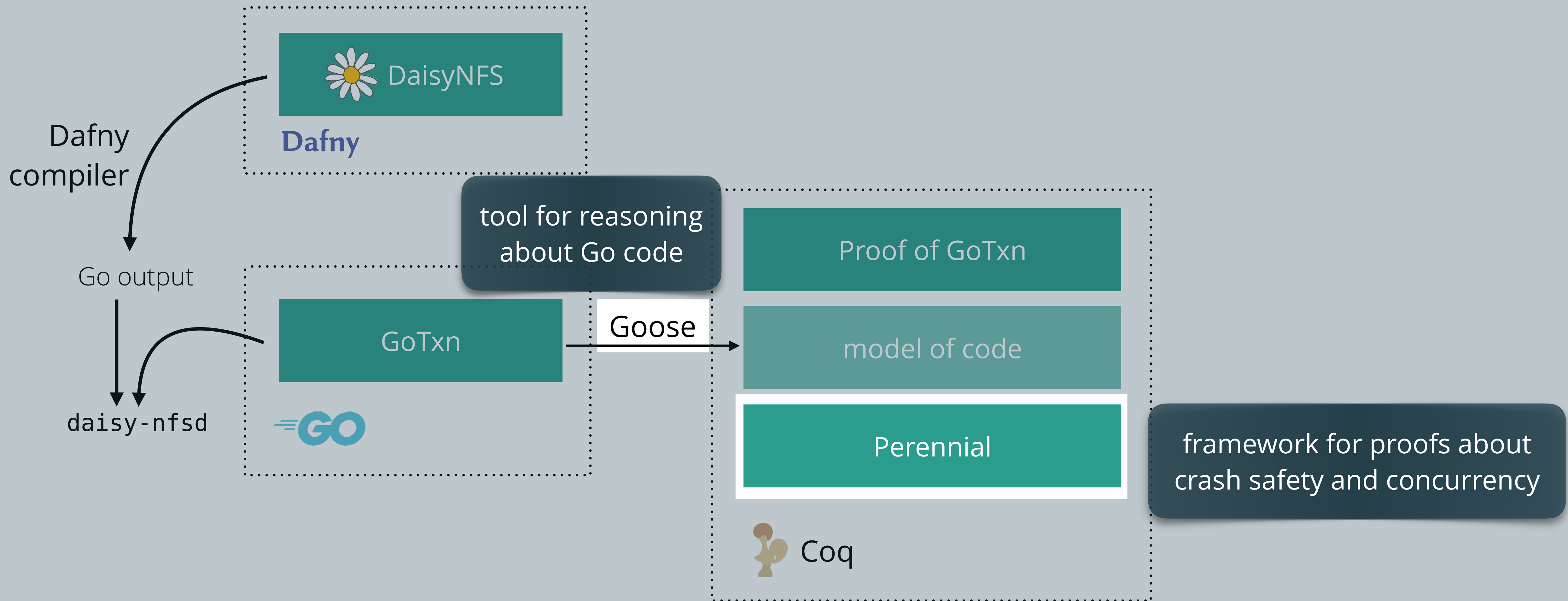
DaisyNFS architecture



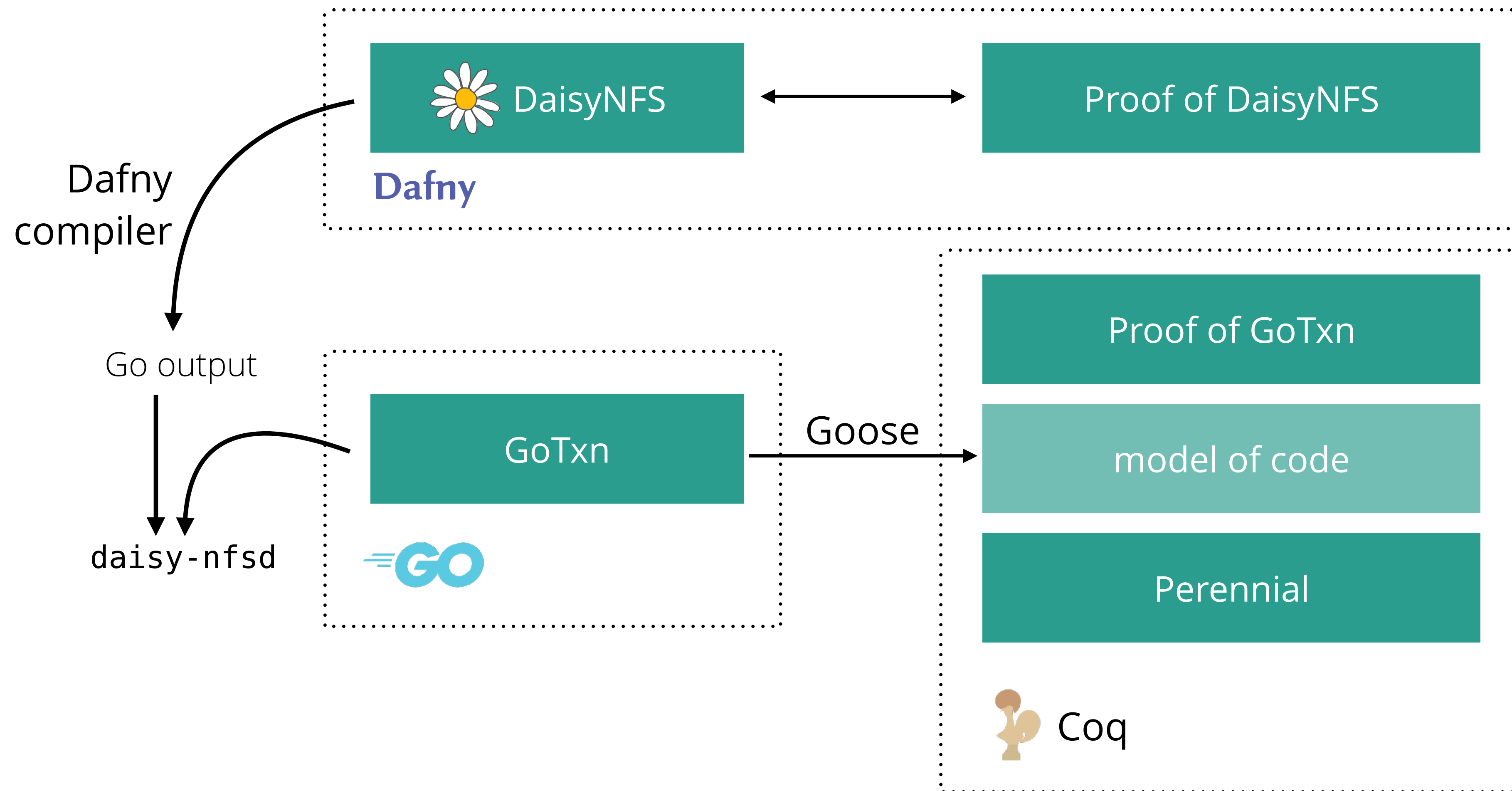
DaisyNFS architecture



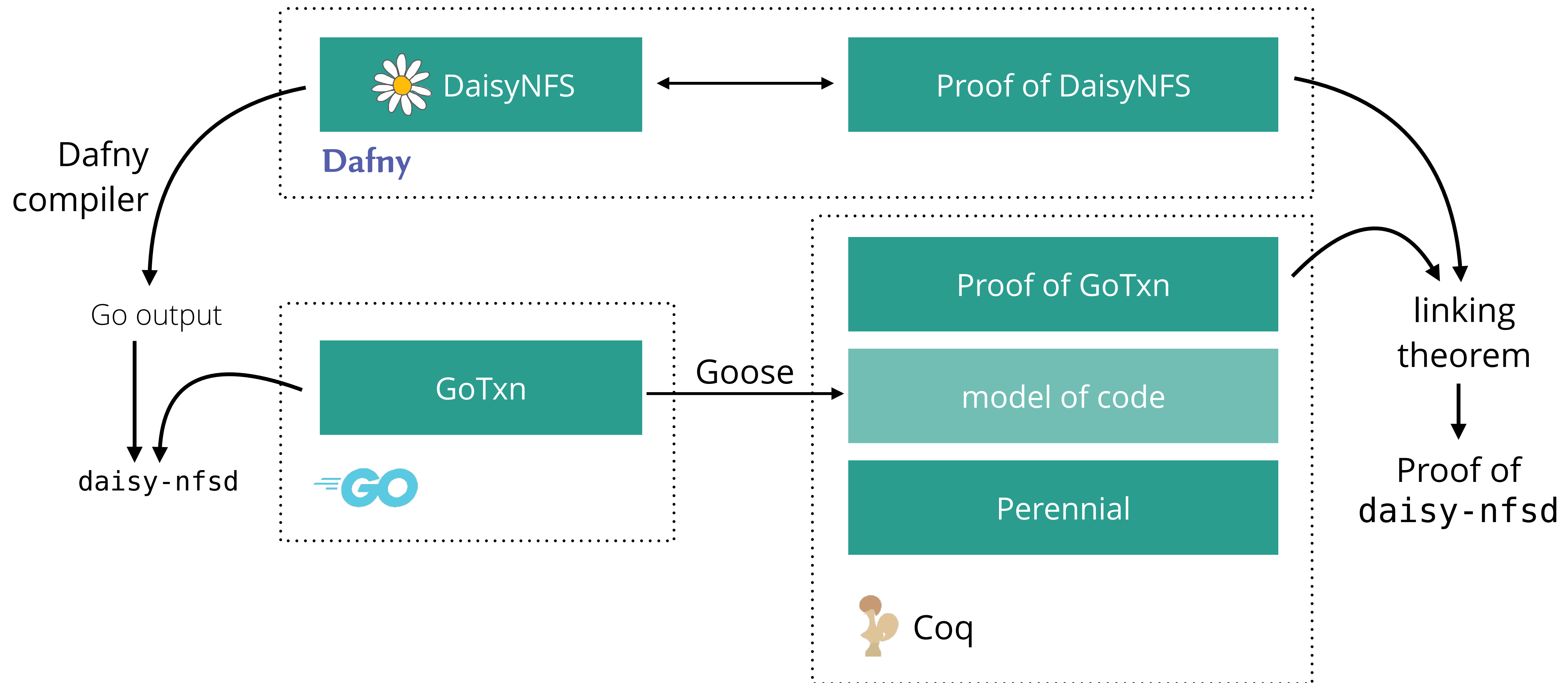
DaisyNFS architecture



DaisyNFS architecture



DaisyNFS architecture



Contributions

Perennial + Goose: foundations for verifying storage systems written in Go

GoTxn: handles crash safety and concurrency to enable sequential reasoning

DaisyNFS: a verified concurrent, crash-safe file system

What did we prove?



DaisyNFS

GoTxn



Read, Write
(of 4KB blocks)

What did we prove?

NFS GETATTR, SETATTR
CREATE, READ, WRITE, REMOVE
MKDIR, LOOKUP, REaddir, RENAME



DaisyNFS

GoTxn



Read, Write
(of 4KB blocks)

What did we prove?

NFS GETATTR, SETATTR
CREATE, READ, WRITE, REMOVE
MKDIR, LOOKUP, REaddir, RENAME



DaisyNFS

GoTxn



Read, Write
(of 4KB blocks)

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

Design and implementation of DaisyNFS

Verifying a high-performance transaction system

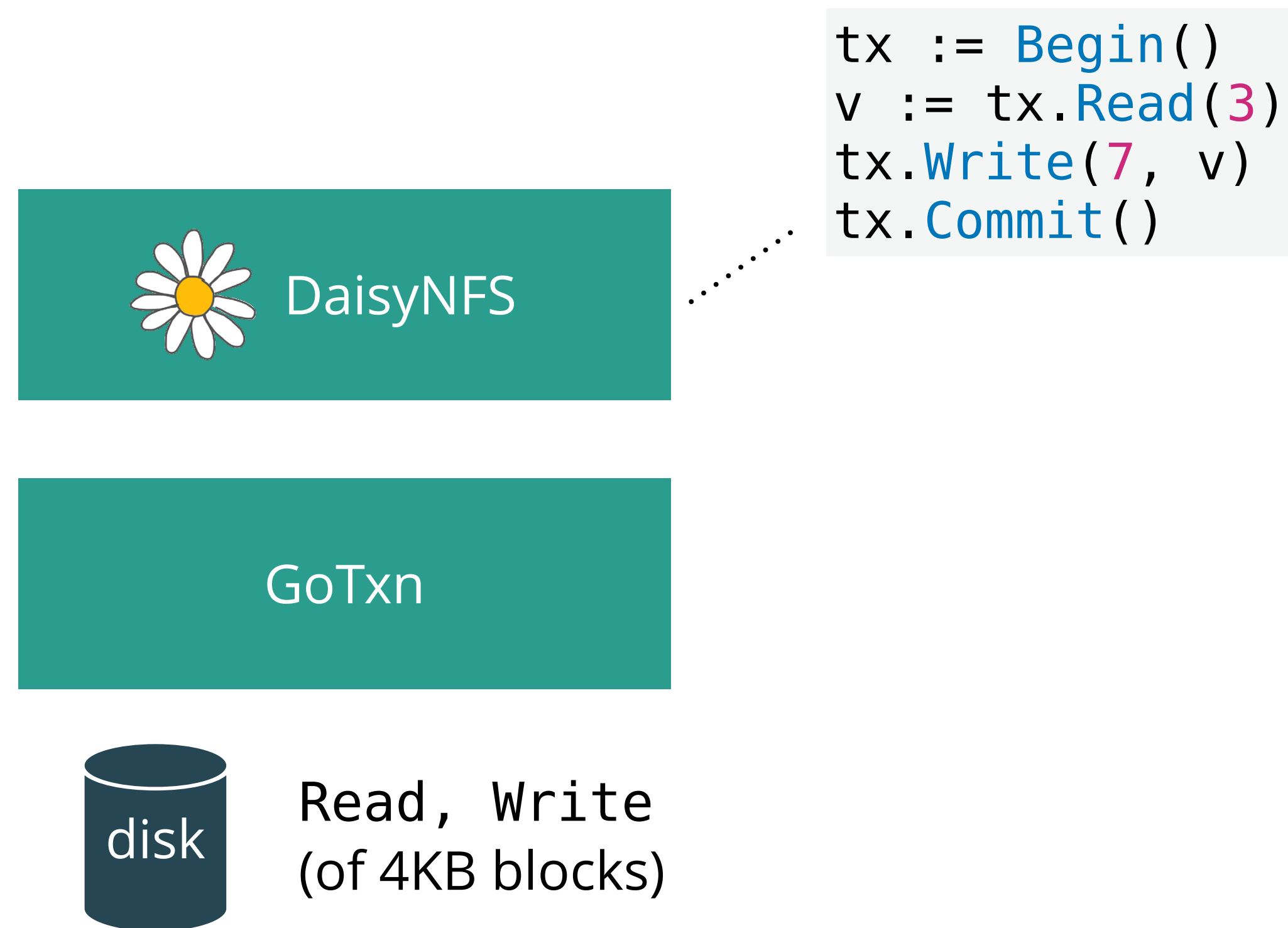
Evaluating DaisyNFS

DaisyNFS accesses the disk through a transaction system



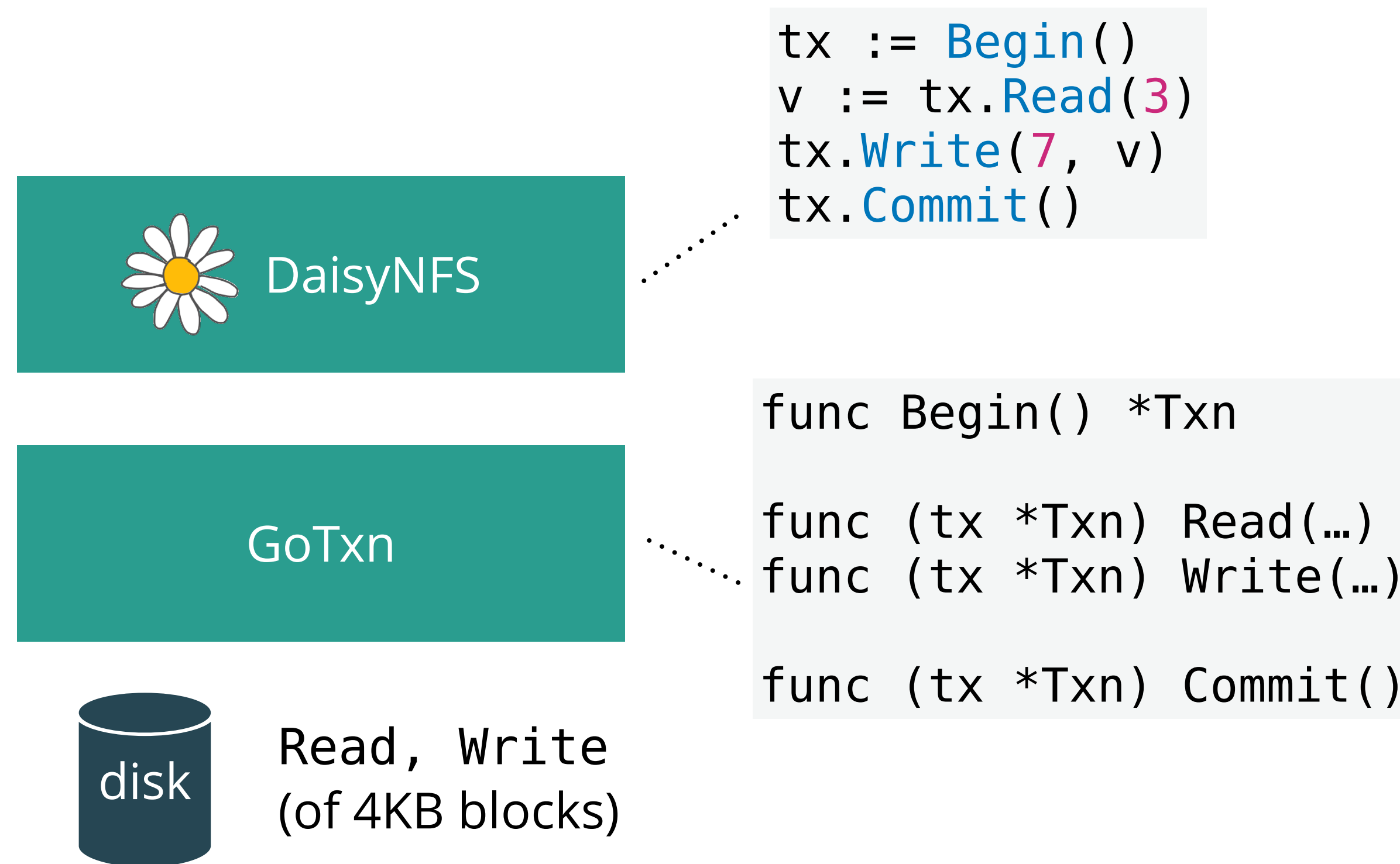
Read, Write
(of 4KB blocks)

