# FSCQ

Sys. verification
  Important code.
  Stable specs.
  Bugs
Crashes
  Logging
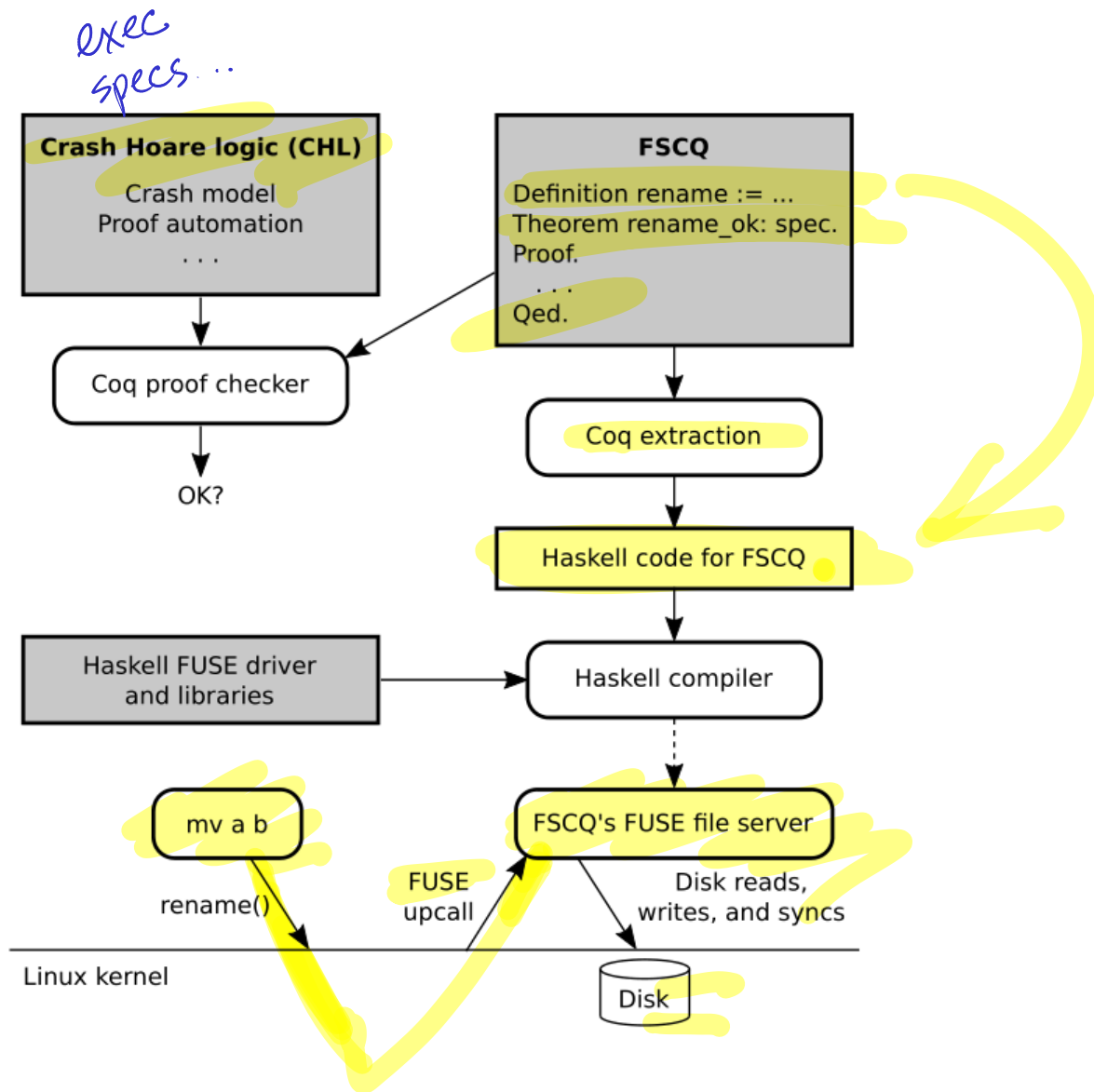  Lab 3, 4.
Non-determinism
  Async disk

