

JEPSEN 8

HPTS



Kyle Kingsbury
@aphyr

I break
databases!

Public
API {



Ruby {



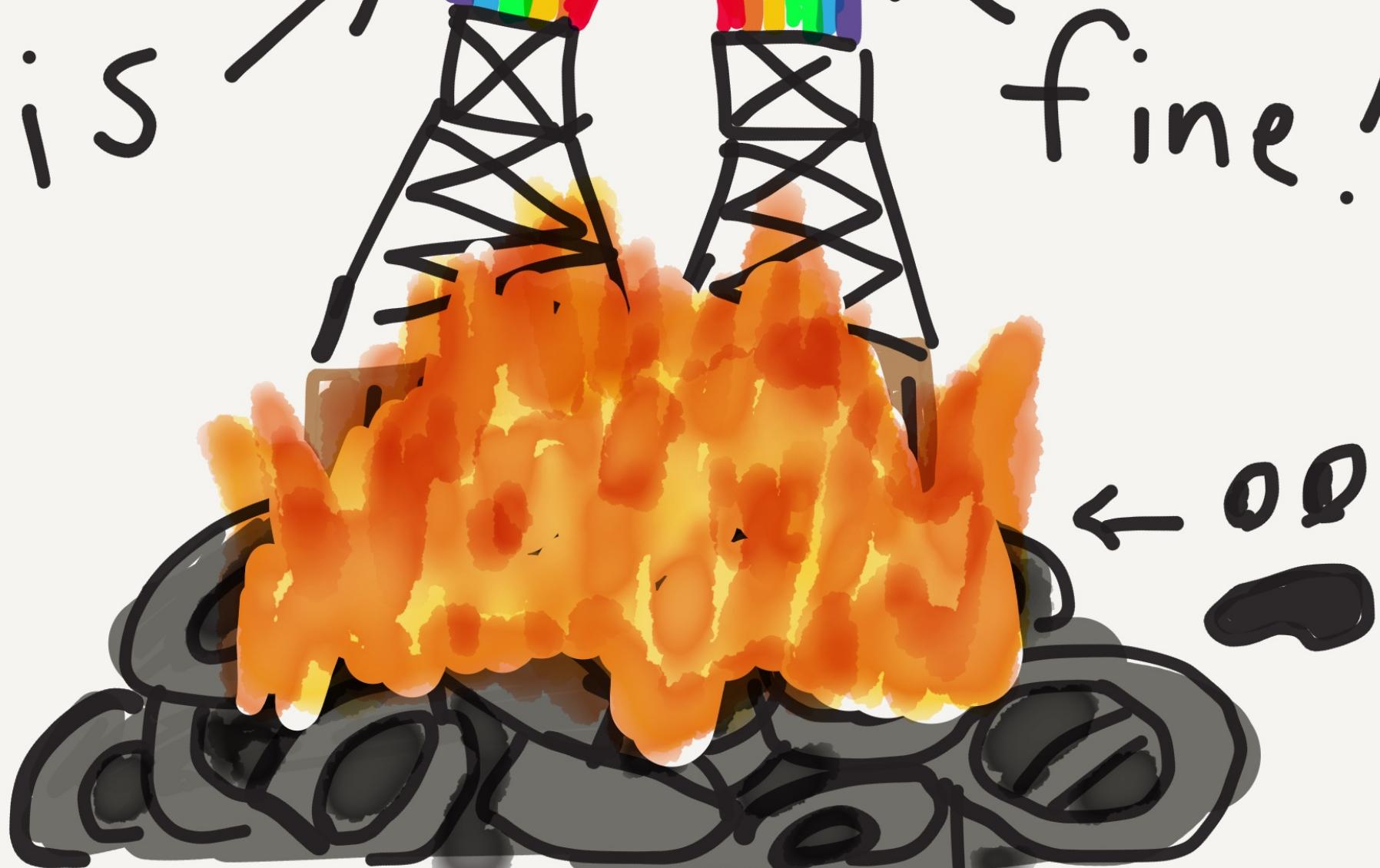
} API code

DBs

Every →
is →

← thing

← fine!

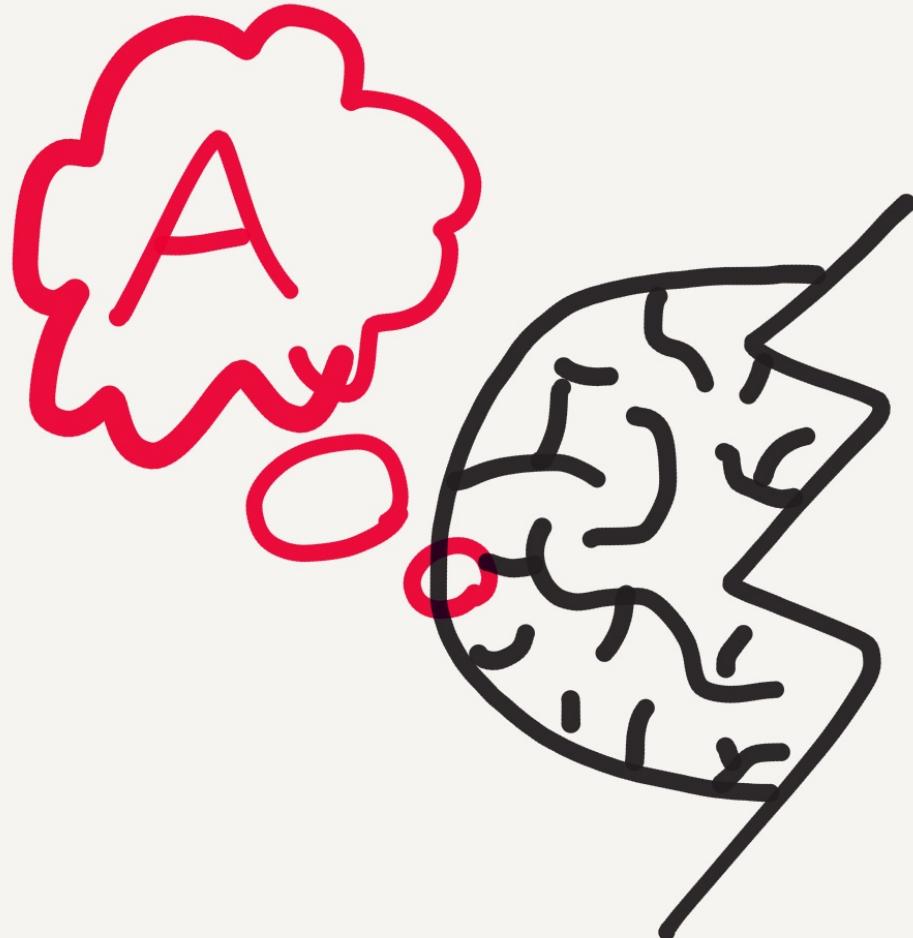


Databases!

Queues!