DaisyNFS accesses the disk through a transaction system



Each operation runs within a transaction

DaisyNFS accesses the disk through a transaction system



Each operation runs within a transaction

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

Transactions isolate difficult reasoning and leave simpler sequential reasoning

```
tx := Begin()  
v := tx.Read(3)  
tx.Write(7, v)  
tx.Commit()
```

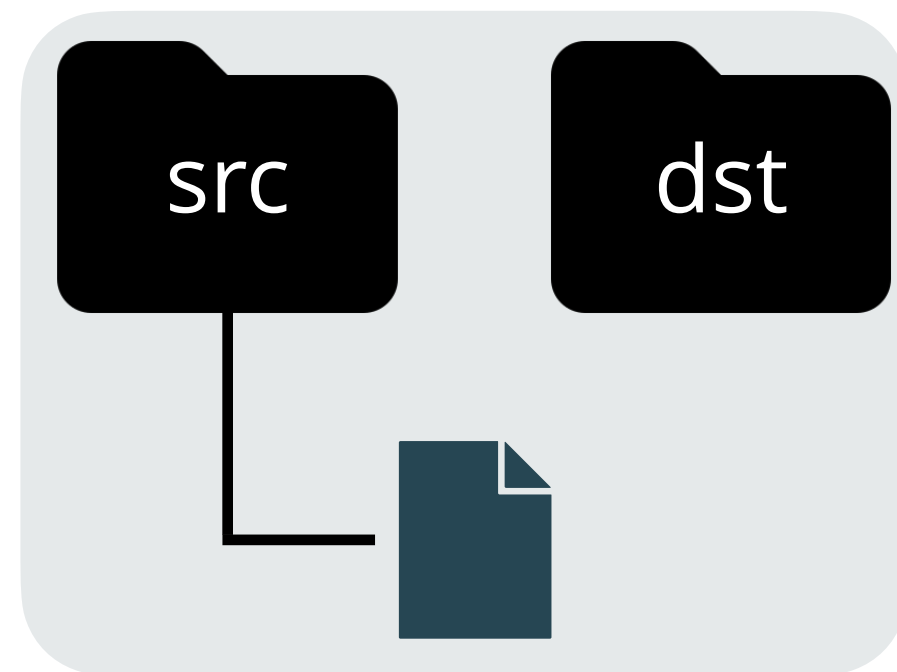
Operations are atomic — without worrying about crash safety or concurrency

GoTxn

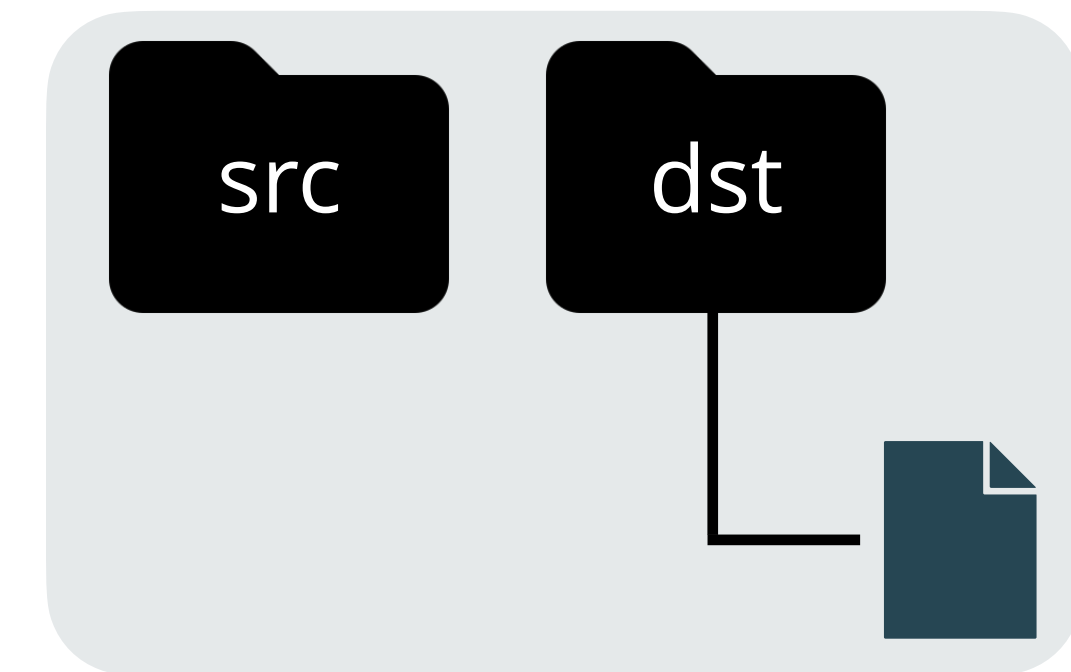
Fine-grained concurrency and crashes mean things are hard

Crash atomicity is a key correctness challenge in file systems

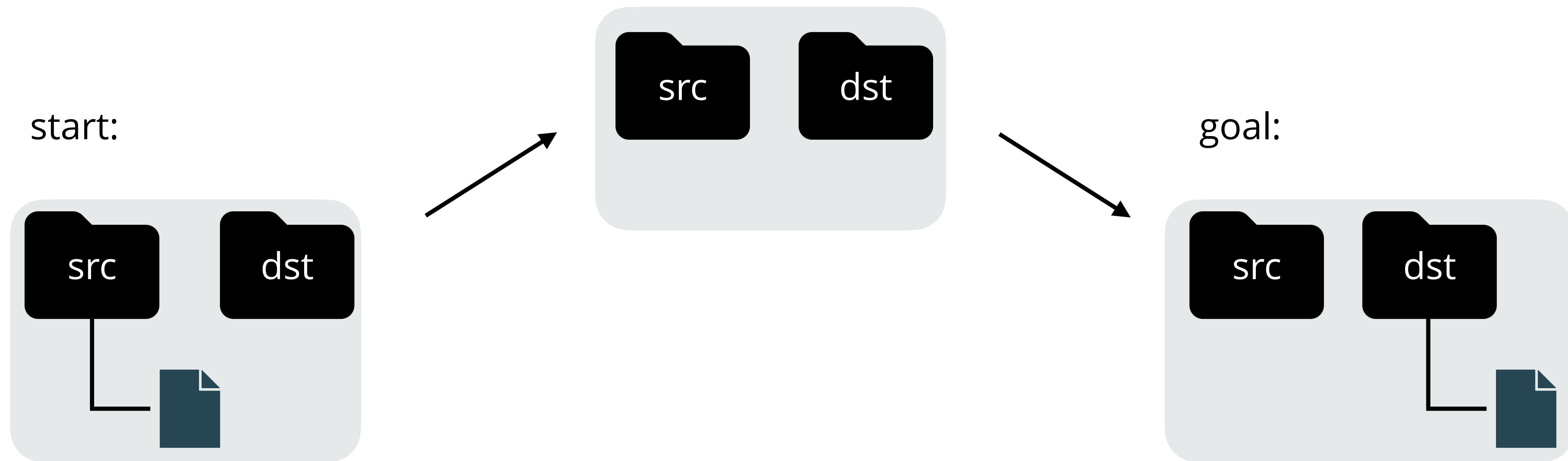
start:



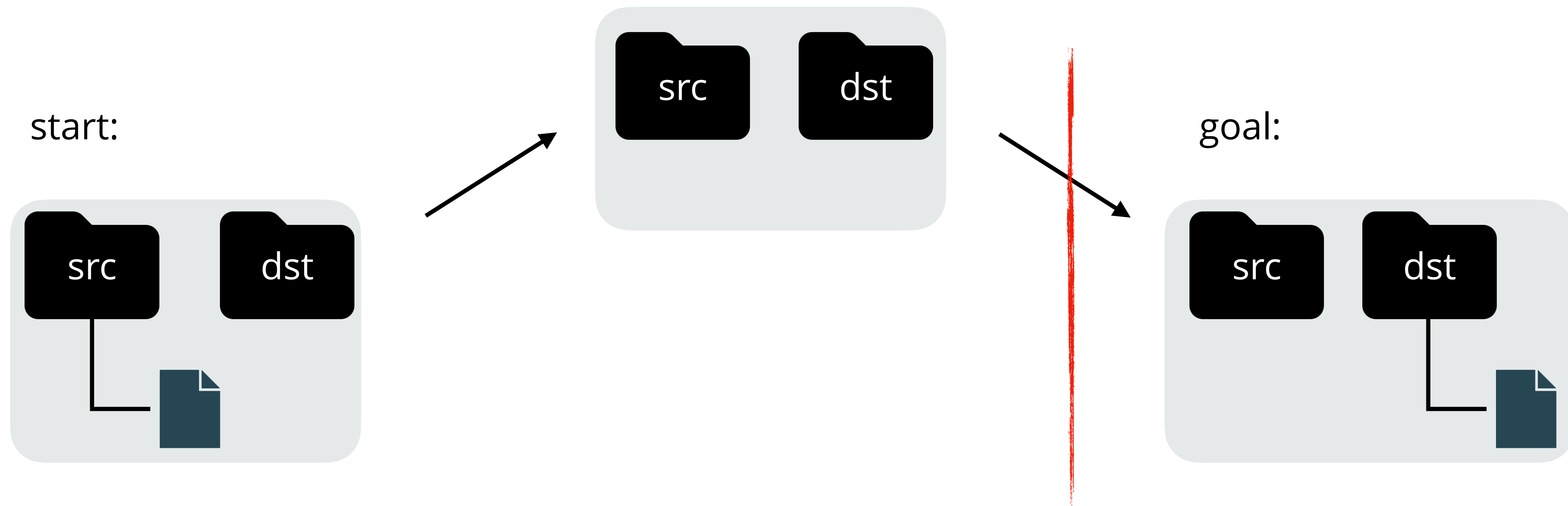
goal:



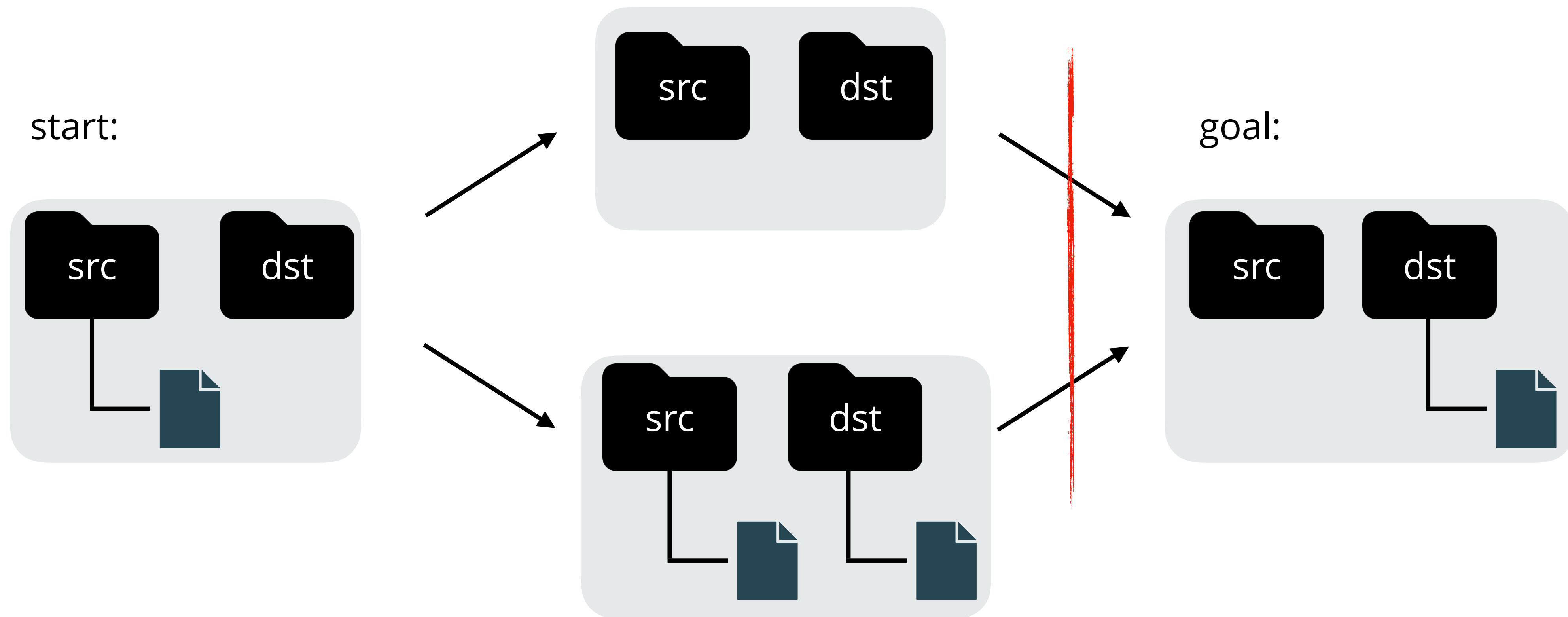
Crash atomicity is a key correctness challenge in file systems



Crash atomicity is a key correctness challenge in file systems



Crash atomicity is a key correctness challenge in file systems



Common approach is to use journaling

One solution: *journaling* is a way to write multiple values atomically

Simplifies crash atomicity but journaling is subtle to use correctly

Journal gathers up writes and issues them at once

```
op := Begin()  
v := op.Read(3)  
op.Write(7, v)  
op.Write(8, v)  
op.Commit()
```

Journal gathers up writes and issues them at once

```
op := Begin()  
v := op.Read(3)  
op.Write(7, v)  
op.Write(8, v)  
op.Commit()
```

Journal gathers up writes and issues them at once

```
op := Begin()  
v := op.Read(3)  
op.Write(7, v)  
op.Write(8, v)  
op.Commit()
```



both writes go to disk together on Commit()

Journal gathers up writes and issues them at once

```
op := Begin()  
v := op.Read(3)  
op.Write(7, v)  
op.Write(8, v)  
op.Commit()
```

both writes go to disk together on Commit()

Journal gathers up writes and issues them at once

```
op := Begin()  
v := op.Read(3)  
op.Write(7, v)  
op.Write(8, v)  
op.Commit()
```

code needs to guarantee other threads don't touch 3, 7, 8

both writes go to disk together on Commit()

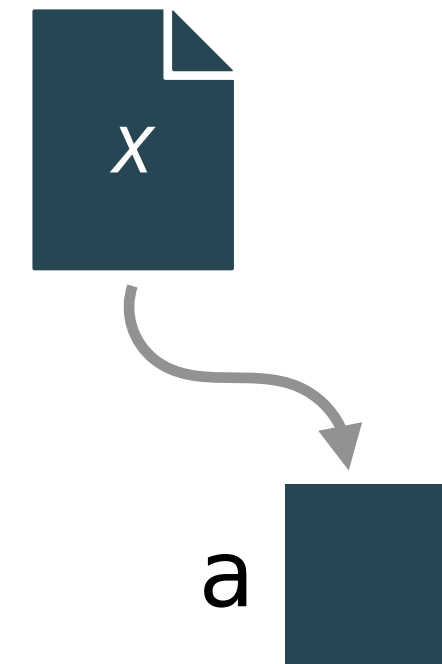
Journal gathers up writes and issues them at once

```
op := Begin()  
v := op.Read(3)  
op.Write(7, v)  
op.Write(8, v)  
op.Commit()
```

code needs to guarantee other threads don't touch 3, 7, 8

both writes go to disk together on Commit()

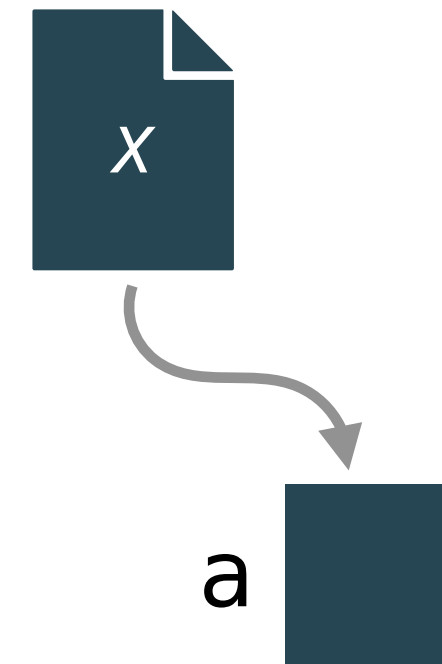
Potential bug even with journaling



Potential bug even with journaling

time ↓

```
deleting file x:  
op := Begin()  
...  
free(a)  
  
op.Commit()
```



Potential bug even with journaling

time
↓

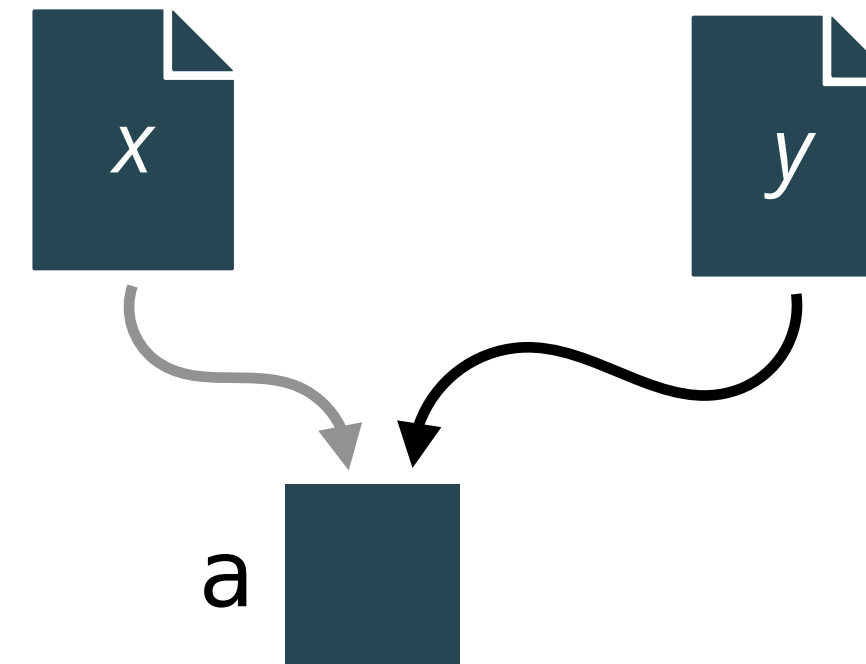
deleting file x:

```
op := Begin()  
...  
free(a)
```