Academic papers
  Prototype: show idea.
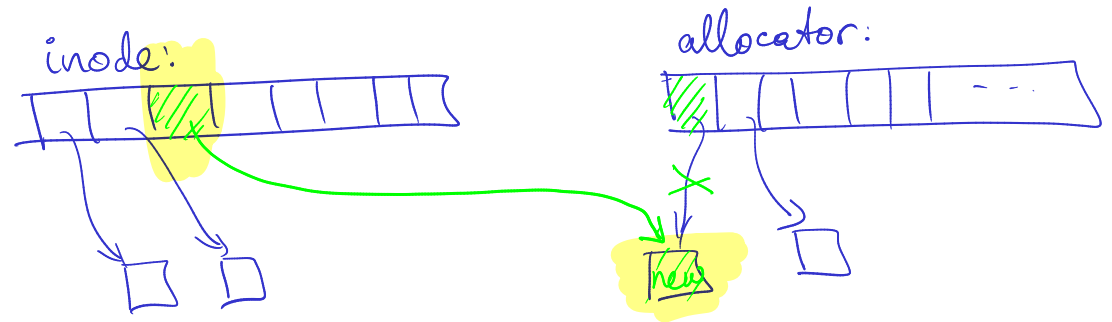
exec
specs . . .

**Crash Hoare logic (CHL)**

Crash model
Proof automation
. . .

**FSCQ**

Definition rename := ...
Theorem rename_ok: spec.
Proof.
. . .
Qed.

Coq proof checker

OK?

Coq extraction

Haskell code for FSCQ

Haskell FUSE driver
and libraries

Haskell compiler

mv a b

FSCQ's FUSE file server

rename()

FUSE
upcall

Disk reads,
writes, and syncs

Linux kernel

Disk

# Goal: Crash safety

reboot

OS boots ✓
FS intact ✓
User data present

# Challenge: Crash any time

append to file:   write(f, new data)

inode:

allocator:

## Logging

Never update in-place
Write to log,   apply log.
Recovery:   apply log.

# Verify logging

### Spec?
### ~~Model?~~
#### Disk
#### Crash
#### Recover
### Proof?
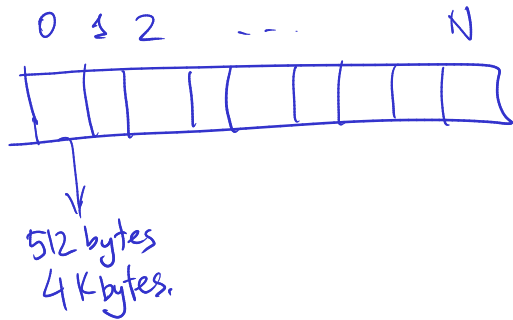#### No crash: func. correctness
##### SibylFS.
### Crashes: ...

# Spec: transactions
### Atomicity: all-or-none updates.
### Durability: persist after crash.
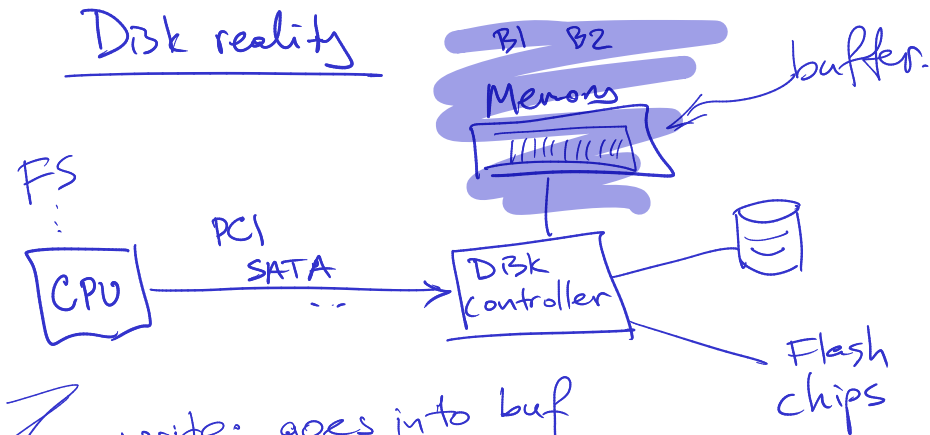### Concurrency: sequential.