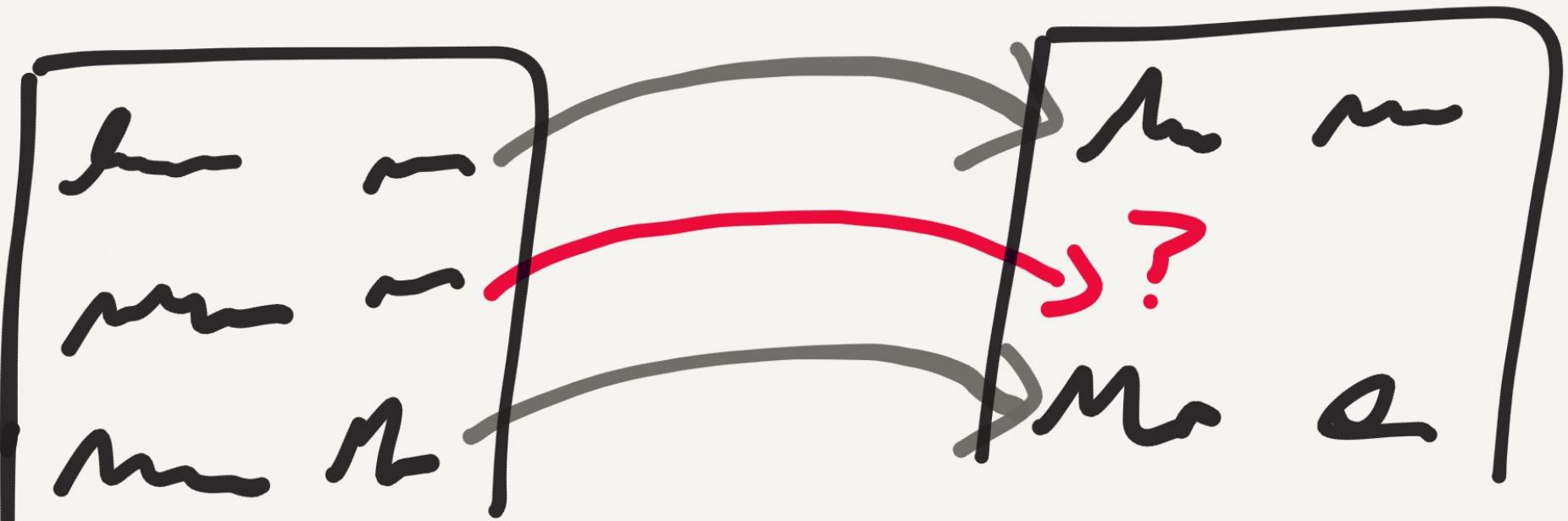
Discovery!

THE HORROR

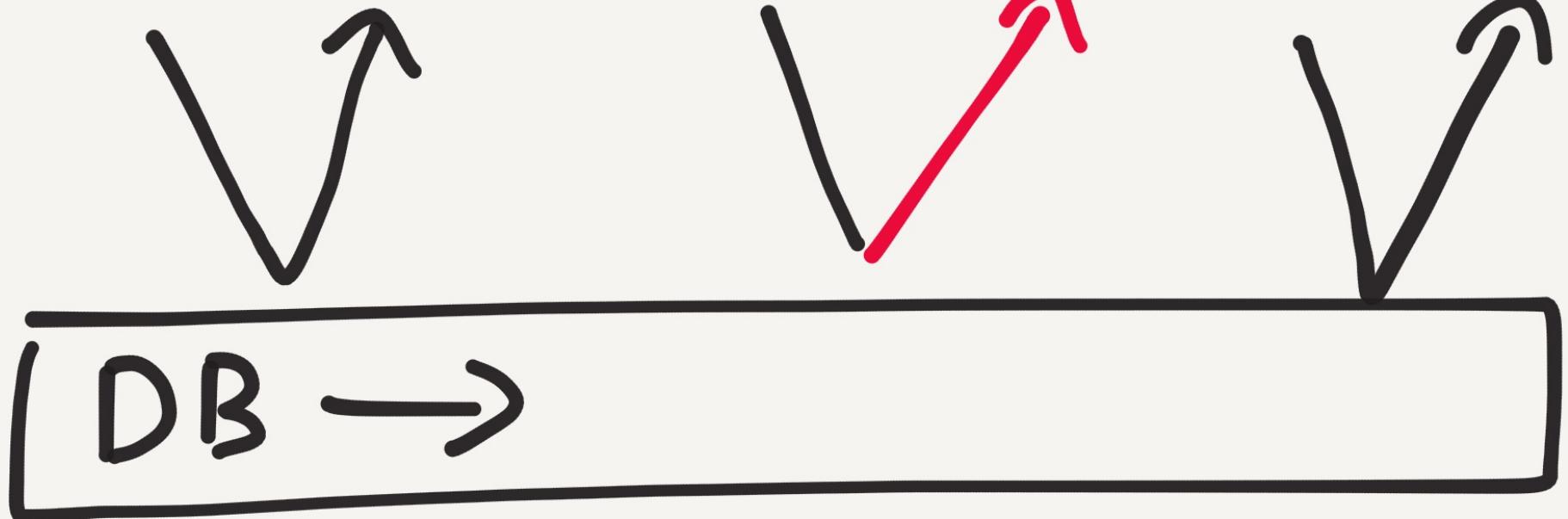


Split Brain

Broken Foreign Keys



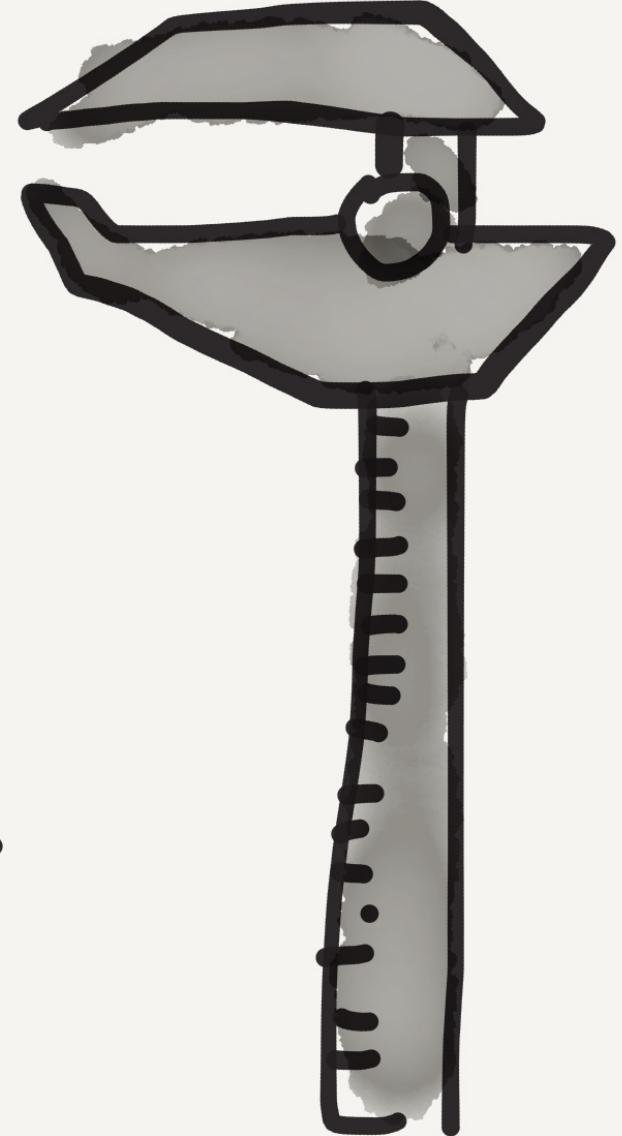
write ok read ⚡ read ok



Anomalies

How do you
know if a
system is safe?