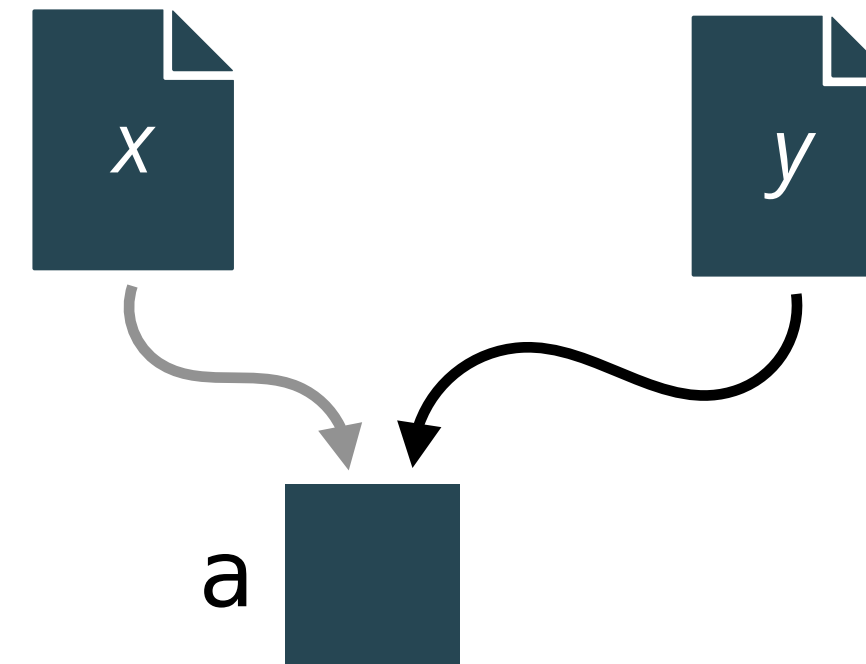
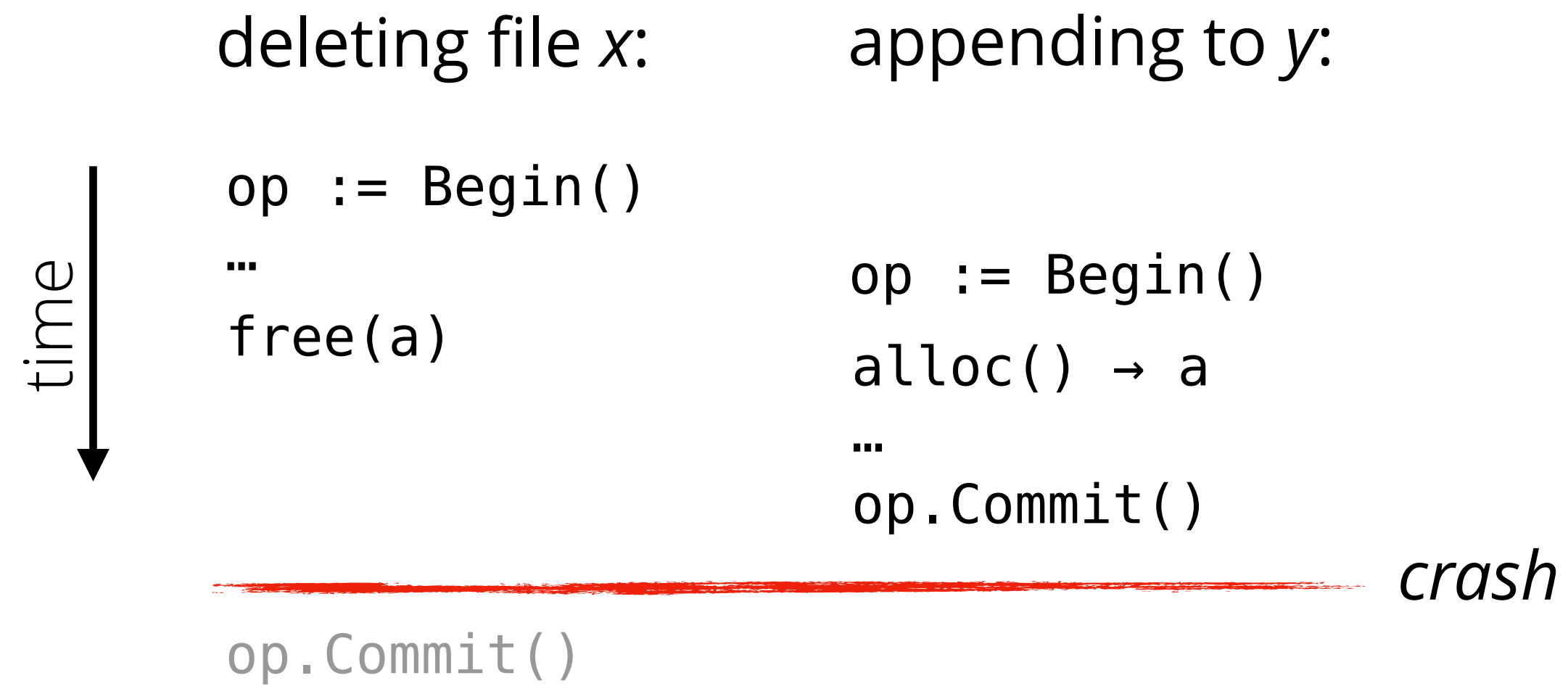
```
op.Commit()
```

appending to y:

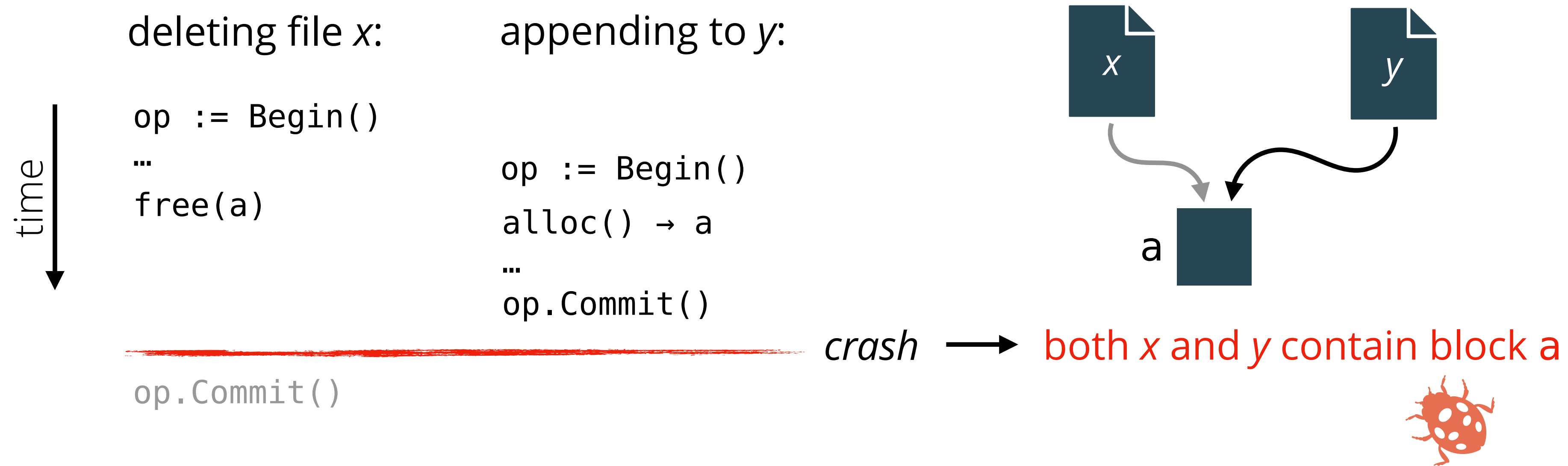
```
op := Begin()  
alloc() → a  
...  
op.Commit()
```



Potential bug even with journaling



Potential bug even with journaling



Designed a file system around transactions



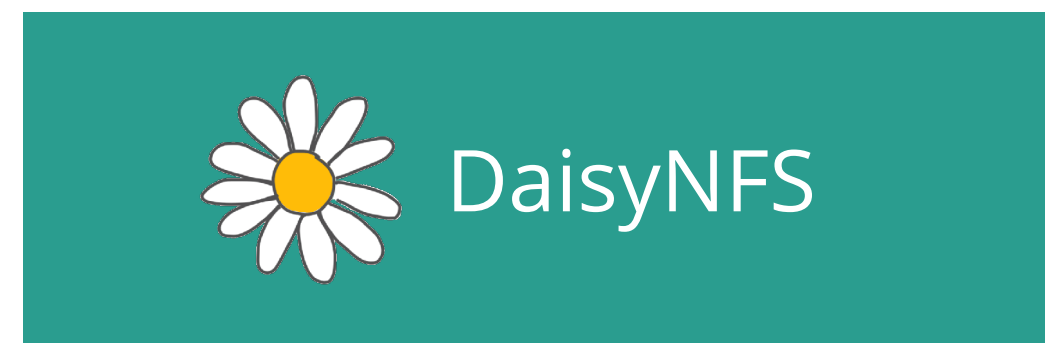
DaisyNFS

GoTxn

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

Unlike journaling, provides strong **atomicity** guarantee

Designed a file system around transactions



```
tx := Begin()  
v := tx.Read(3)  
tx.Write(7, v)  
tx.Commit()
```

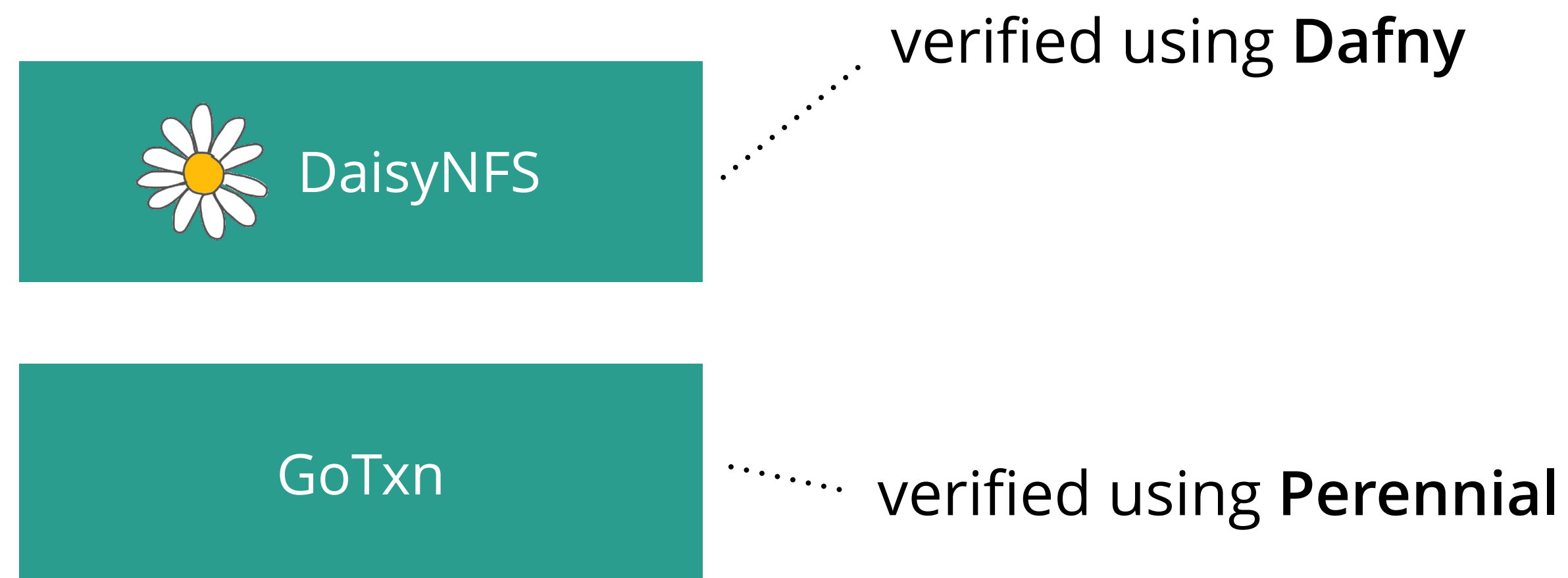
Design that fits all file-system code into transactions



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

Unlike journaling, provides strong **atomicity** guarantee

Transactions are so sequential that we verify them without a concurrency framework



Transactions are so sequential that we verify them without a concurrency framework

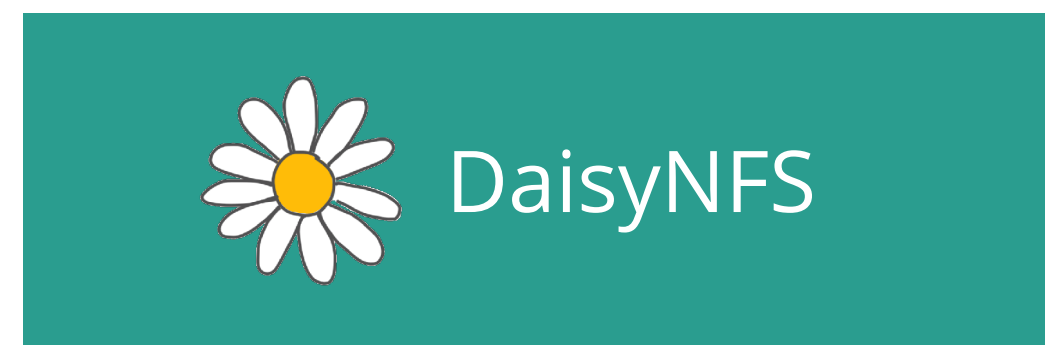


verified using **Dafny**
existing, widely-used verification system



verified using **Perennial**
our own custom infrastructure

Transactions are so sequential that we verify them without a concurrency framework

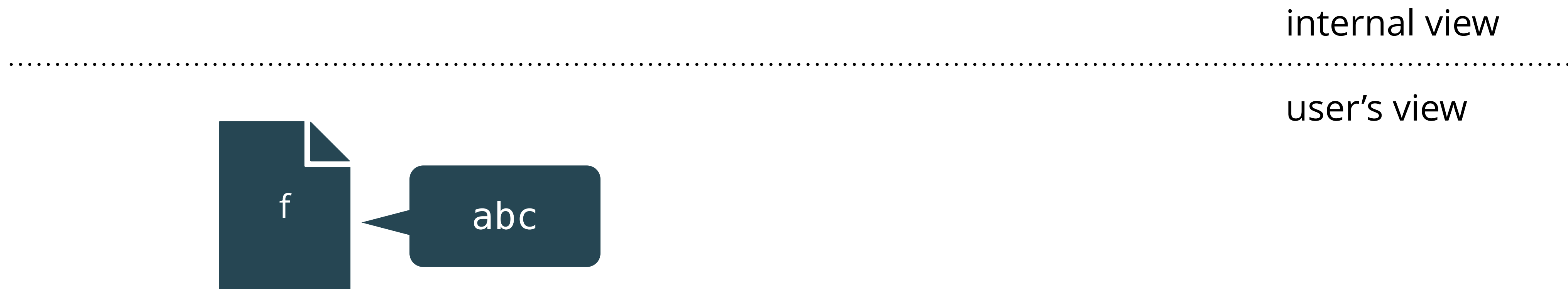


verified using **Dafny**
existing, widely-used verification system
2x as much proof as code



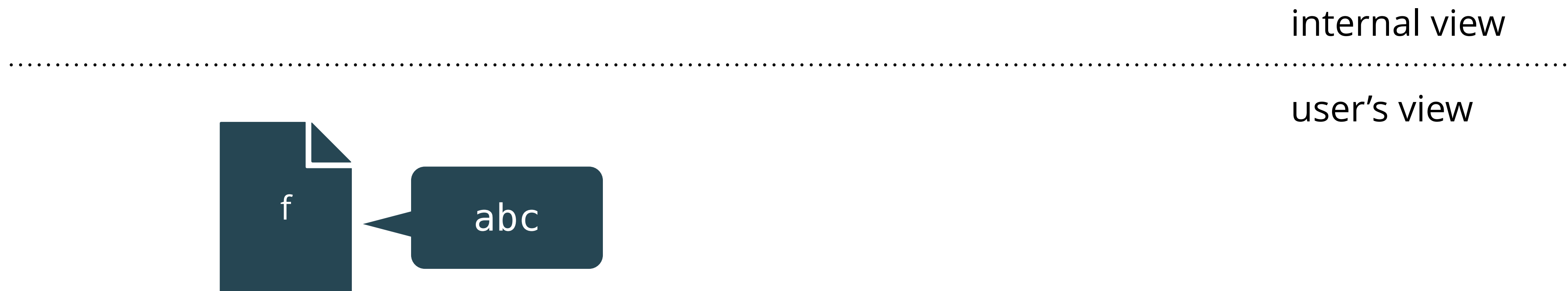
verified using **Perennial**
our own custom infrastructure
20x as much proof as code