# Disk model

0 1 2 ⋯ N

512 bytes
4 K bytes.

$\text{labs} \begin{bmatrix} \text{read}(a) \\ \text{write}(a,v) \end{bmatrix}$

sync()

$\{a \mapsto v_0\} \; \text{write}(a,v) \; \{a \mapsto v\}$

SYNC
DISK

# Disk reality

B1  B2
Memory
buffer.

FS
:
CPU $\xrightarrow{\substack{PCI \\ SATA \\ -- }}$ Disk controller

Flash chips

write: goes into buf
read: buf, or durable
sync: flush buf to durable.
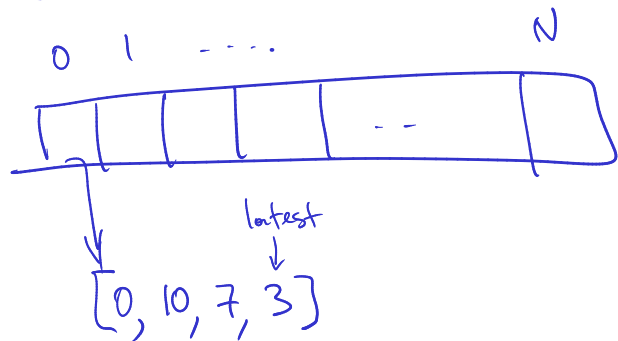crash: lose mem / buf.
background: buf → durable, without sync.

Q: FSCQ model?

write(a, B1)
write (a, B2)
crash

read (a) → B0
B2
B1

# Async disk model

```
 0  1   ....           N
┌──┬──┬──┬──┬─────┬──┐
│  │  │  │  │ --  │  │
└──┴──┴──┴──┴─────┴──┘
   │        │
   │      latest
   ↓        ↓
```

$[0, 10, 7, 3]$

```
    { a ↦ [0] }                 ⟶        a ↦ 0.
        write(a, 10);
    { a ↦ [0, 10] }
        read(a);
    { r = last([0,10])  ⟹  a ↦ [0,10] }
        write(a, 11);                  crash?  → a ↦ [0]
    { a ↦ [0, 10, 11] }                       → a ↦ [10]
        sync()                                → a ↦ [11]
    { a ↦ [11] }           crash   → a ↦ [11].
```

# Crashes in code

1. <u>Where could we crash?</u>

2. How to recovery?

1 = crash condition.

```
    { a ↦ [0] }
        write(a, 10)                    Atomic
    { a ↦ [0, 10] }                     Sector
        CRASH                           write.
    { a ↦ [0]      ∨
      a ↦ [0, 10]      }
```
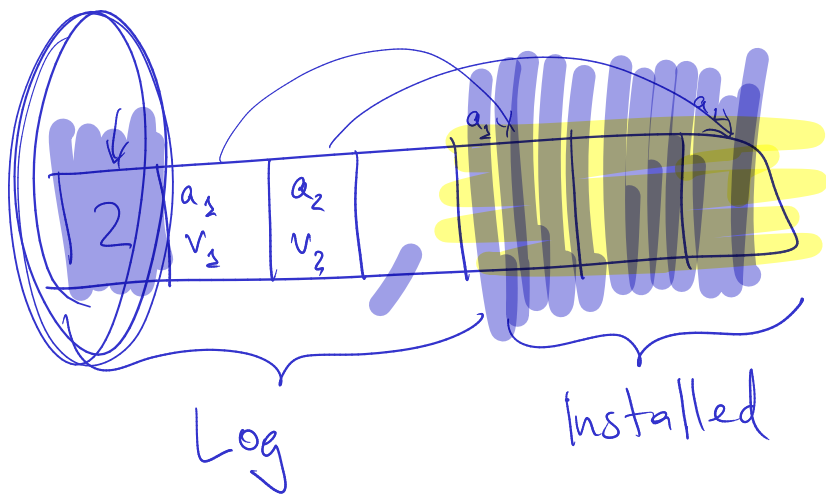
2 = recovery.