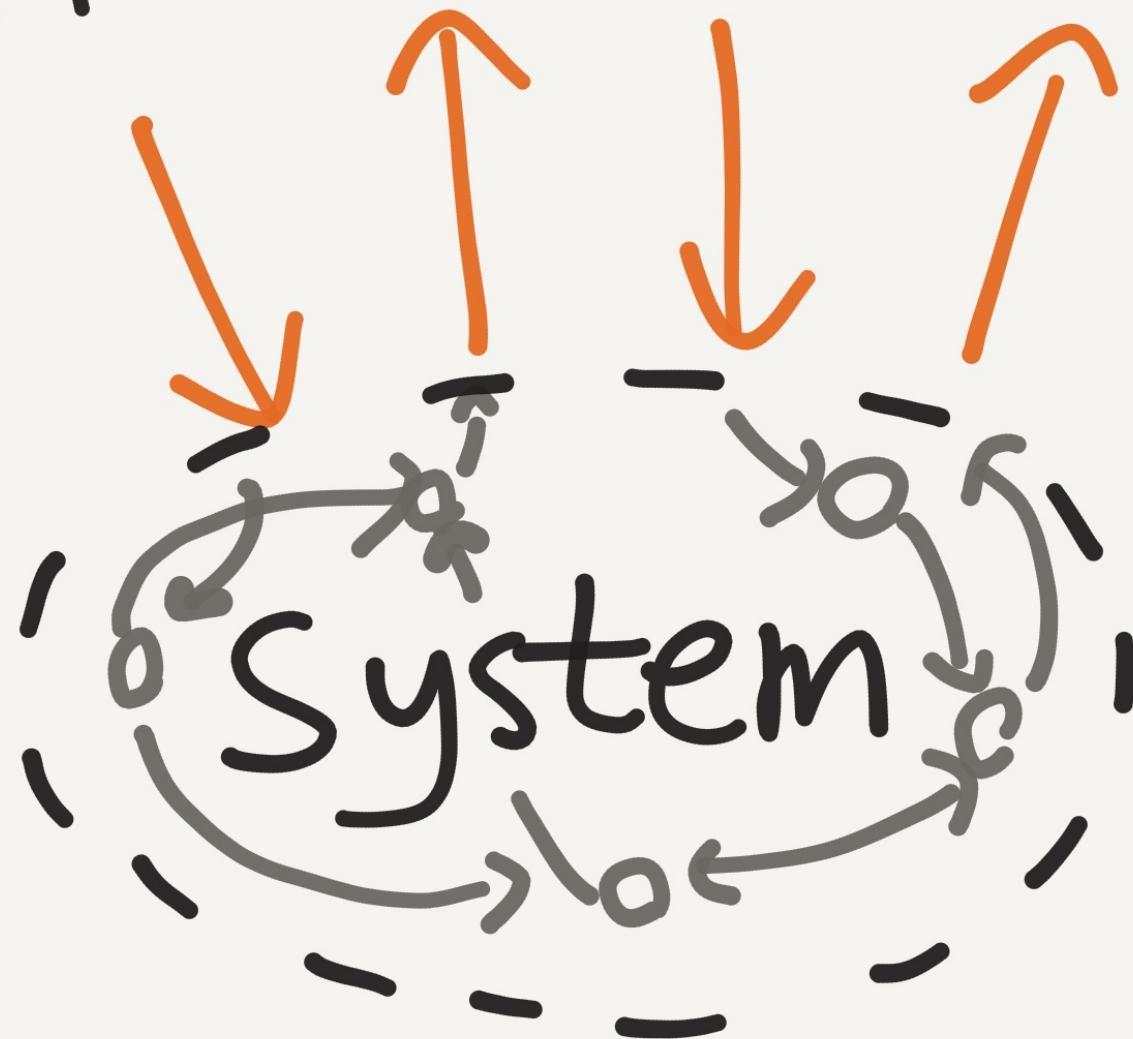
Measure
your
Systems



Jepsen

github.com/aphyr/jepsen

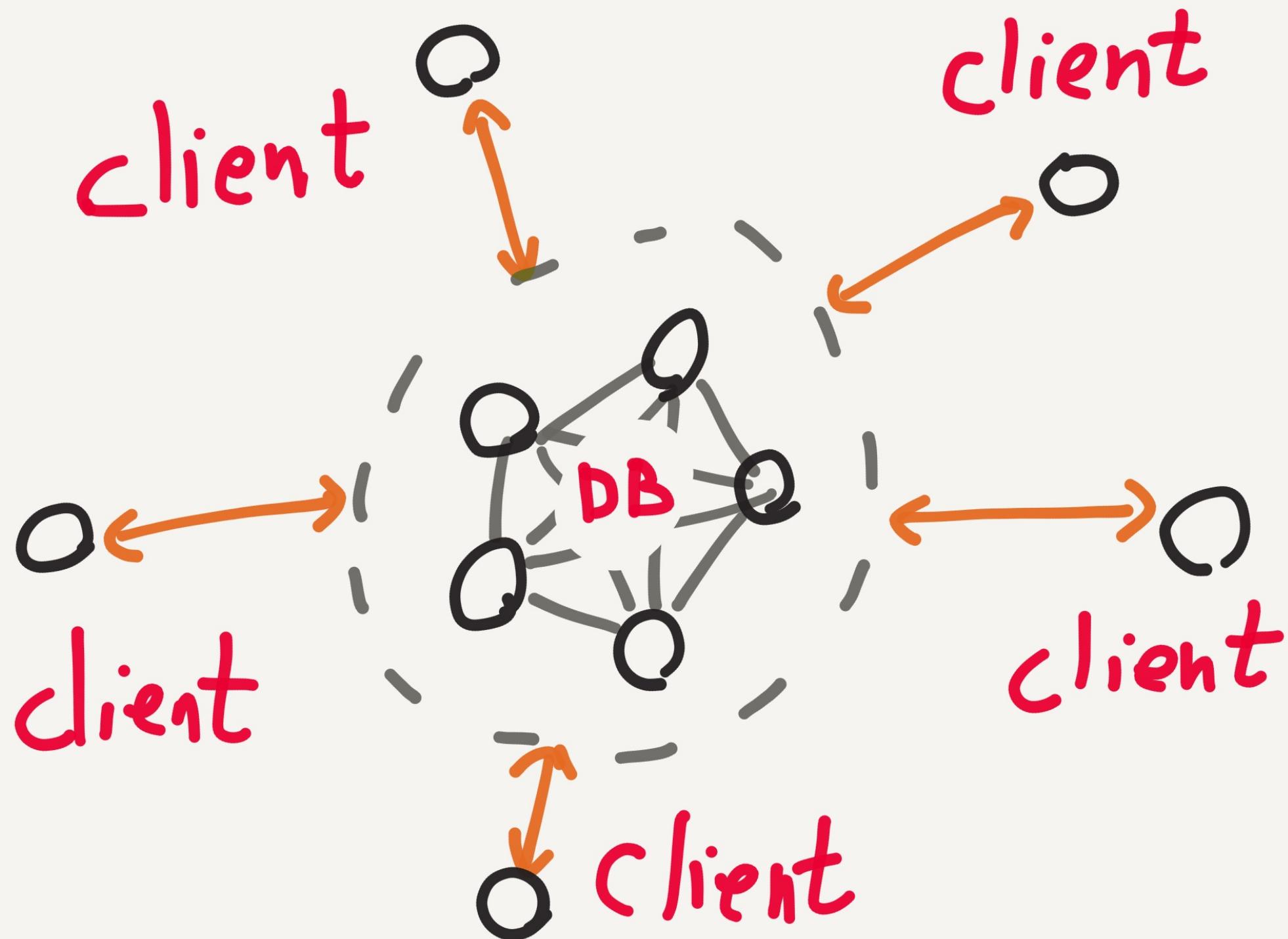
Environment





-INVARIANTS-





client: $w-w'$ $r-r'$ $w-w'$

The diagram illustrates three client requests originating from a single client node. Each request is represented by a horizontal line segment with arrows at both ends, labeled $w-w'$, $r-r'$, and $w-w'$. These segments point towards a central, blurred node labeled "DB".