Sequential reasoning helps because each operation needs to do a lot

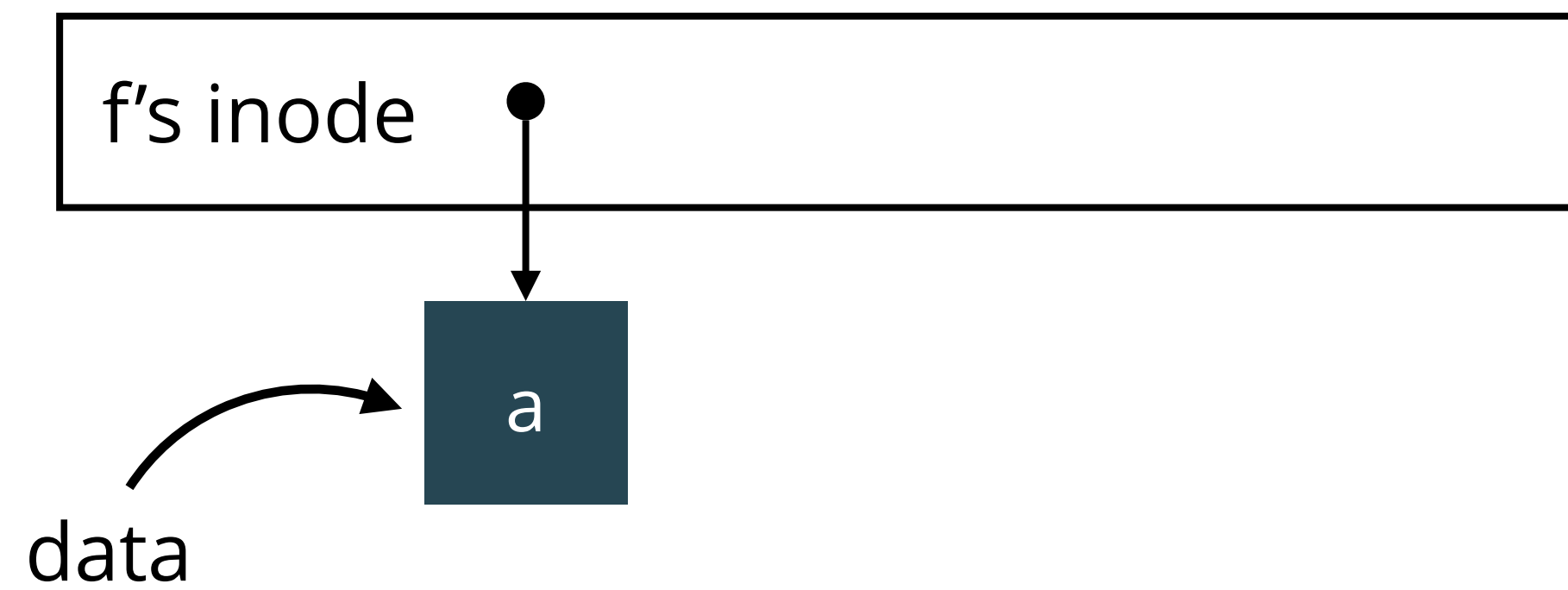


Sequential reasoning helps because each operation needs to do a lot

f's inode



Sequential reasoning helps because each operation needs to do a lot



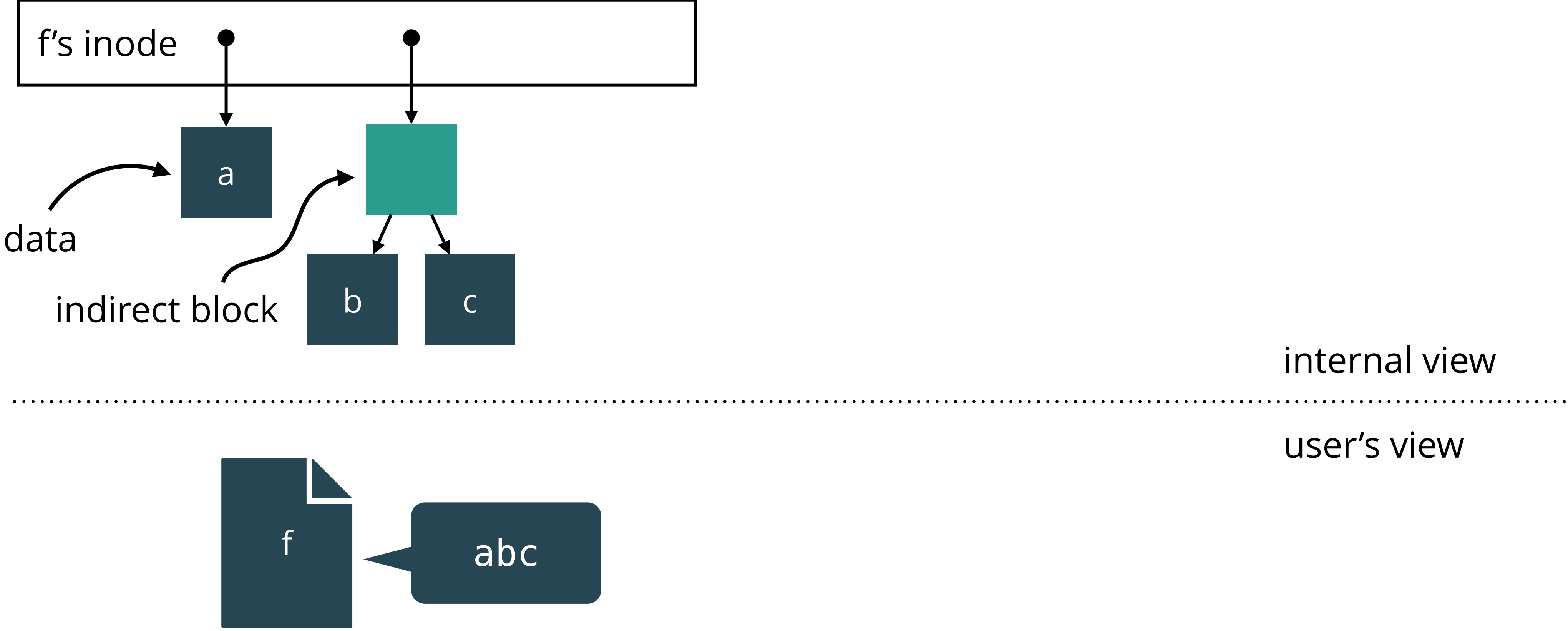
internal view



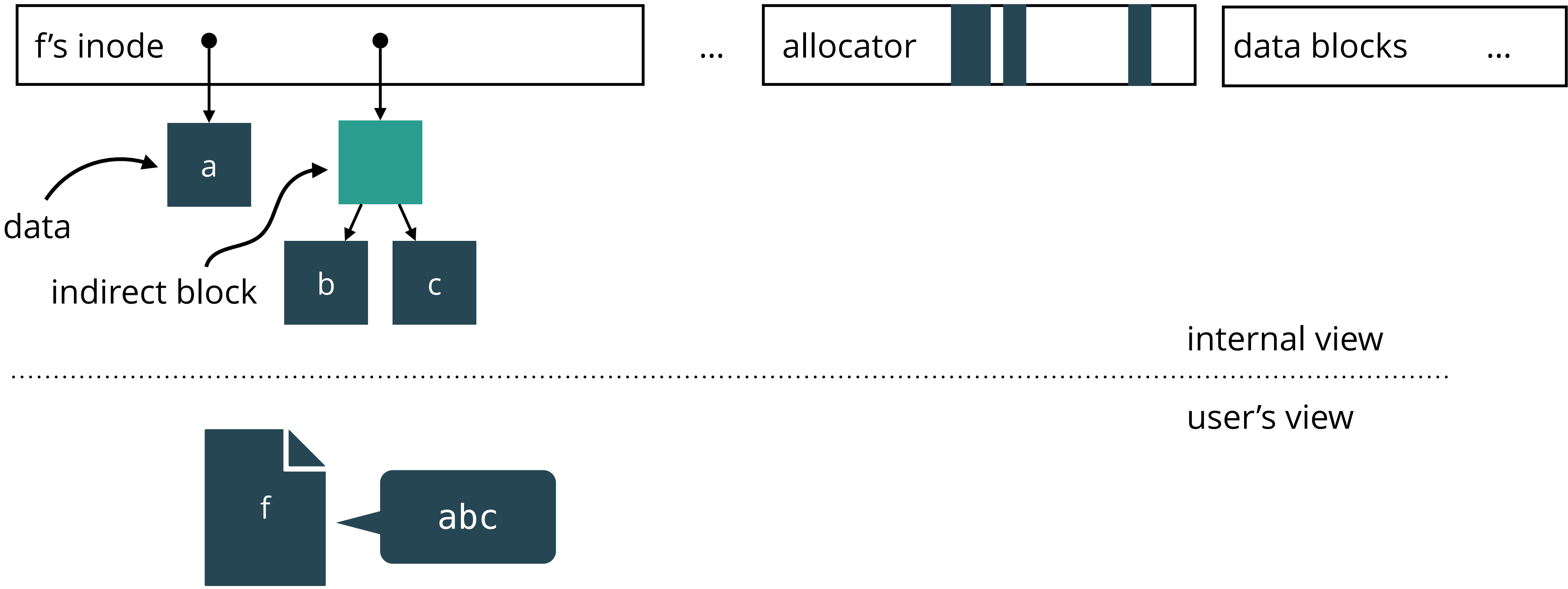
user's view



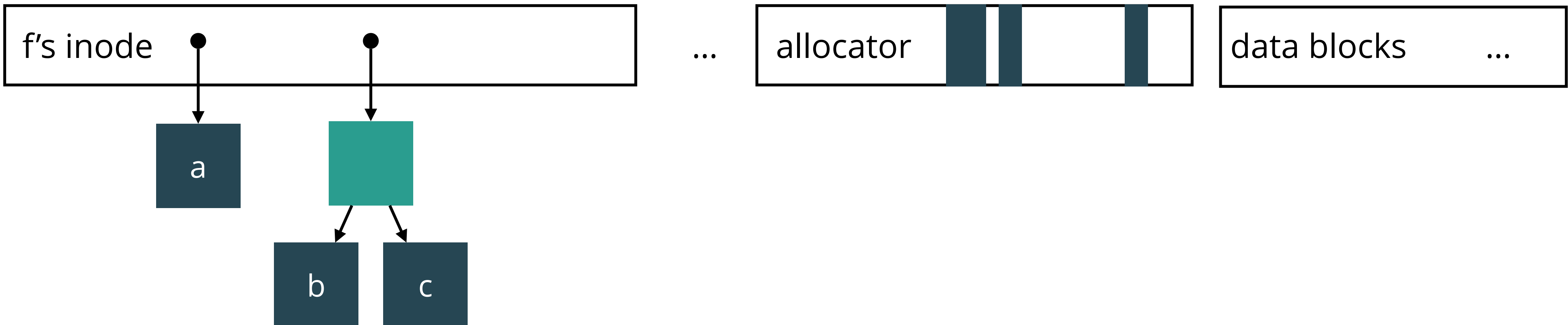
Sequential reasoning helps because each operation needs to do a lot



Sequential reasoning helps because each operation needs to do a lot



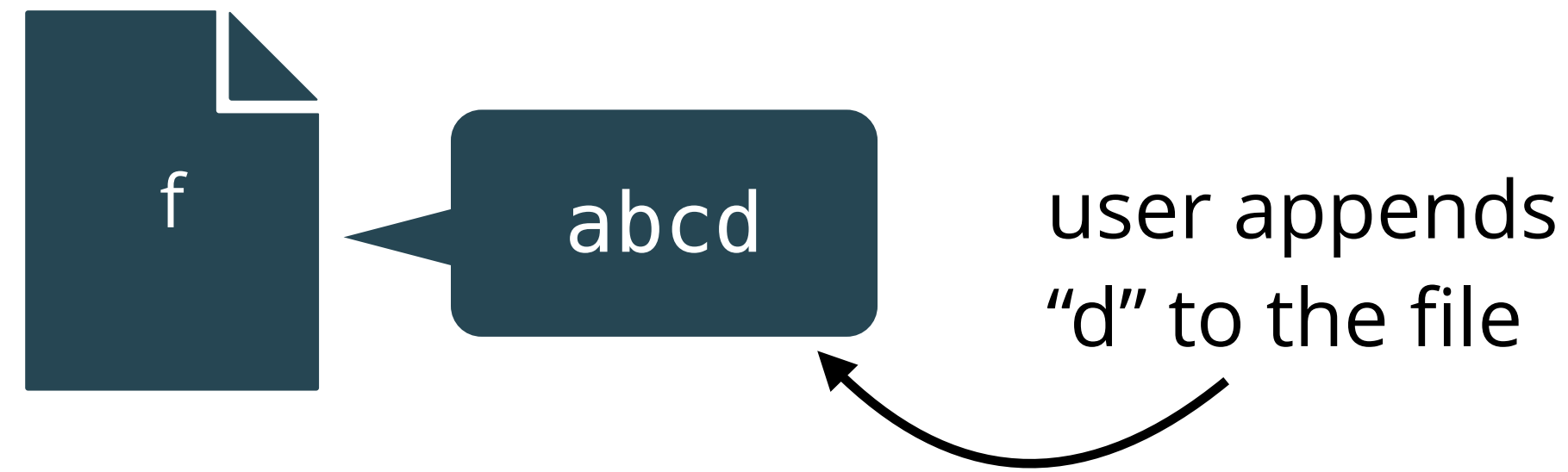
Sequential reasoning helps because each operation needs to do a lot



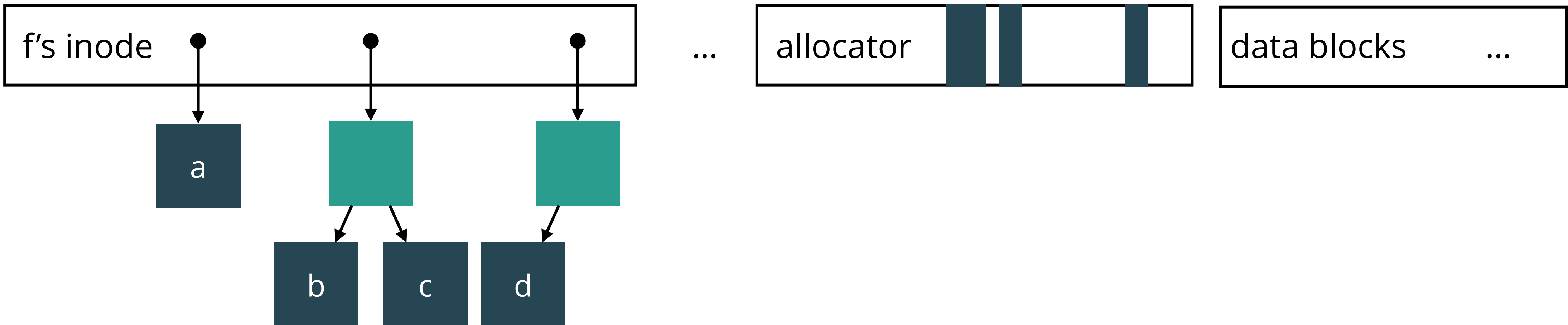
internal view



user's view



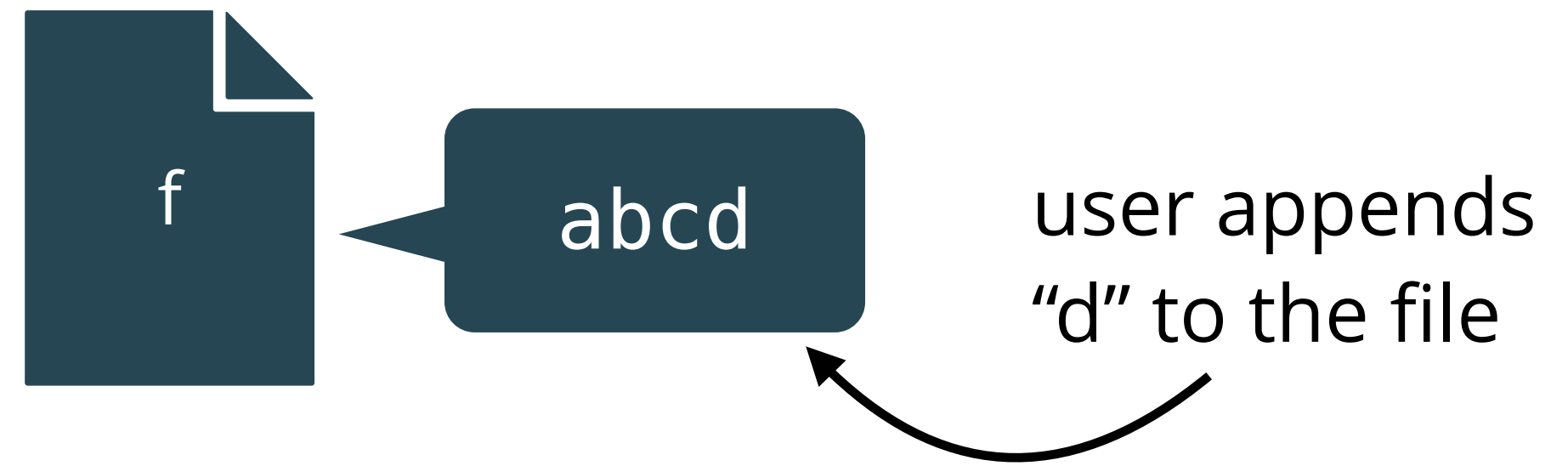
Sequential reasoning helps because each operation needs to do a lot



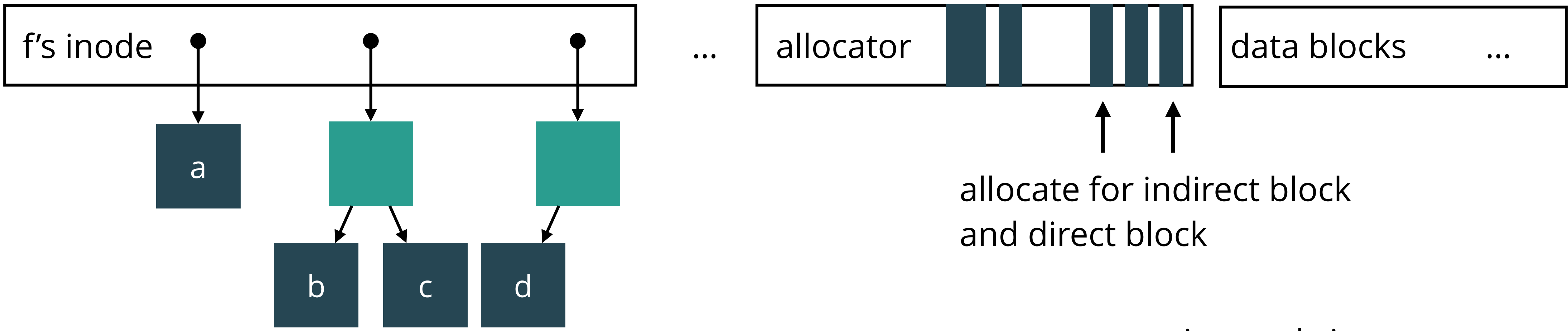
internal view



user's view

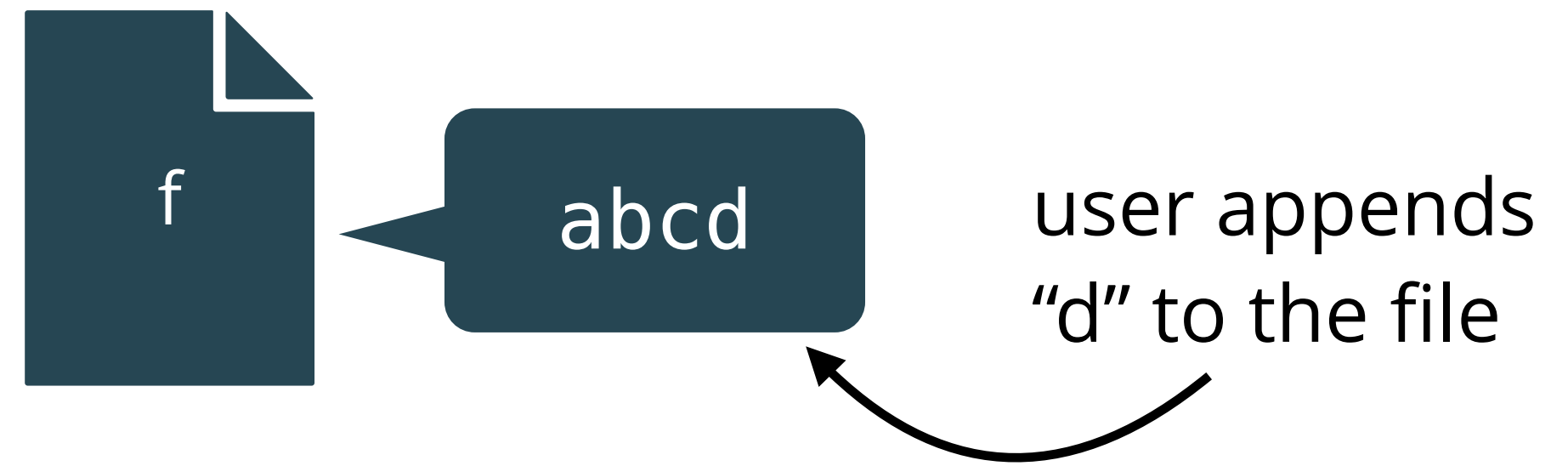


Sequential reasoning helps because each operation needs to do a lot



internal view

user's view



Careful specification of the transaction system enables this division of proof

code

```
tx := Begin()  
v := tx.Read(3)  
tx.Write(7, v)  
tx.Commit()
```

Intuitively, think of this block of code as being atomic

Careful specification of the transaction system enables this division of proof

code

```
tx := Begin()  
v := tx.Read(3)  
tx.Write(7, v)  
tx.Commit()
```

Intuitively, think of this block of code as being atomic

spec

```
atomically {  
  v ← Read(3);  
  Write(7, v);  
}
```

Specification formalizes this by relating code programs to simpler spec programs

Design and implementation of DaisyNFS

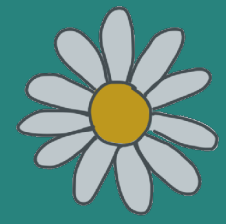
Verifying a high-performance transaction system