reboot: choose one block val
for each addr.
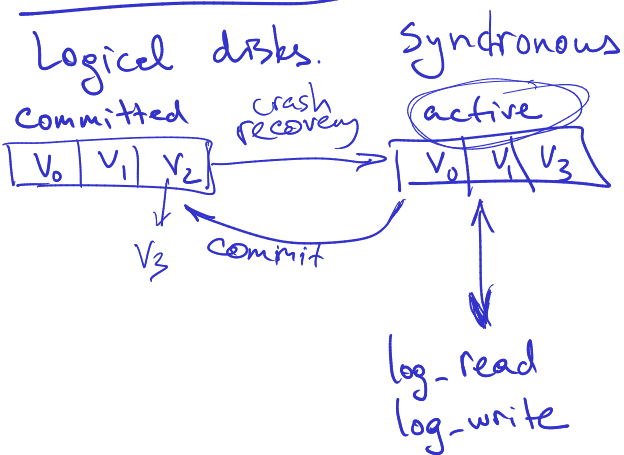
# Logging spec

```
def atomic_two_write(a1, v1, a2, v2):
    log_begin()
    log_write(a1, v1)    ←
    log_write(a2, v2)    ←
    log_commit()
```



Log          Installed

## Abstract state

Logical disks.          Synchronous

Committed    crash recovery    active

State:   $|V_0|V_1|V_2|$  →  $|V_0|V_1|V_3|$

$V_3$  commit
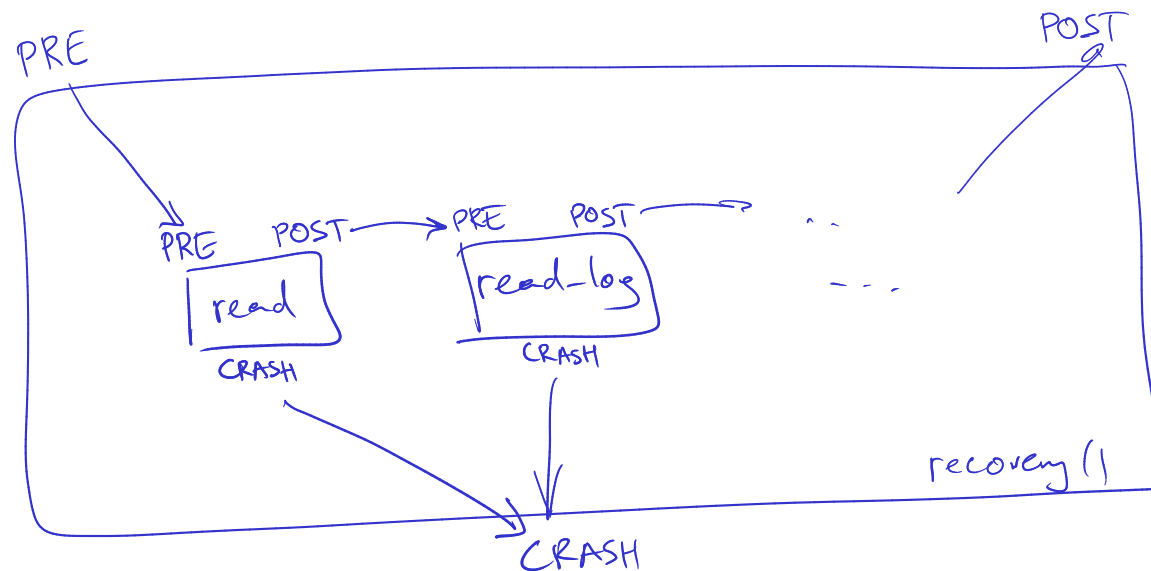
log-read
log-write

Type of state          Abs. state

log_rep(ActiveTxn, $start\_state$, $cur\_state$) :=
   COMMITBLOCK $\mapsto \langle 0, \varnothing \rangle \star (\forall a,\ start\_state[a] = v \rightarrow a \mapsto \langle v, \varnothing \rangle)$
   $\wedge$ replay($start\_state$, inMemoryLog) = $cur\_state$