DB :

client : w — w' w — ...

The diagram shows a central "DB" node receiving three requests from a client. The client node is labeled "client : w — w' w — ...". Three arrows originate from the client and point to the DB node. From the DB node, three arrows point back to the client, labeled w , w' , and w respectively.

Clients Generate
random operations (w)
and apply them
to the system (w')



invoke

ok

invoke

fail

invoke ?

?

info

?

?

?

invoke

ok

invoke

fail

invoke

ok

invoke info

invoke

ok

invoke

fail

invoke

ok

invoke

info

invoke ok

invoke fail

invoke ok

invoke info

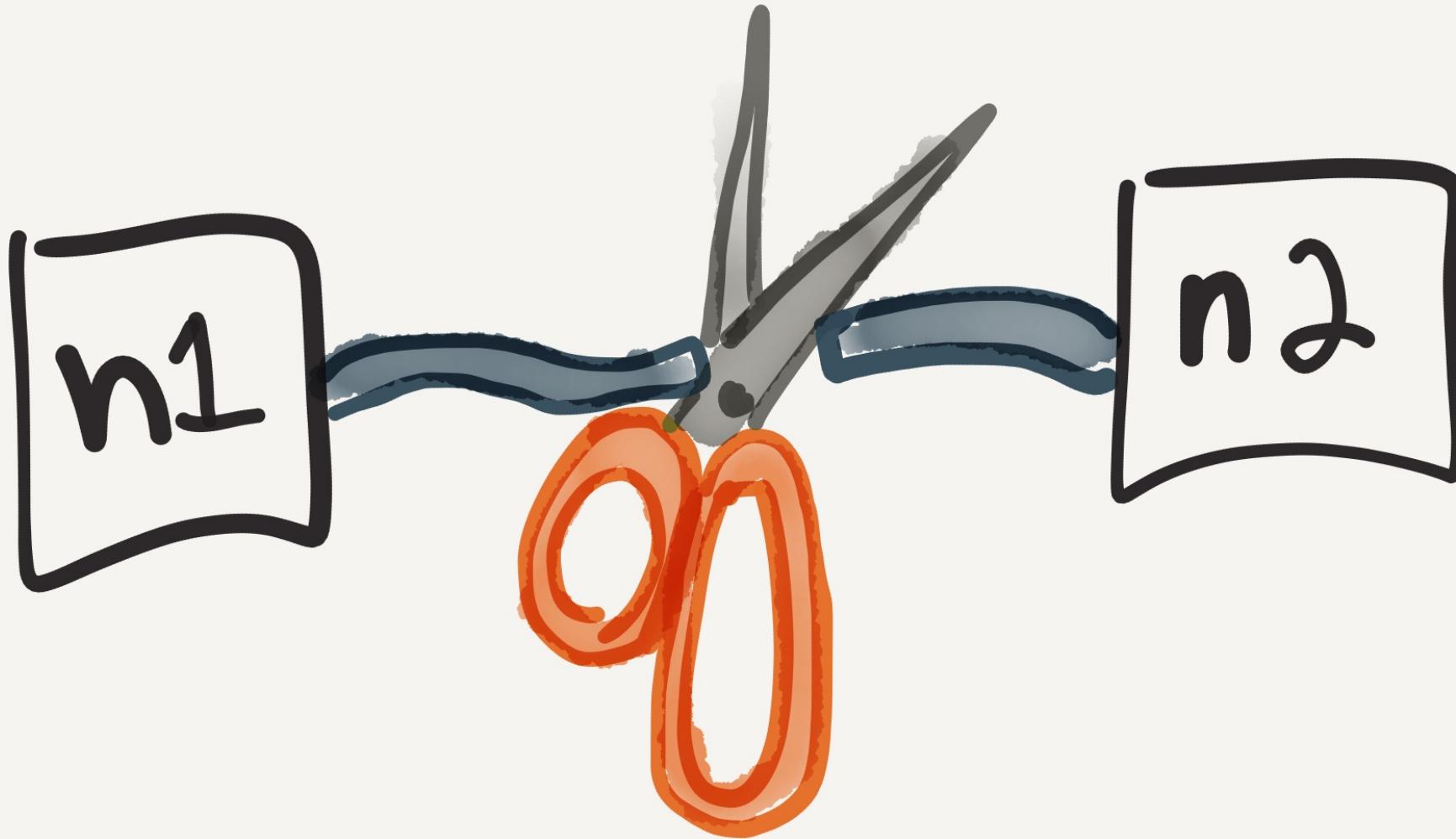
invoke ok

invoke fail

invoke ok

invoke info

1. Generate ops
2. Record history
3. Verify history is
consistent w/ model



Partitions!

“So, what
have you
found?”

Riak

5/13

LWW → lost writes

CREDITs → safe

5/13

Mongo

Data loss at all
Write concerns

5/13

Redis Sentinel

Split brain,
massive write loss

Cassandra

9/13

- LWW write loss
- Row isolation broken
- Transaction deadlock
data loss

9/13

NuoDB

Beat CAP by buffering all requests in IRAM during partition

9/13

Kafka

In-sync Replica Set

could shrink to 0

nodes, causing msg

loss.

9/13

Zookeeper

Works.

etcd / Consul

6/14

Stale reads

6/14

Elastic search

Loses documents in
every class of partition
tested.

6/14

RabbitMQ

Split brain, massive

message loss

5/15

Aerospike

Claims "ACID", was
really LWW.

Elasticsearch 1.5.0

5/15

Still loses data
in every test case

MongoDB 2.6.7

stale reads

dirty reads

Chronos

8/15

- Breaks forever after losing quorum

Percona XtraDB/Galera 9/15

- "Snapshots" weren't
- First-committer-wins
not preserved
- Read locks broken

RethinkDB

1/16

- Basic tests passed
- Reconfiguration could destroy cluster in rare cases

Crate.io

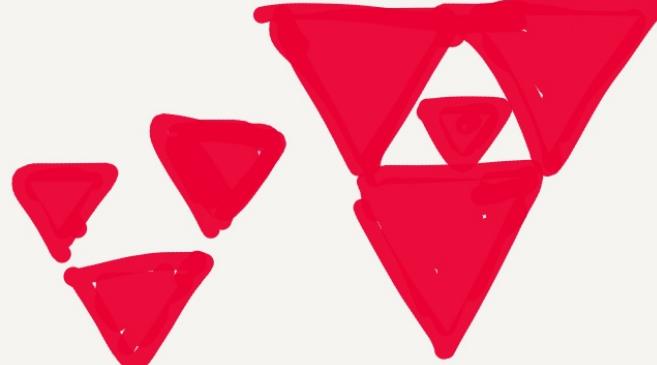
2016

- Stale reads
- Dirty reads
- Lost/corrupt updates
- Lost inserts

Cockroach DB

10/2016

- Double-insert
- Phantom Anomaly



VOLTDB

6.3

Distributed SQL

Database

All data in RAM

(but persistent)