Evaluating DaisyNFS



DaisyNFS

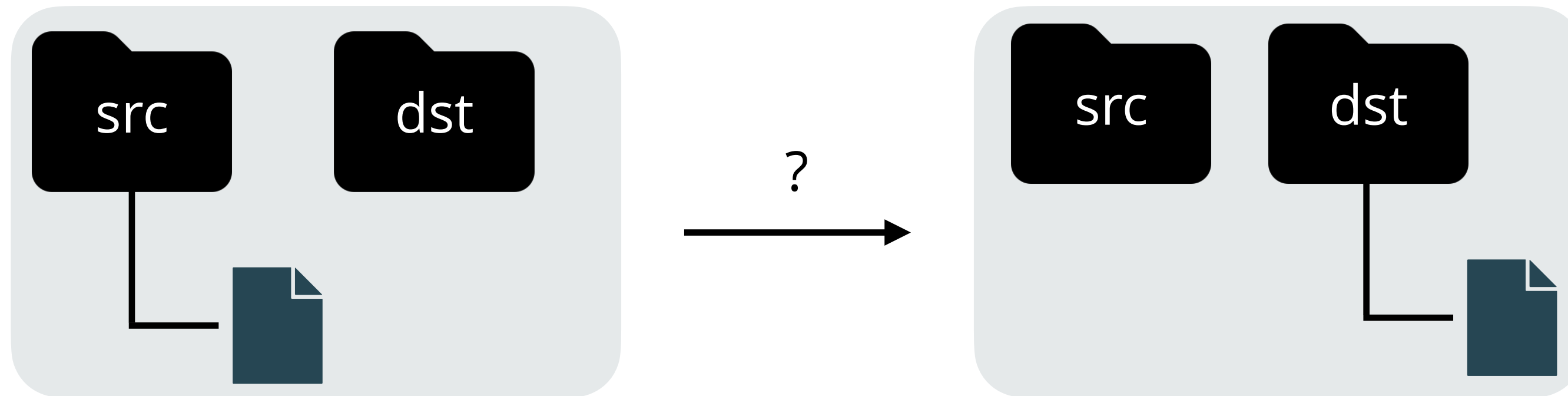
GoTxn



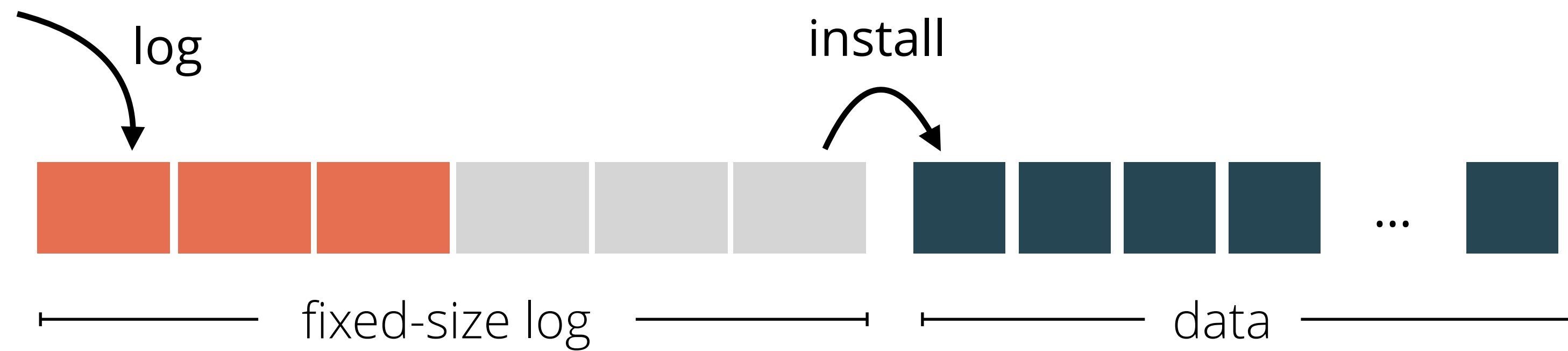
DaisyNFS

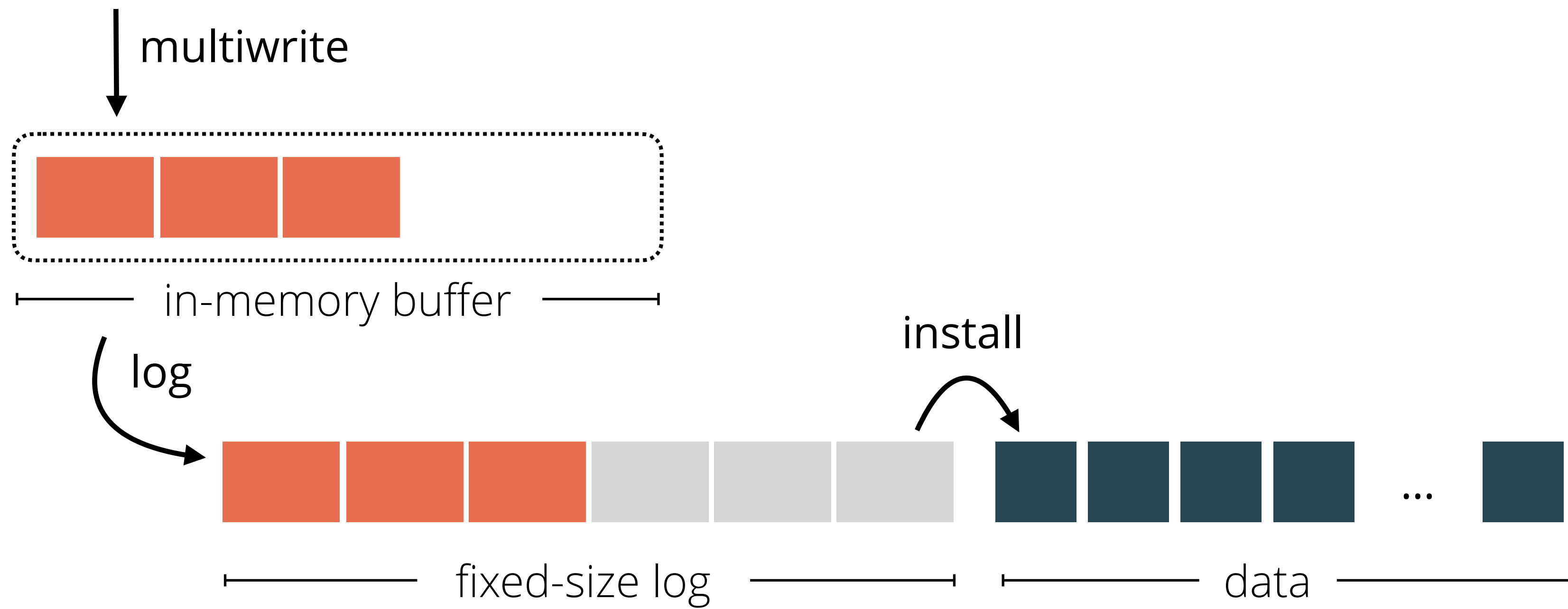
GoTxn

Recall: RENAME needs to update two things atomically

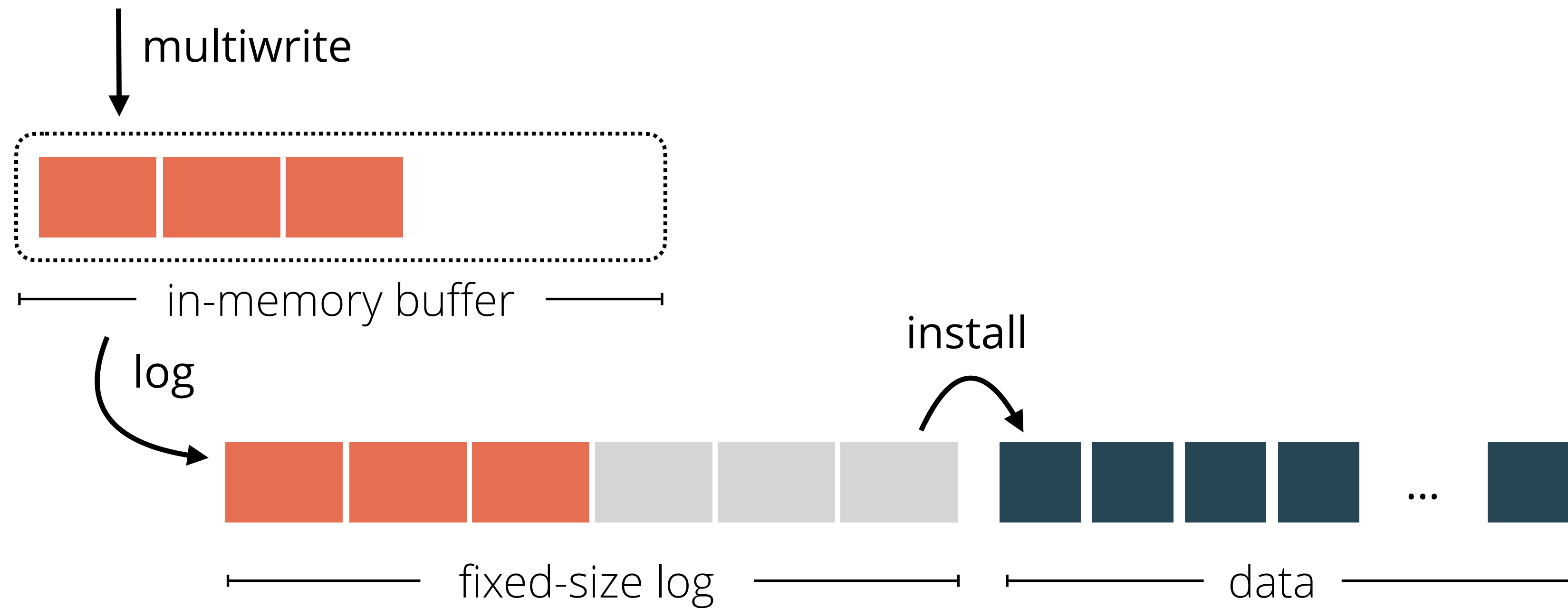


Write-ahead logging is the core atomicity primitive



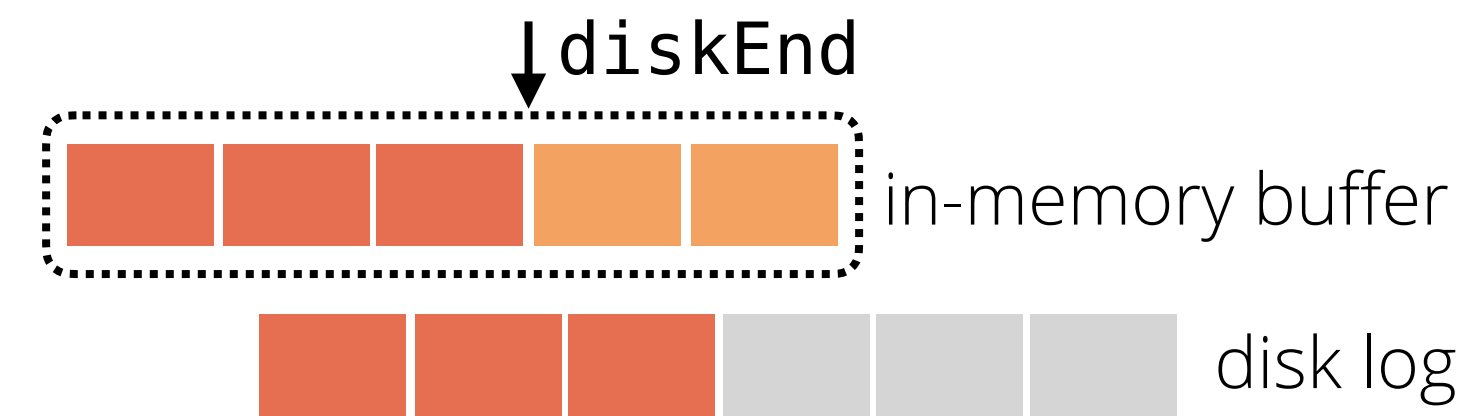


Writes, logging, and installation are all concurrent



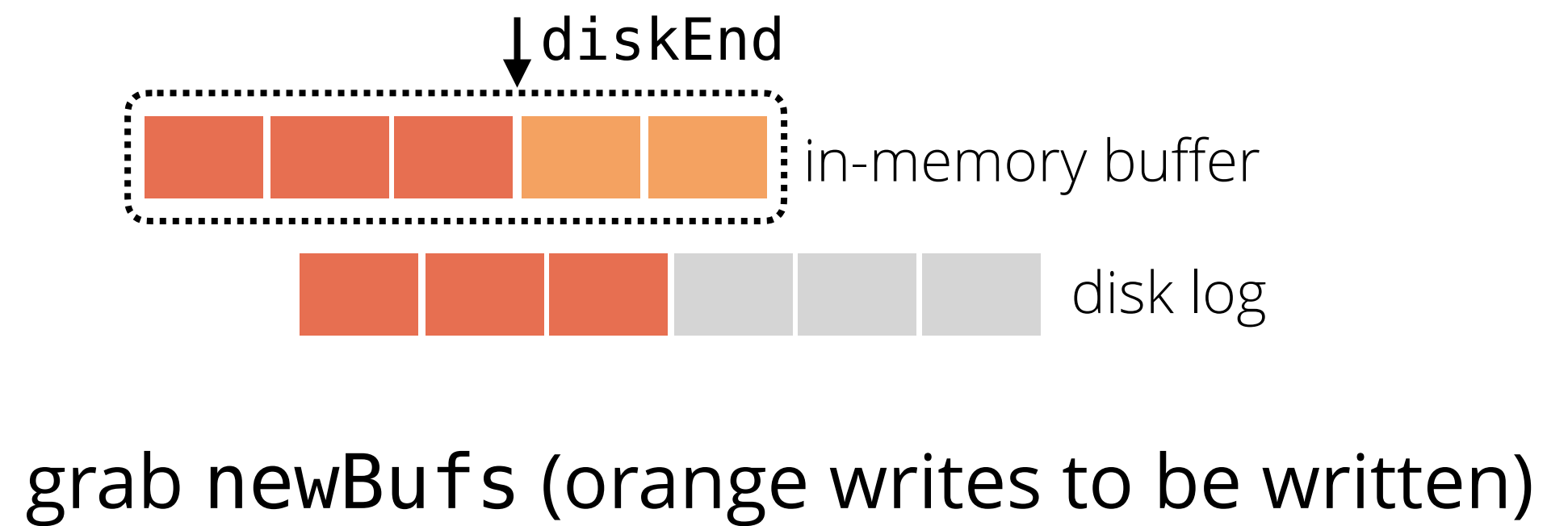
Logging happens lock-free in the background

```
var diskEnd uint64
for {
}
}
```



Logging happens lock-free in the background

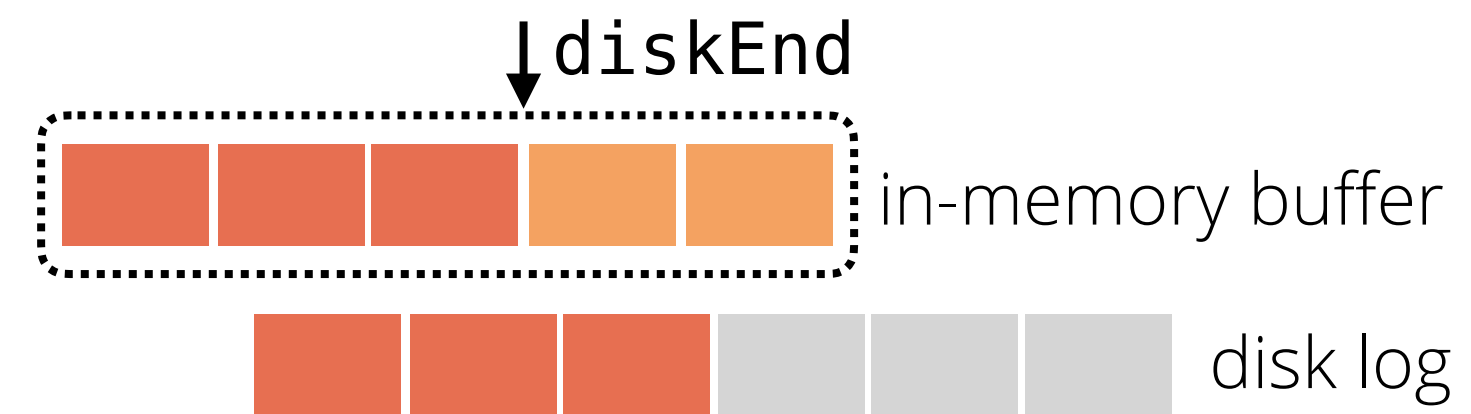
```
var diskEnd uint64
for {
    l.memLock.Lock()
    newBufs := l.memLog.takeFrom(diskEnd)
    l.memLock.Unlock()
}
```