full disk (async)

```
def recovery():
  r = read(LOG_HDR)
  if r != 0:
    l = read_log()
    for a, v in l:
      write(a, v)
    sync()
    write(LOG_HDR, 0)
```
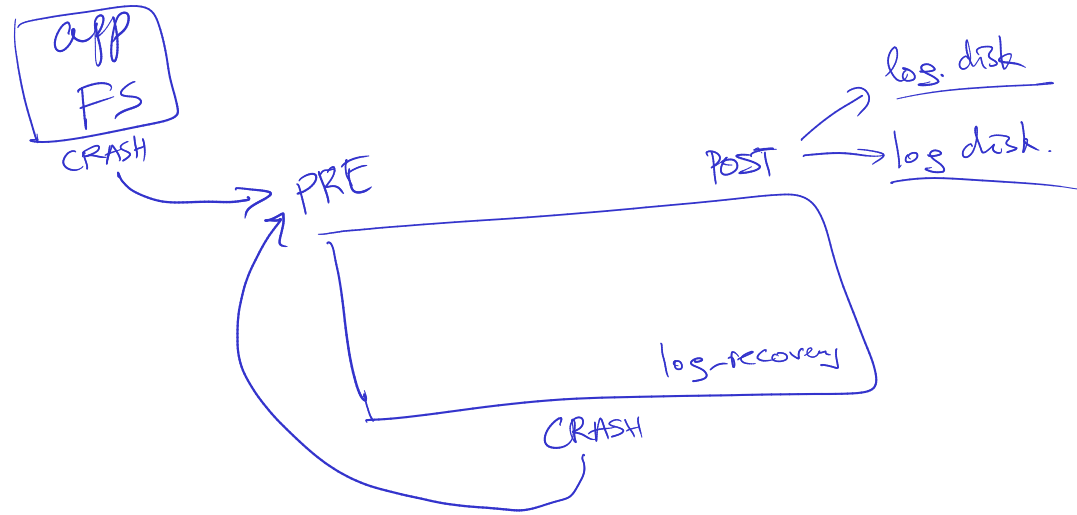


step_proc

PRE

POST

PRE   POST →   PRE   POST

read

read_log

CRASH

CRASH

CRASH

recovery ()

GH2

# Idempotence

SPEC    log_recover()
PRE     **disk**: log_intact(*last_state*, *committed_state*)
POST    **disk**: log_rep(NoTxn, *last_state*) ∨
             log_rep(NoTxn, *committed_state*)
CRASH   **disk**: log_intact(*last_state*, *committed_state*)

app
FS
CRASH

→ PRE

log. disk
POST → log disk.

log-recovery

CRASH

**SPEC** rename(*cwd_ino*, *oldpath*, *newpath*) ≫ fs_recover

**PRE** **disk**: log_rep(NoTxn, *start_state*)

**start_state**: tree_rep(*old_tree*) ∧
  find_subtree(*old_tree*, *cwd*) = *cwd_tree* ∧
  tree_inum(*cwd_tree*) = *cwd_ino*

**POST** **disk**: ((*ret* = (COMPLETED, NoErr) ∨ *ret* = RECOVERED) ∧
    log_rep(NoTxn, *new_state*)) ∨
    ((*ret* = (COMPLETED, Error) ∨ *ret* = RECOVERED) ∧
    log_rep(NoTxn, *start_state*))

**new_state**: tree_rep(*new_tree*) ∧
  *mover* = find_subtree(*cwd_tree*, *oldpath*) ∧
  *pruned* = tree_prune(*cwd_tree*, *oldpath*) ∧
  *grafted* = tree_graft(*pruned*, *newpath*, *mover*) ∧
  *new_tree* = update_subtree(*old_tree*, *cwd*, *grafted*)