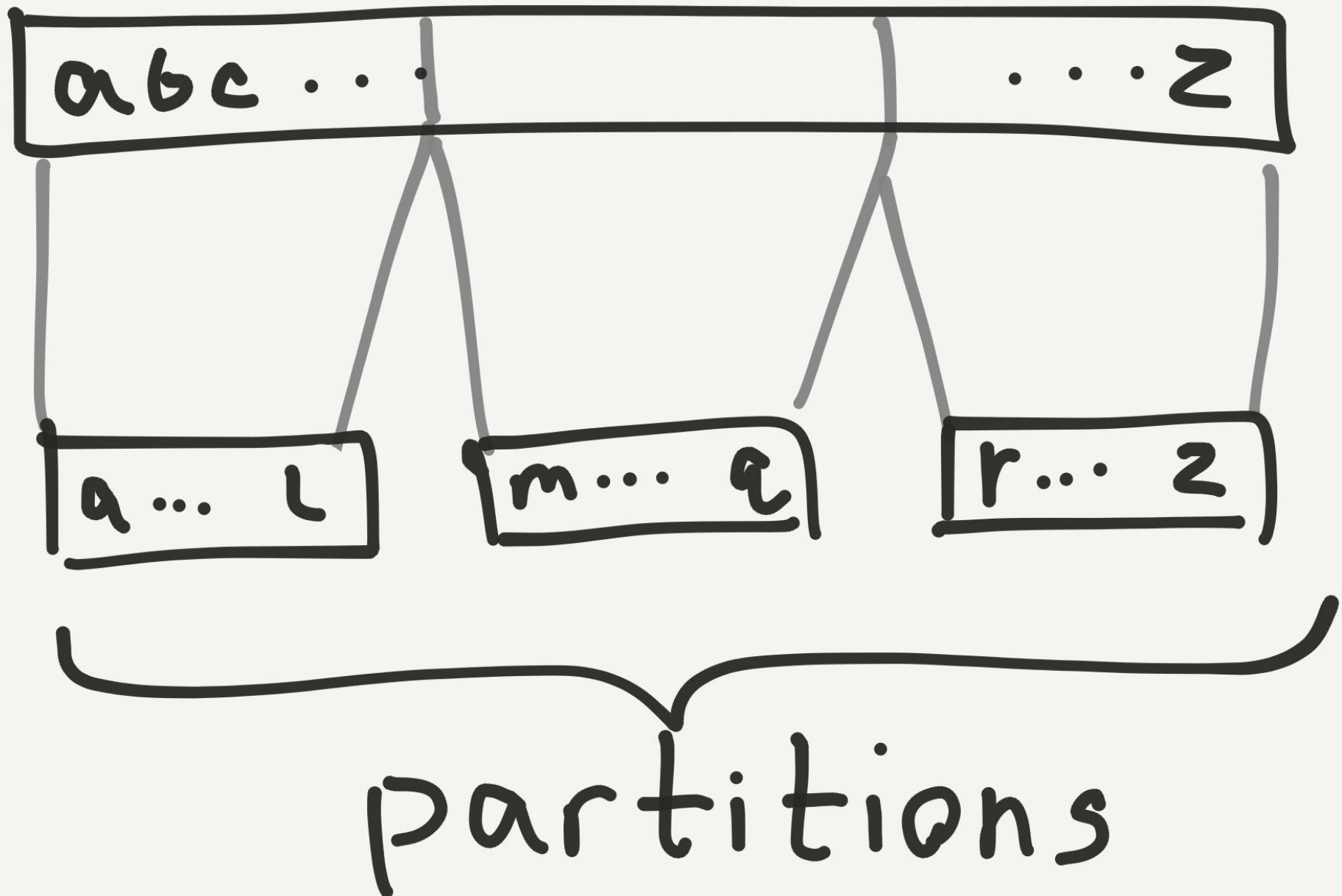
Basic SQL queries



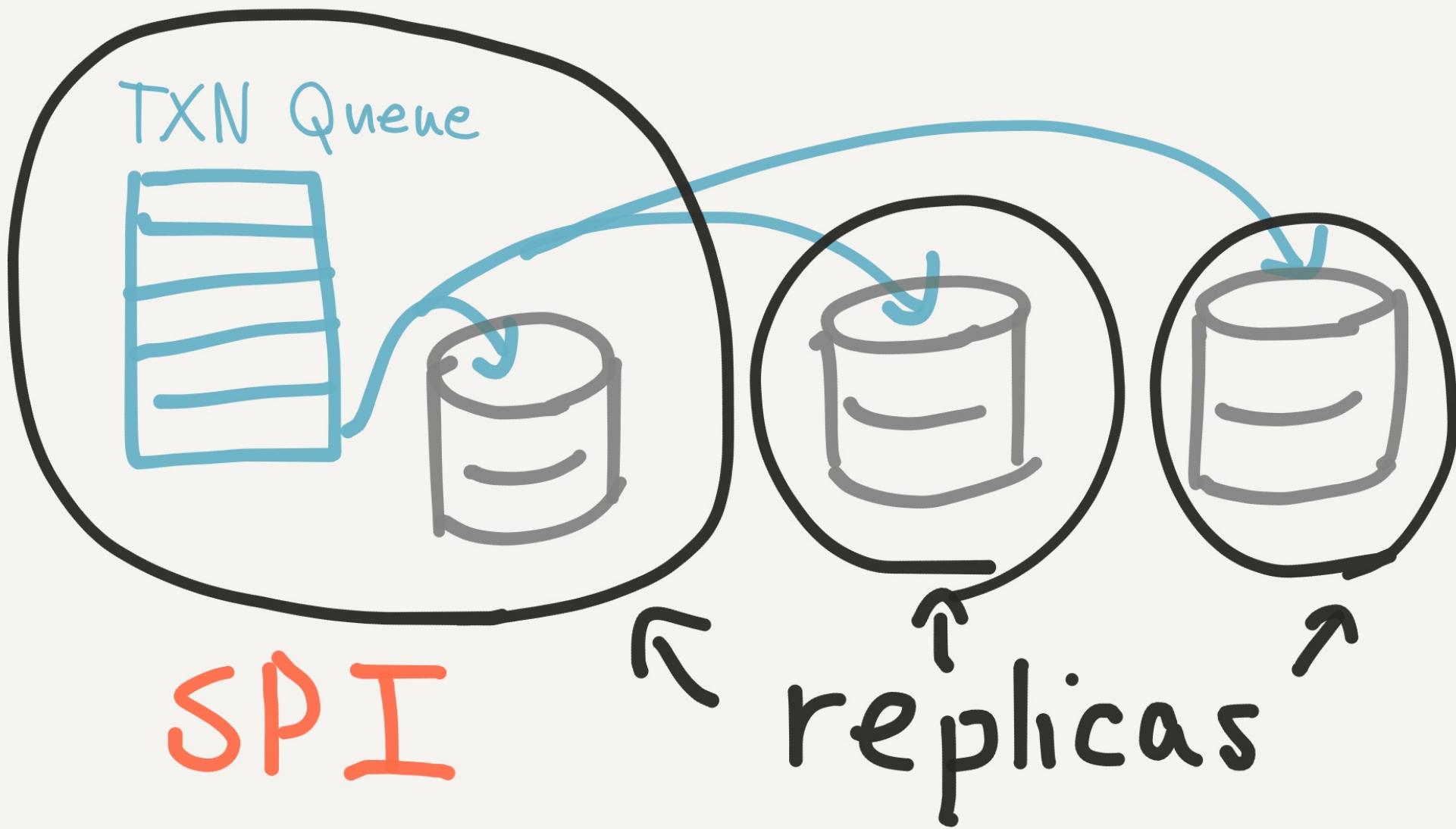
Txns as Java
stored procedures

Intended for
high-throughput,
mostly sharded
transactions