Logging happens lock-free in the background

```
var diskEnd uint64
for {
    l.memLock.Lock()
    newBufs := l.memLog.takeFrom(diskEnd)
    l.memLock.Unlock()

    circ.Append(diskEnd, newBufs)
}
```



grab newBufs (orange writes to be written)

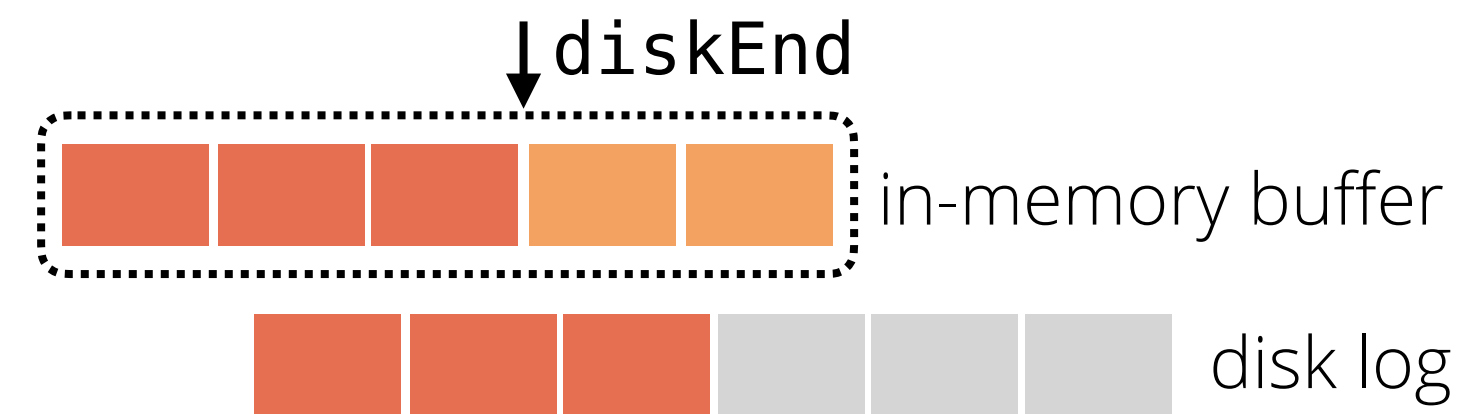
append newBufs to log

Logging happens lock-free in the background

```
var diskEnd uint64
for {
    l.memLock.Lock()
    newBufs := l.memLog.takeFrom(diskEnd)
    l.memLock.Unlock()

    circ.Append(diskEnd, newBufs)

    l.memLock.Lock()
    diskEnd += len(newBufs)
    l.memLock.Unlock()
}
```



grab newBufs (orange writes to be written)

append newBufs to log

record that this batch is durable

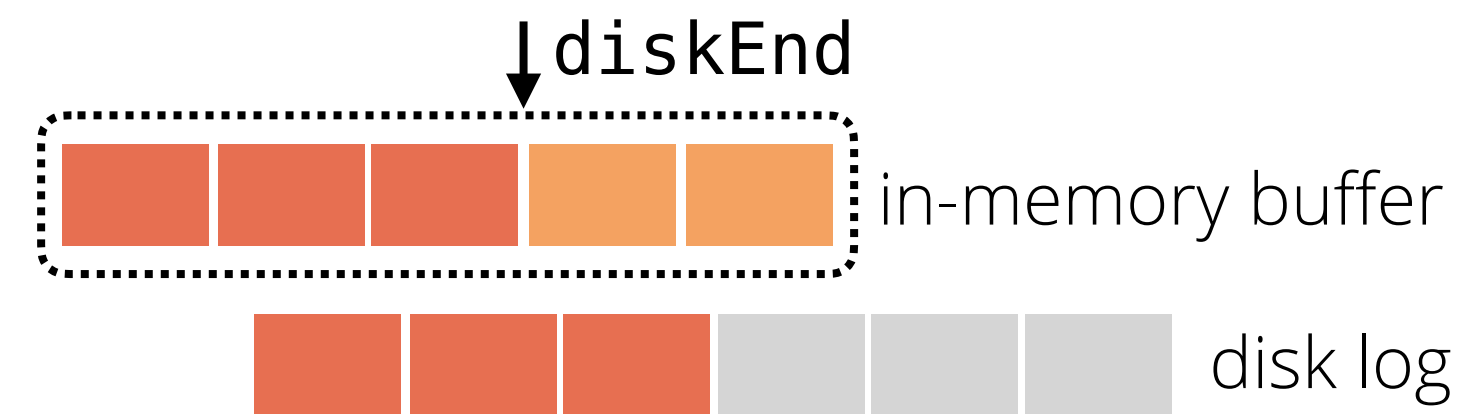
Logging happens lock-free in the background

```
var diskEnd uint64
for {
    l.memLock.Lock()
    newBufs := l.memLog.takeFrom(diskEnd)
    l.memLock.Unlock()

    circ.Append(diskEnd, newBufs)

    l.memLock.Lock()
    diskEnd += len(newBufs)
    l.memLock.Unlock()

    // wait for a bit
}
```



grab newBufs (orange writes to be written)

append newBufs to log

record that this batch is durable

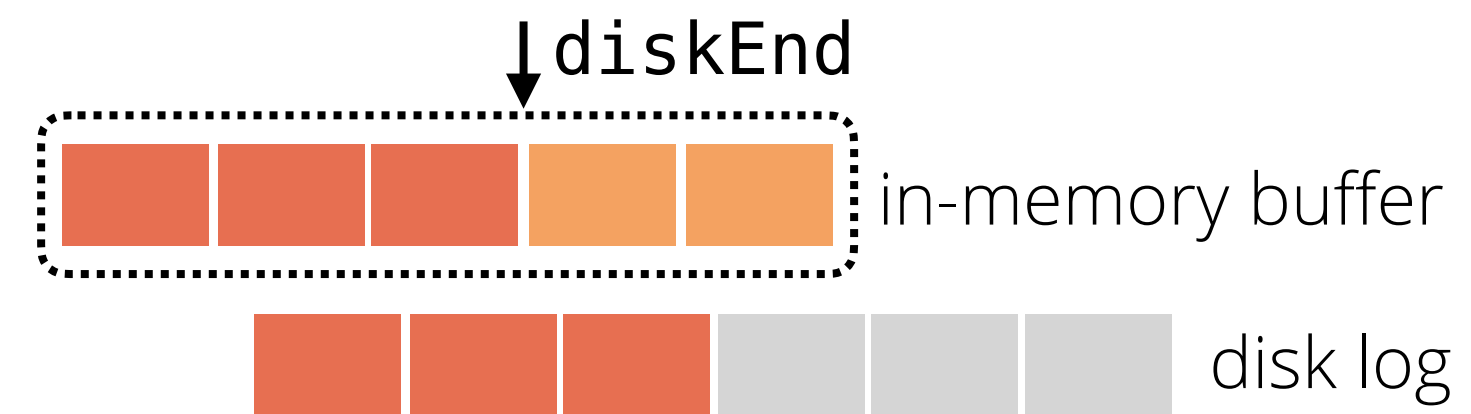
Logging happens lock-free in the background

```
var diskEnd uint64
for {
    l.memLock.Lock()
    newBufs := l.memLog.takeFrom(diskEnd)
    l.memLock.Unlock()

    circ.Append(diskEnd, newBufs)

    l.memLock.Lock()
    diskEnd += len(newBufs)
    l.memLock.Unlock()

    // wait for a bit
}
```

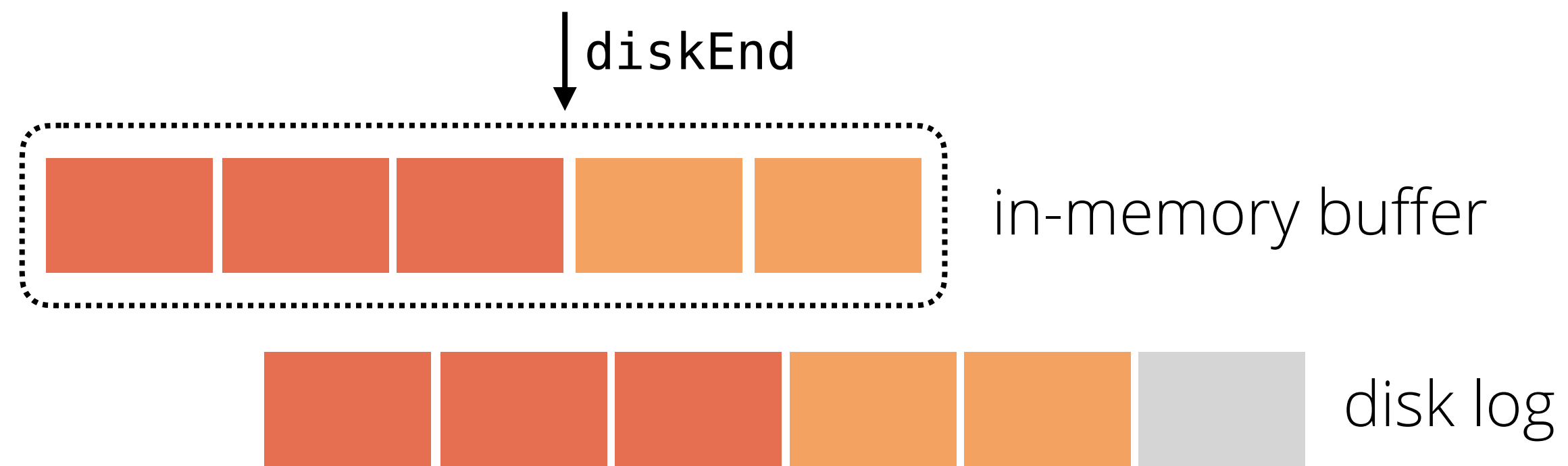


grab newBufs (orange writes to be written)

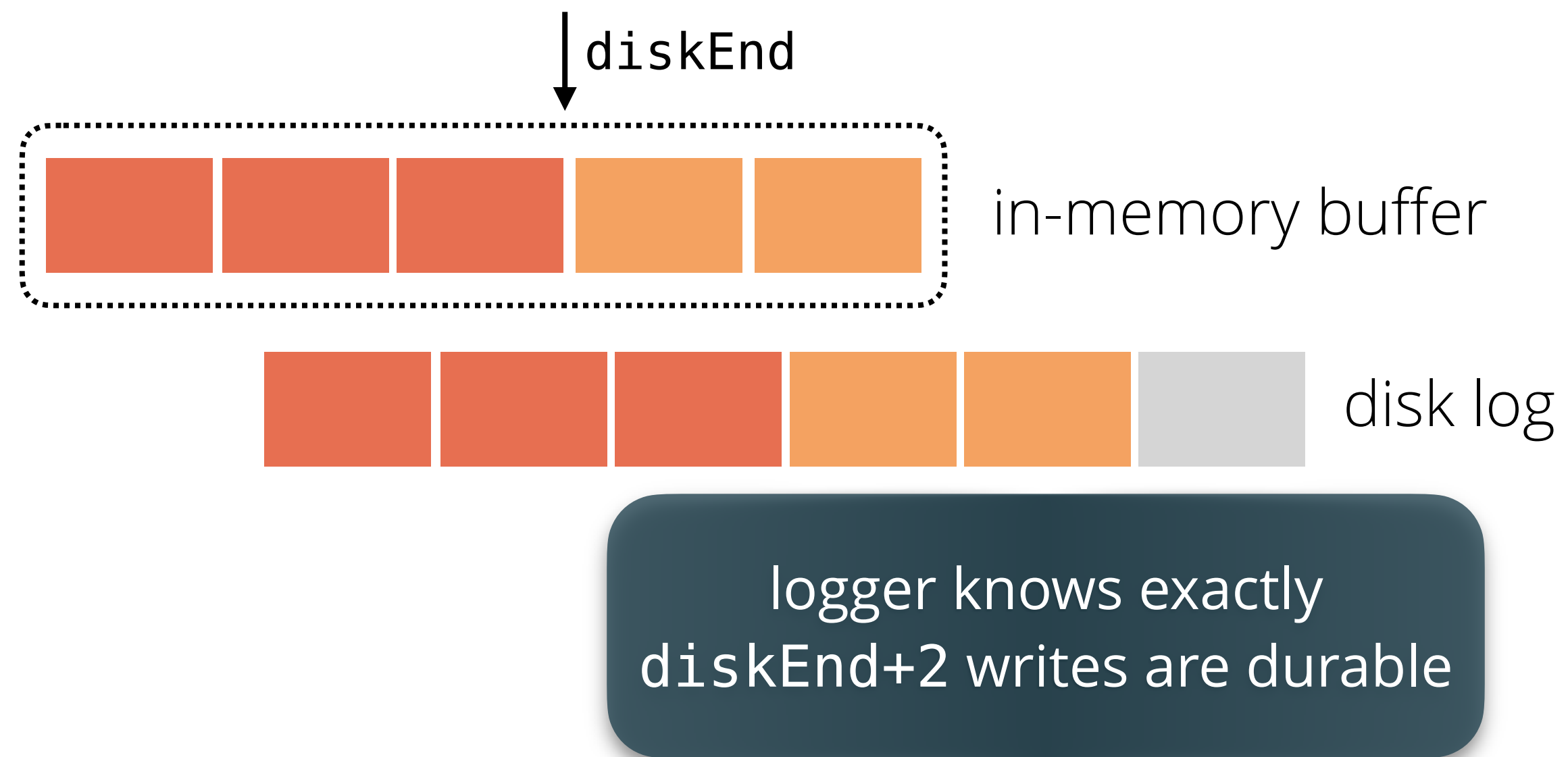
append newBufs to log

after this operation, diskEnd
record that this batch is durable
doesn't reflect what is durable

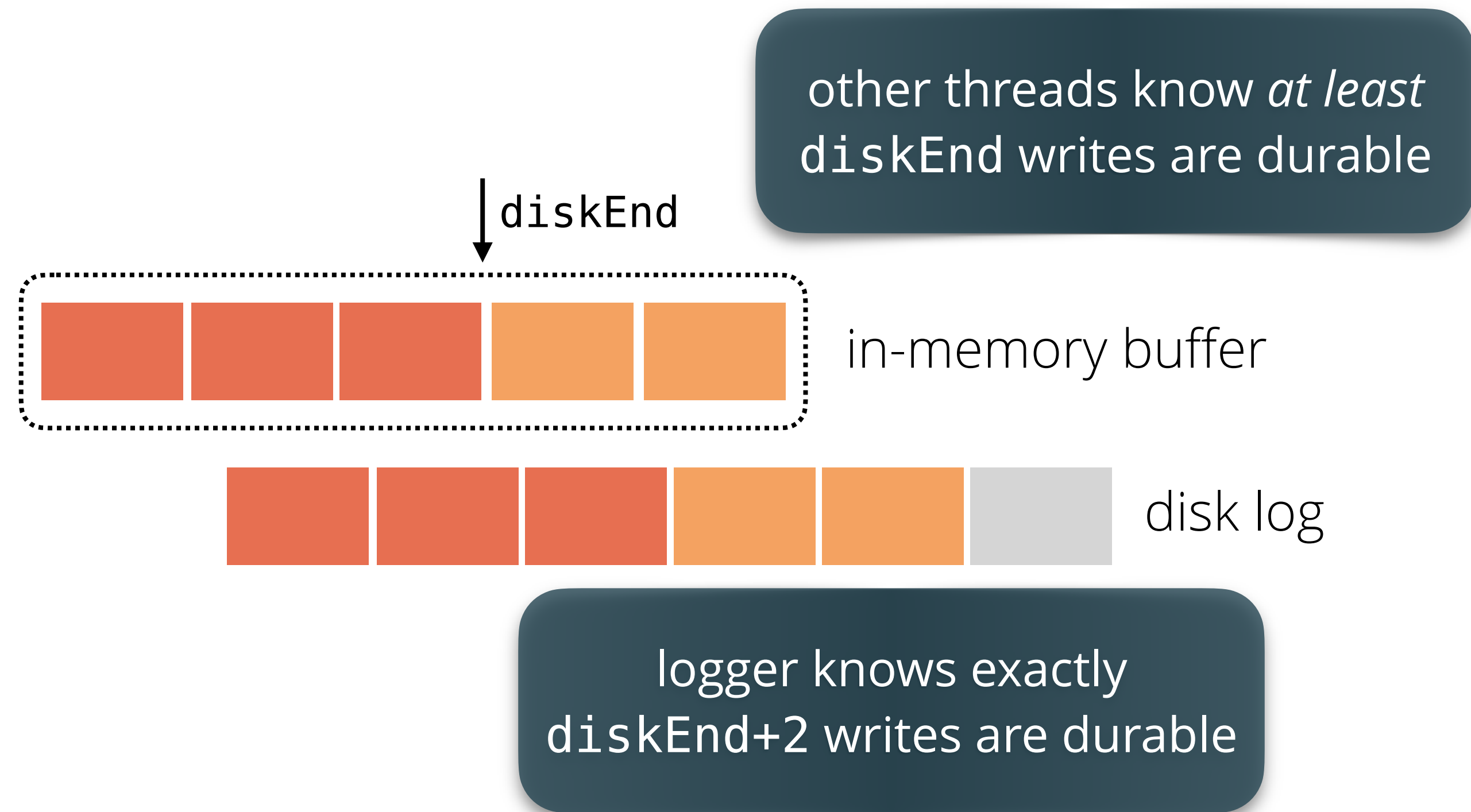
Proof uses general concurrency techniques to reason about lock-free region



Proof uses general concurrency techniques to reason about lock-free region



Proof uses general concurrency techniques to reason about lock-free region



Many other challenges in GoTxn proof

What's the specification for each internal layer?