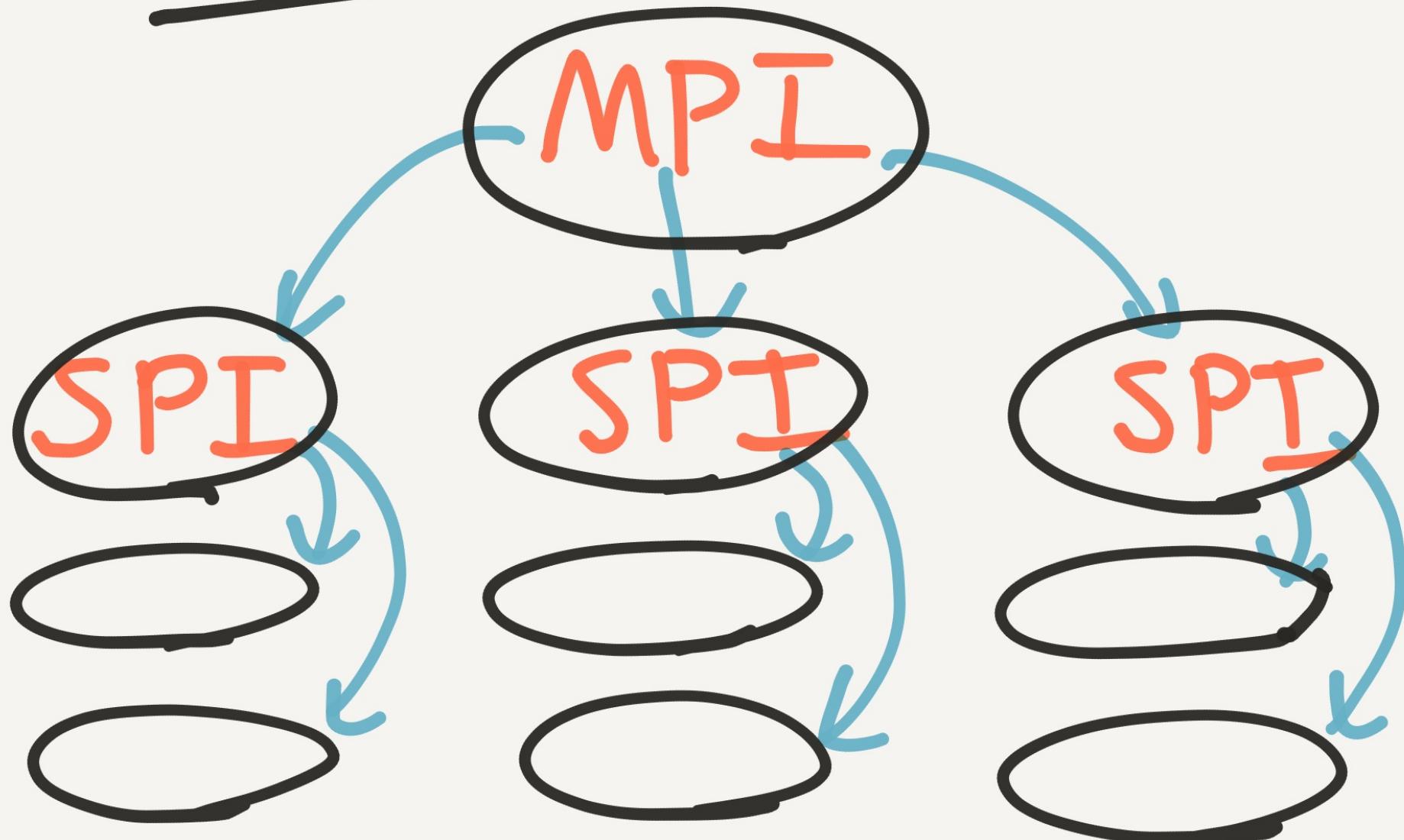
Keyspace



A partition



Across Partitions



- Txns in 1 partition

- scale ~linearly

- Txns across partitions

- have const throughput

Claim: all txns

are strict

Serializable!

Strict Serializable

Linearizable

Sequential

Causal

single-obj

Serializable

RR

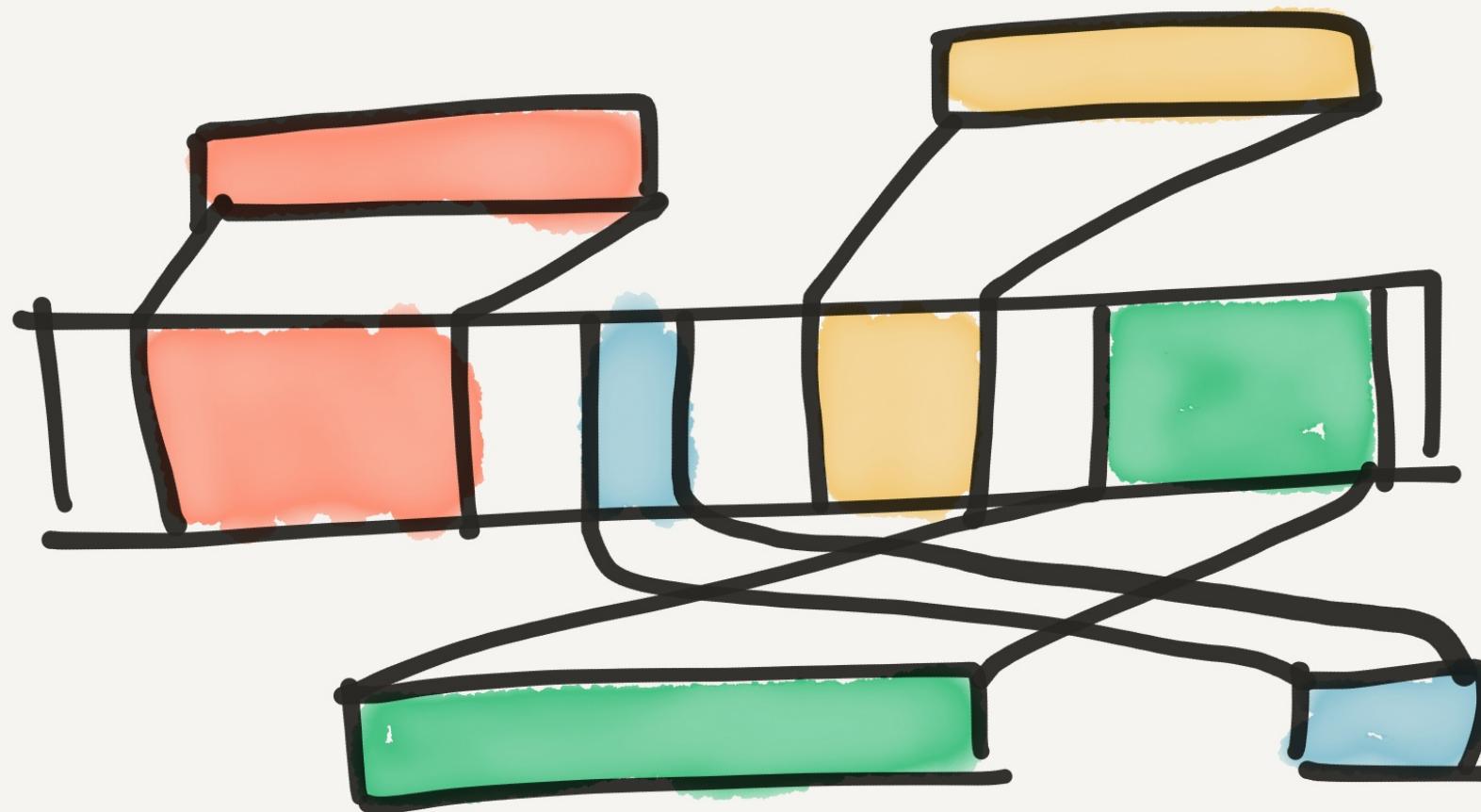
SI

RC

RU

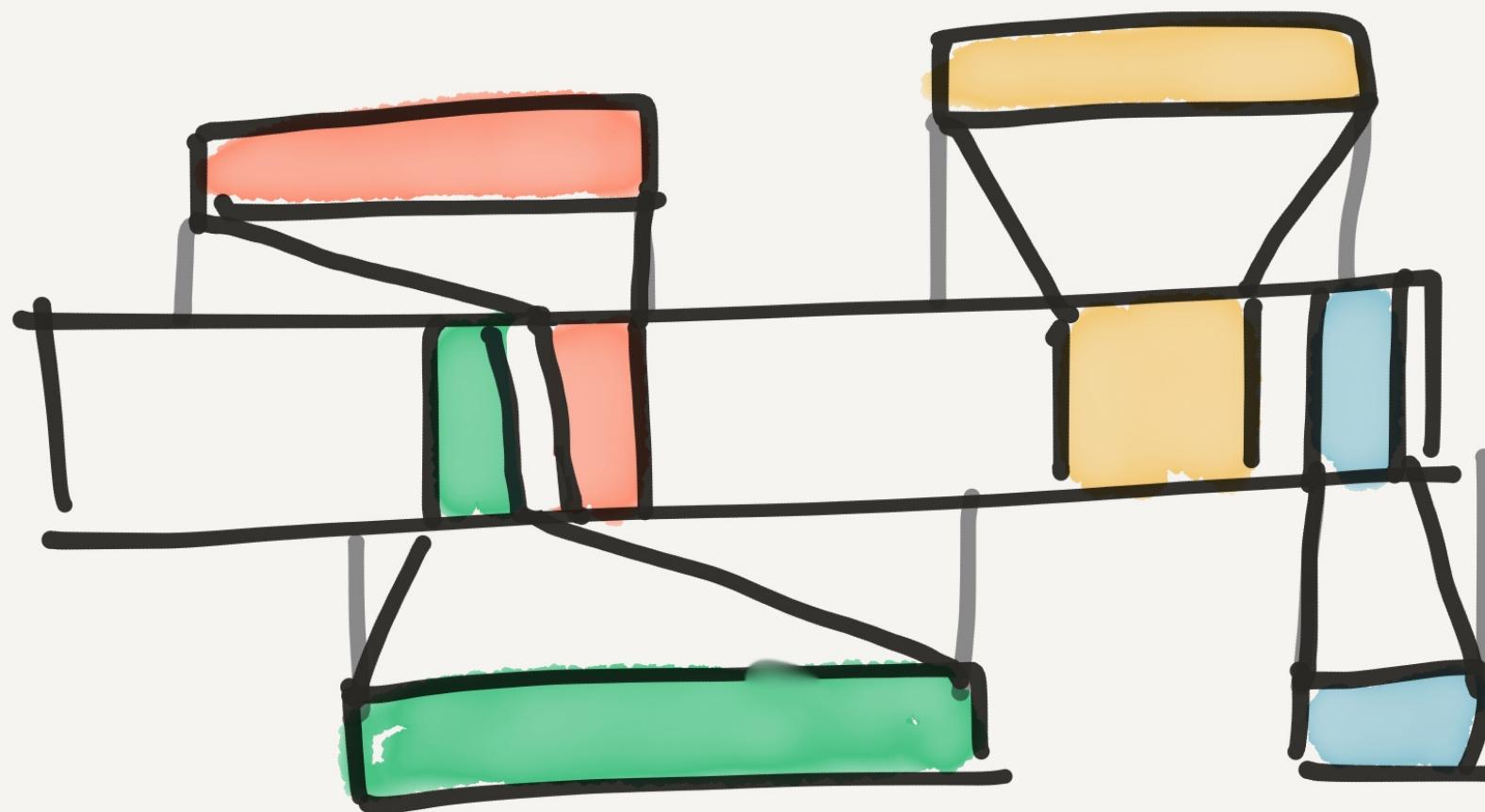
multi-obj

Serializability



time →

Linearizability



time →

Strict Serializable

- Real-time
- Multi-object

VoltDB 6.3

exhibits several