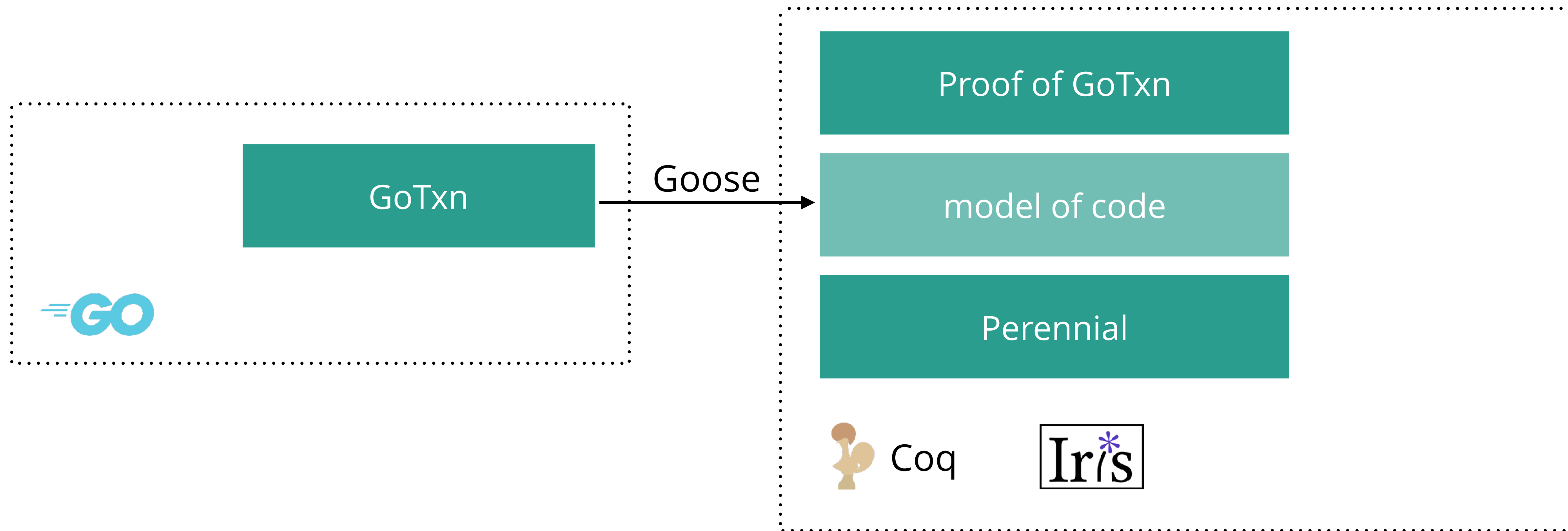
Invariants for lock-free installation, concurrency within a block, two-phase locking

Design and implementation of DaisyNFS

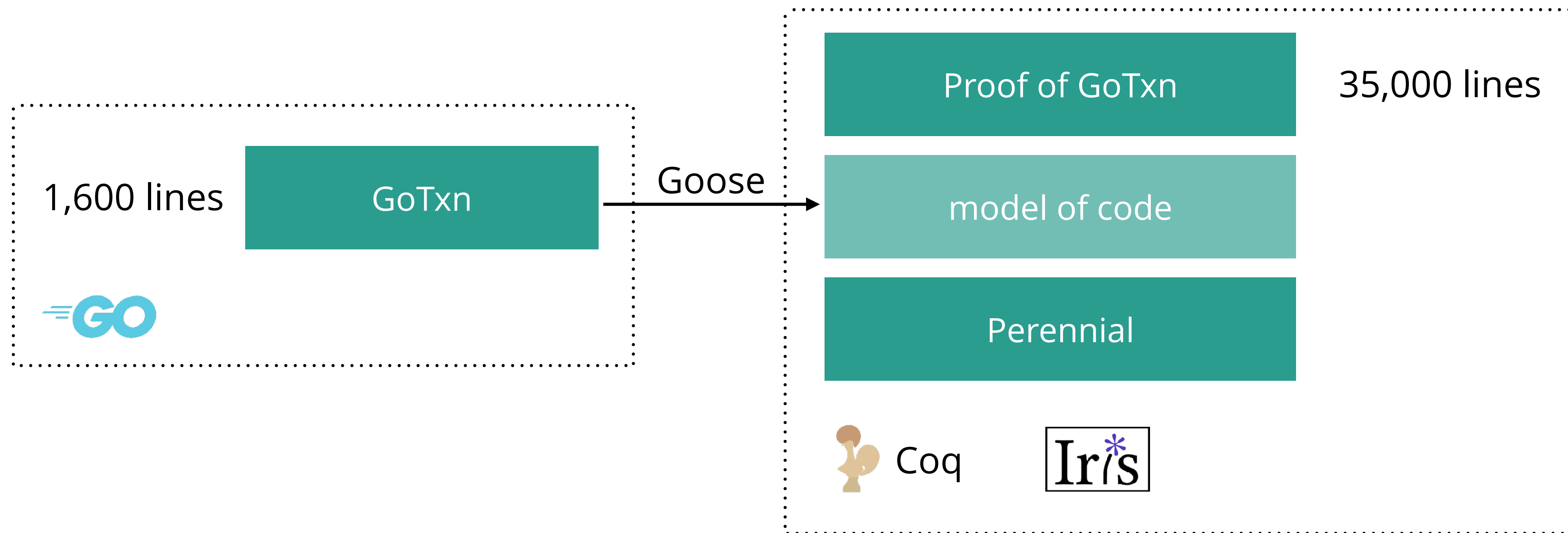
Verifying a high-performance transaction system

Evaluating DaisyNFS

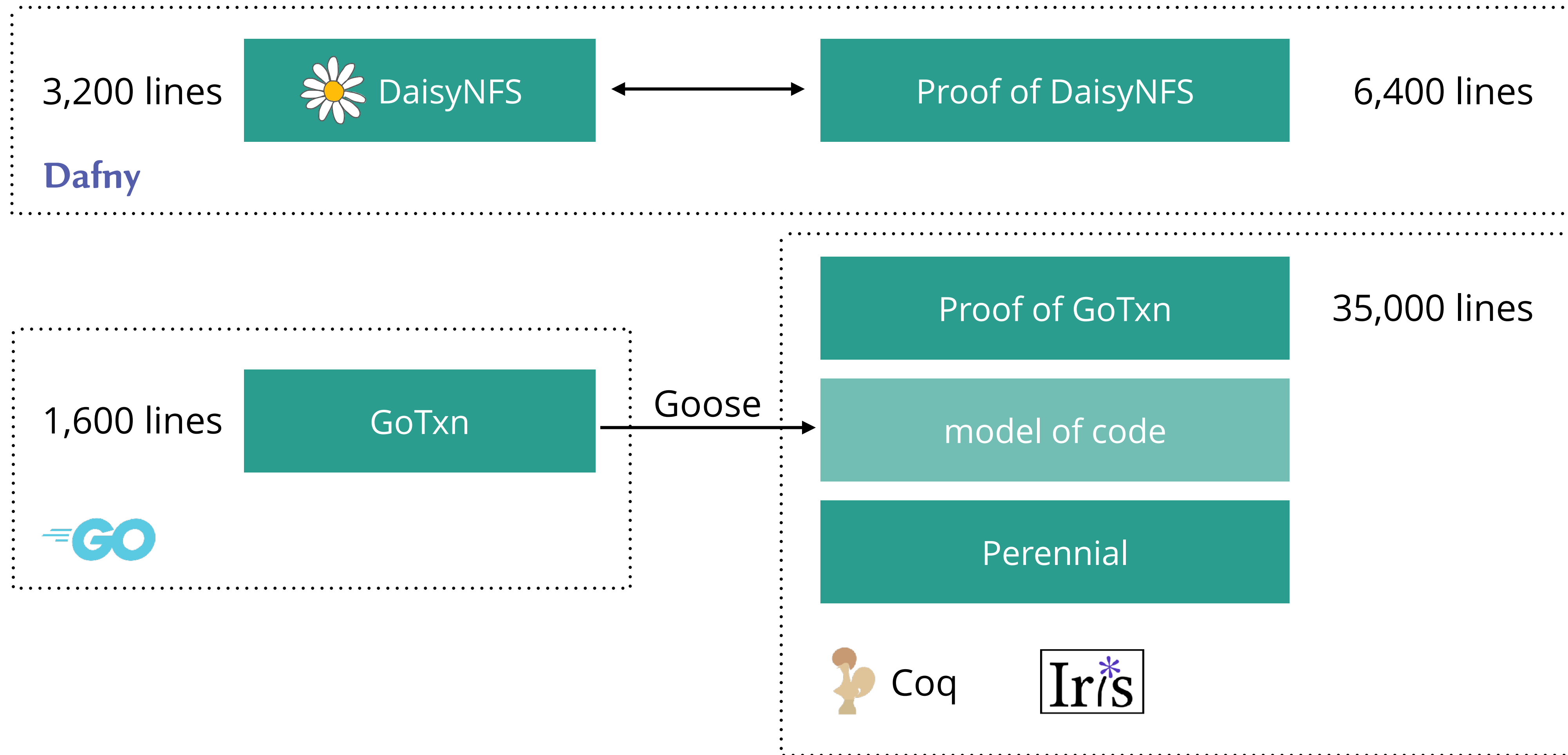
Most of the proof is for GoTxn



Most of the proof is for GoTxn



Most of the proof is for GoTxn



Limitations

Limitations

limit performance



Only synchronous Commit

Must use transactions from Dafny

Limitations

limit performance

Only synchronous Commit

Must use transactions from Dafny

limits to proof

Could still have deadlock

Linking theorem proven on paper

Proof assumptions

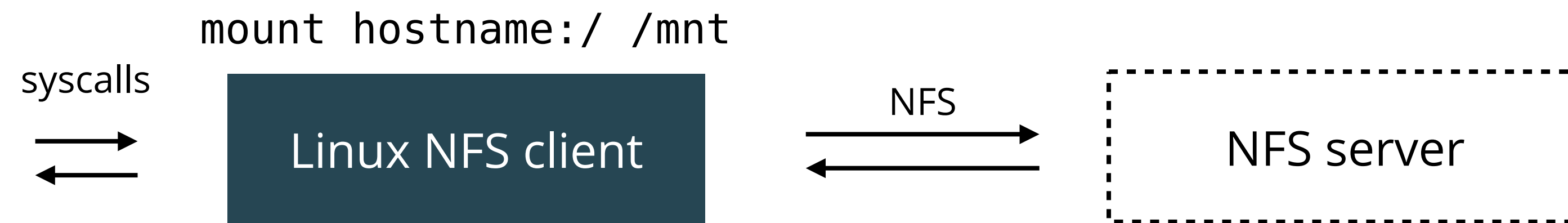
We assume that:

Goose accurately models Go

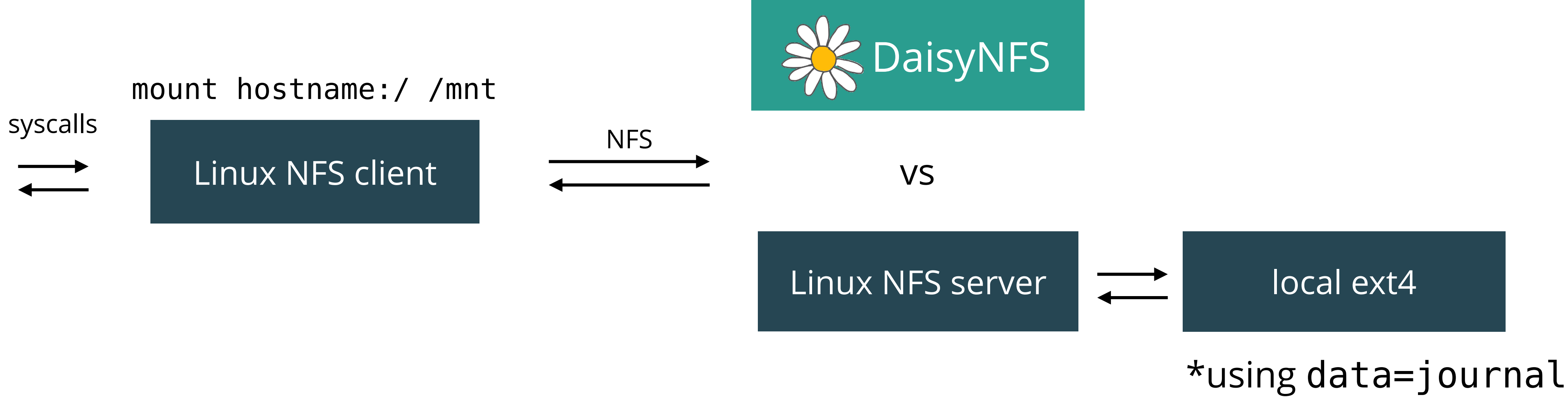
Disk has atomic 4KB reads and writes

NFS specification is written correctly

Evaluate performance using an NFS client



Compare against Linux NFS

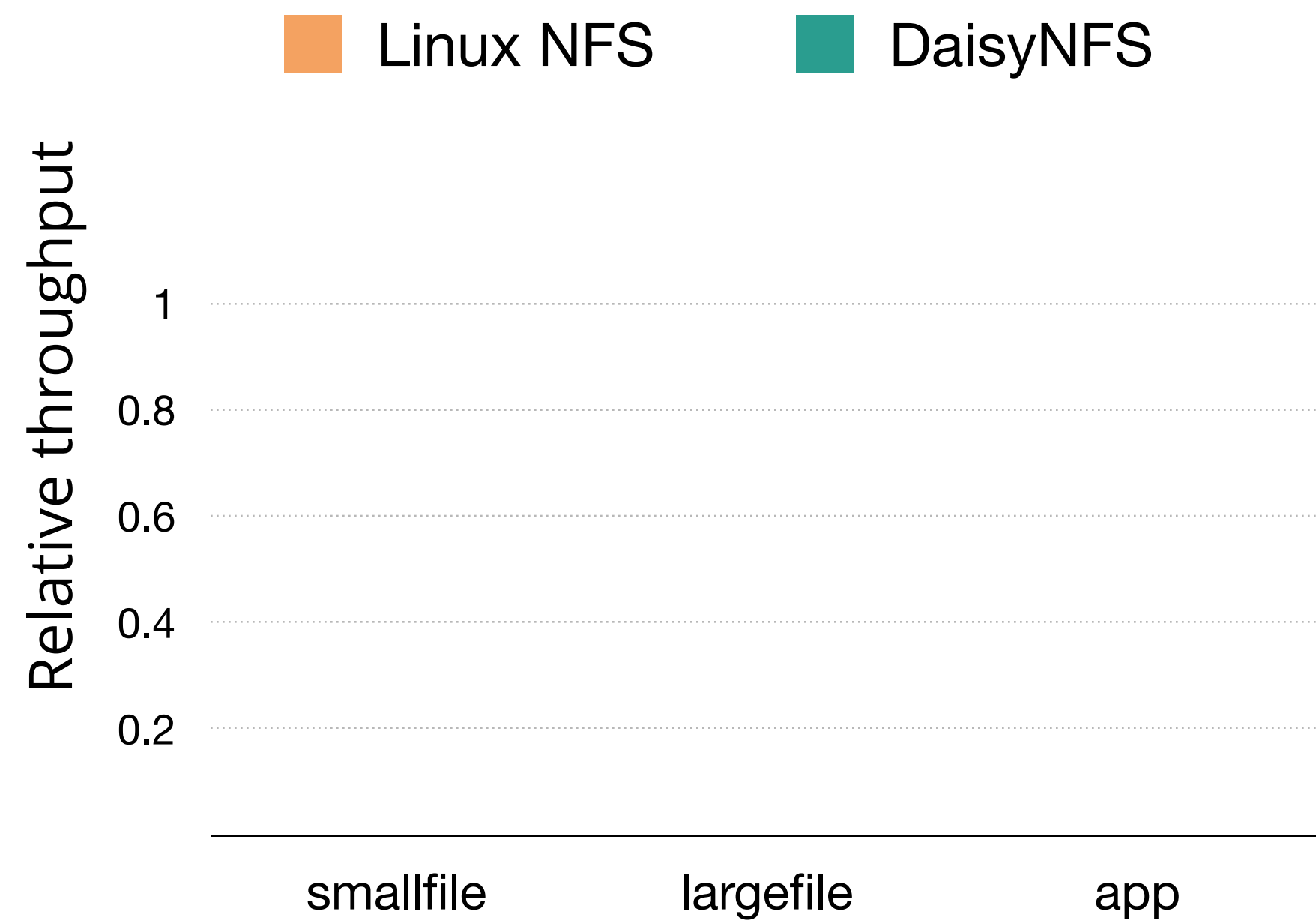


Performance evaluation setup

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

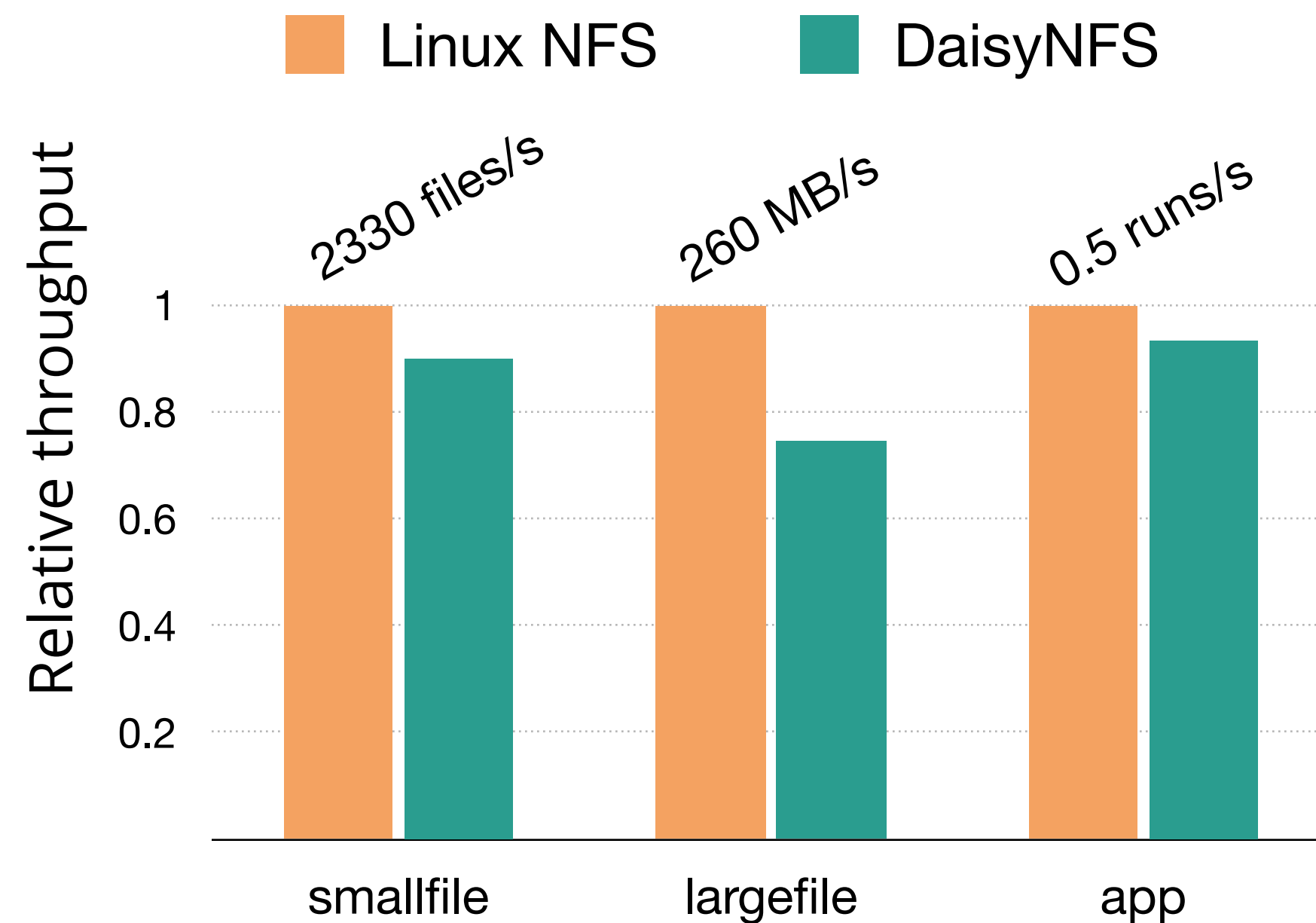
Benchmarks:

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

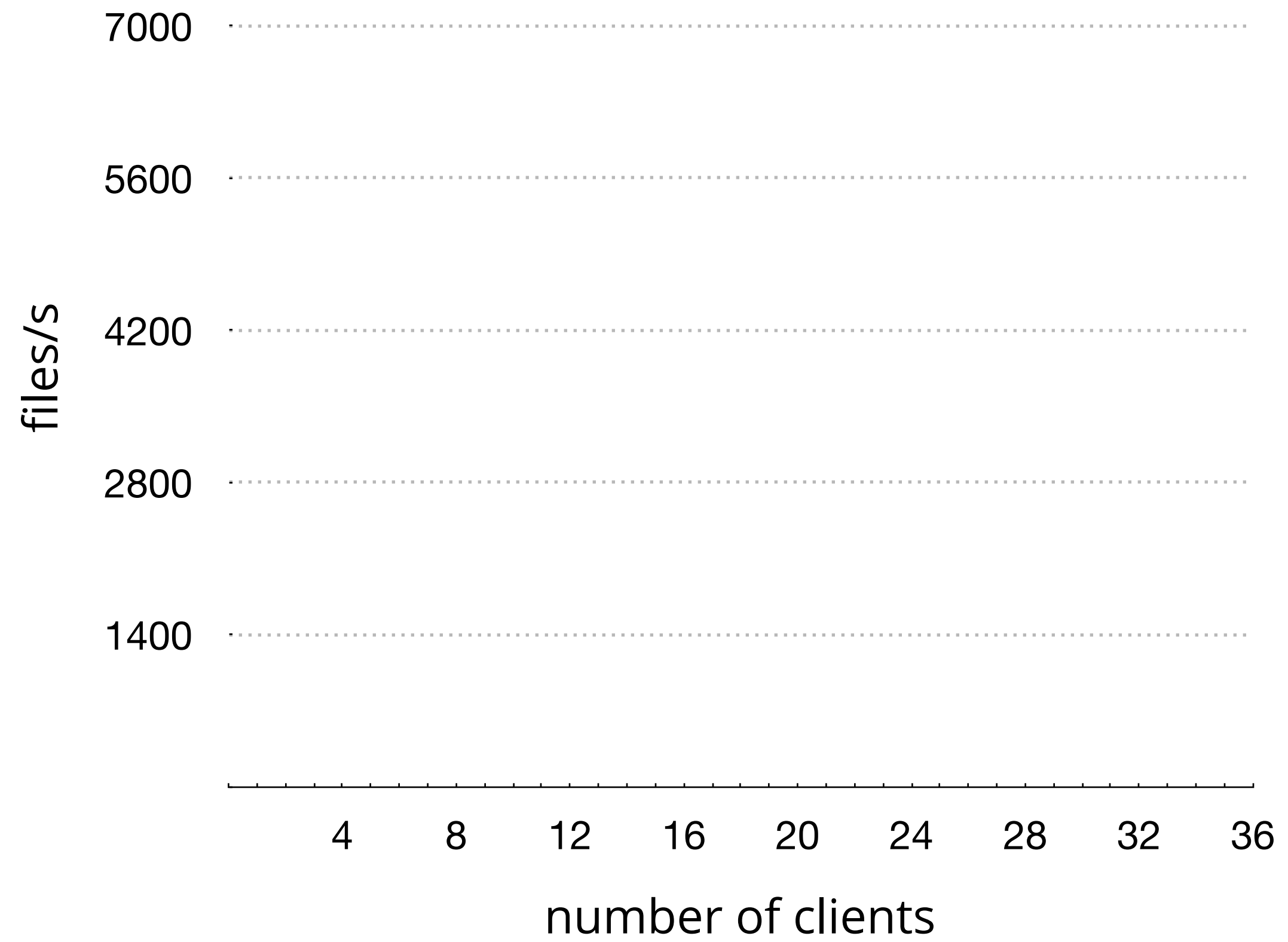


Compare DaisyNFS throughput to Linux,
running on an in-memory disk

DaisyNFS gets comparable performance even with a single client

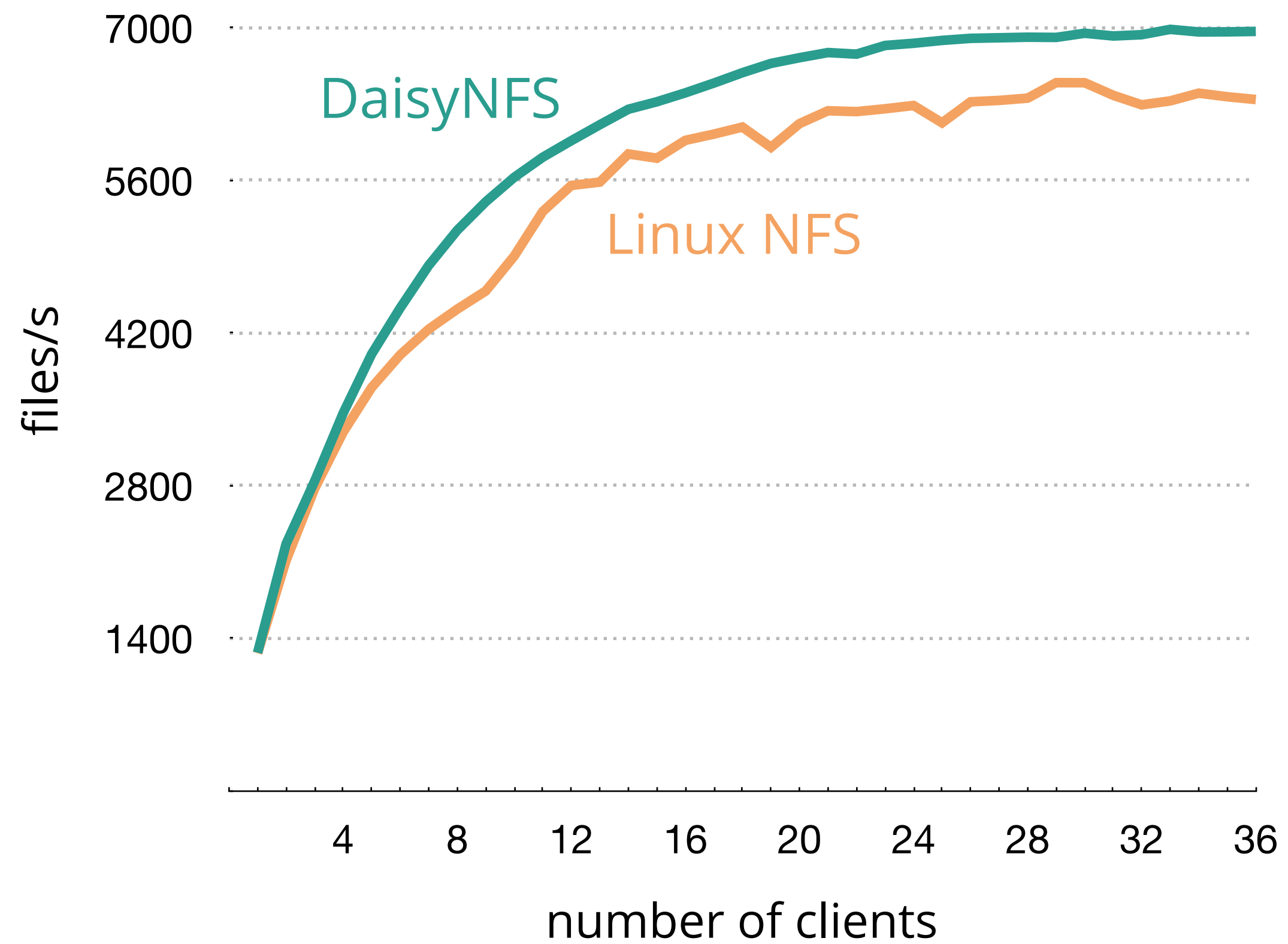


Compare DaisyNFS throughput to Linux, running on an in-memory disk



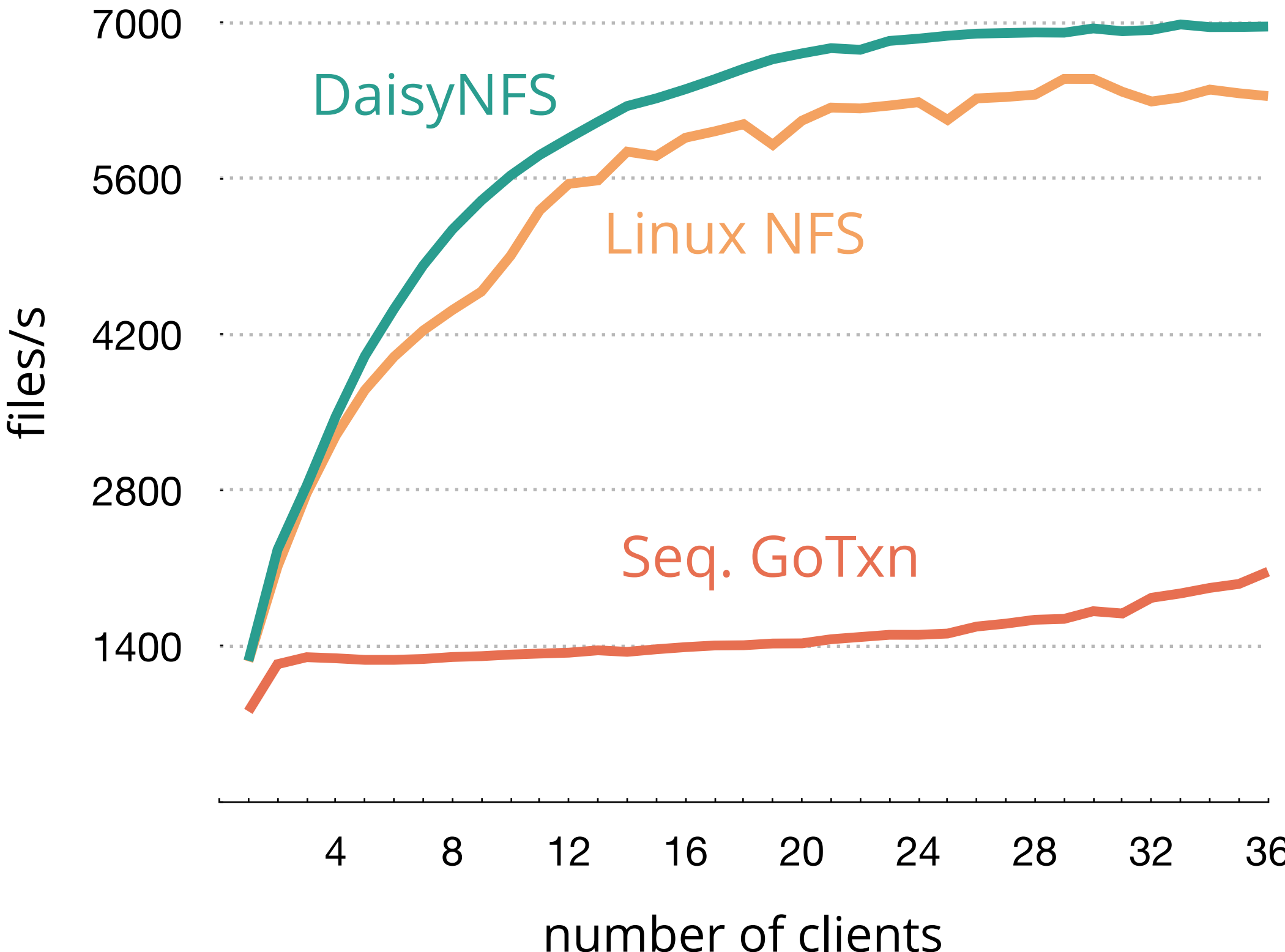
Run smallfile with many clients on an NVMe SSD

DaisyNFS can take advantage of multiple clients



Run smallfile with many clients on an NVMe SSD

Concurrency in the transaction system matters



Seq. GoTxn is DaisyNFS but with locks around tricky concurrent parts of WAL



Related work

Related work

crash safety and concurrency:

Flashix concurrent file system, ShardStore

crash safety:

FSCQ, Yggdrasil, VeriBetrFS

concurrency:

Concurrent GC, CertiKOS, AtomFS

Other related work

Goose: VST and CH20 for reasoning about C

Perennial: builds on top of Iris

GoTxn: verified transaction algorithms but not systems

DaisyNFS: builds upon DFSCQ and Yggdrasil

Each system is general-purpose

Perennial and Goose can be applied to other storage systems, languages, and hardware

GoTxn can be used to build other storage systems, in Dafny or Perennial

Summary



New foundations (**Perennial** and **Goose**) make verification of concurrent storage systems possible

GoTxn isolates the difficult reasoning so proofs on top use sequential reasoning

Verified **DaisyNFS**, a concurrent, crash-safe file system with performance comparable to Linux



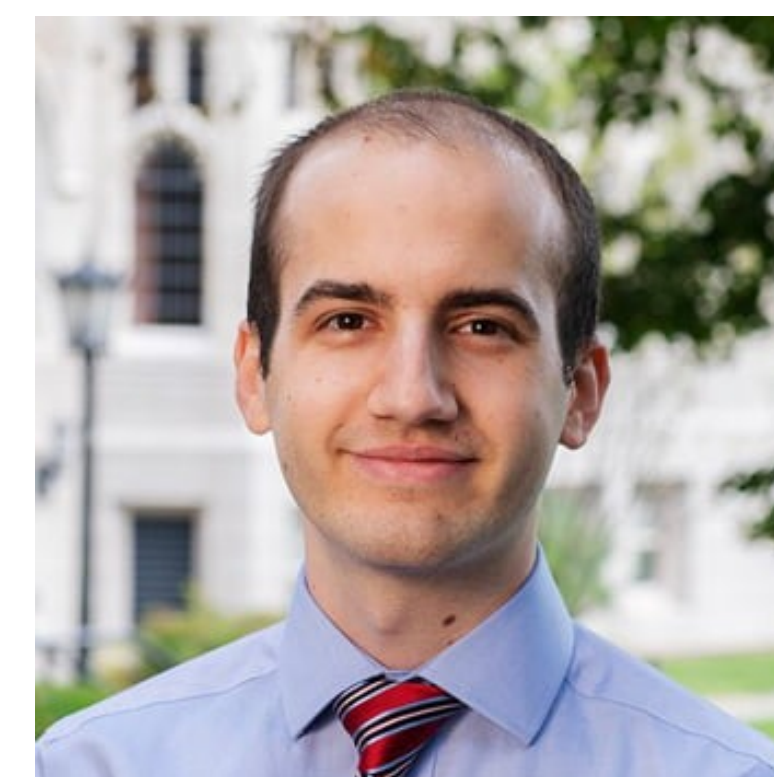
Acknowledgments



Frans Kaashoek



Nickolai Zeldovich



Joe Tassarotti





Neha, Austin, Srivatsa, Julian,
Jelle, Haogang, Shoumik, Cody,
Frank, Amy, Malte, Joe, David,
Jon

Atalay, Anish, Derek, Jonathan,
Akshay, Lily, Josh, Inho, Zain,
Ariel, Kevin, Alex, Upamanyu,
Ralf, Yun-Sheng

Mentees



Daniel Ziegler
Alex Konradi
Lef Ionnadis
Sydney Gibson
Sharon Lin

MIT PL



Sara Achour, Clément Pit-Claudiel, Ben Sherman,
Sam Gruetter, Thomas Bourgeat, and many others

CSC board games

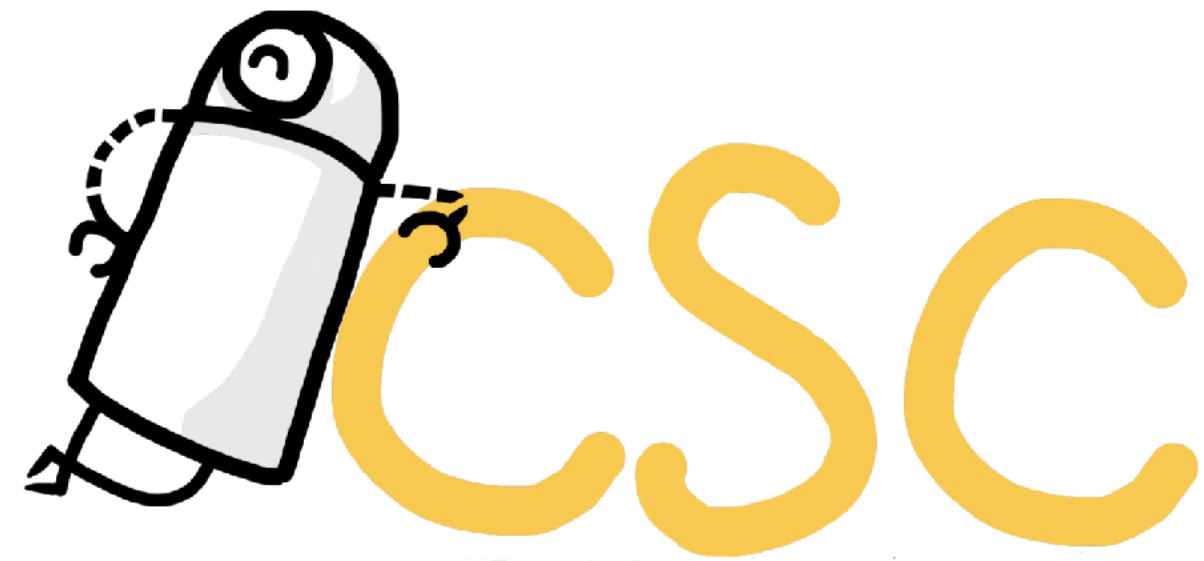


especially Leilani, Nathan, Max, Jon, and Ajay

CSC board games



especially Leilani, Nathan, Max, Jon, and Ajay



250 Elm



Many others