violations of
strict serializability

Stale Reads

Dirty Reads

Lost Updates

Stale Reads

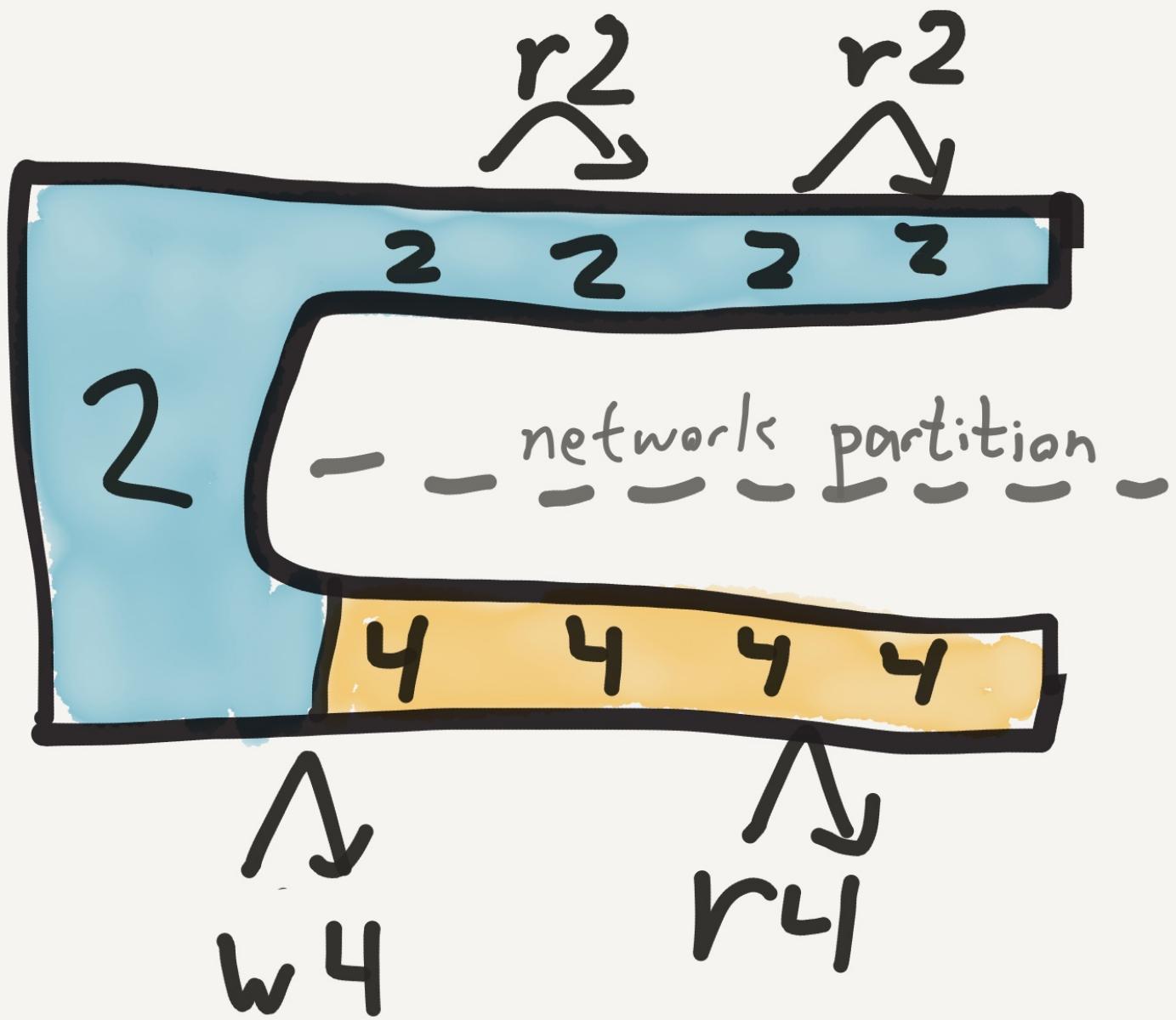
read 2

read 4

write 0

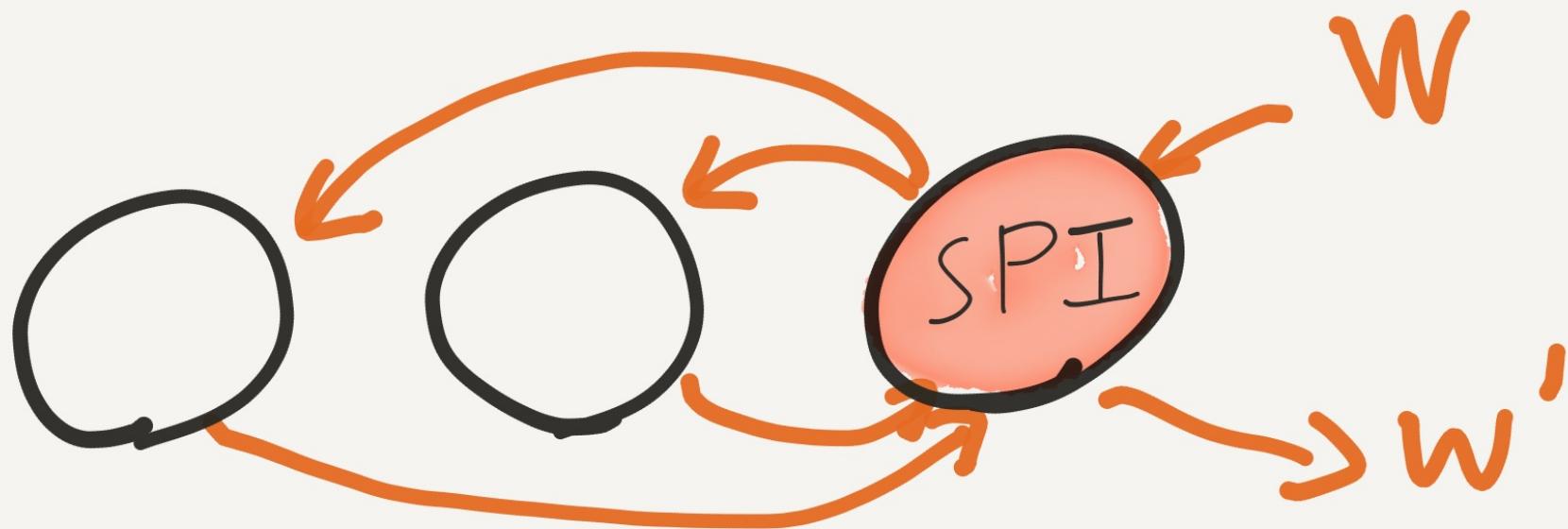
time →

node 1



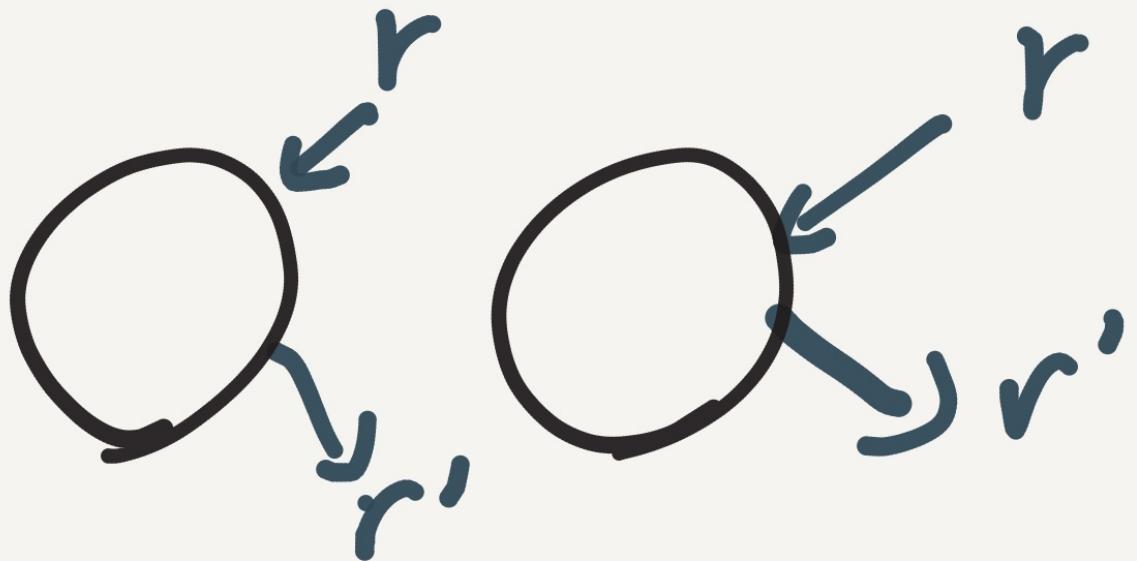
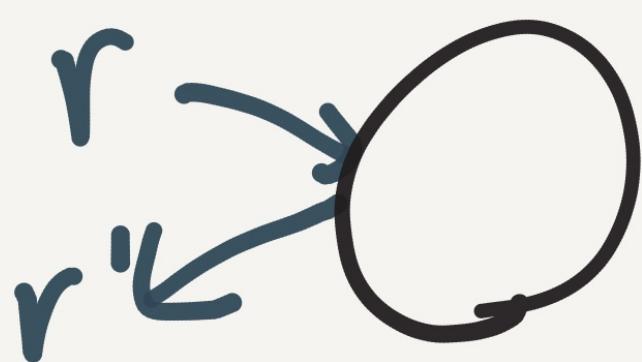
node 2

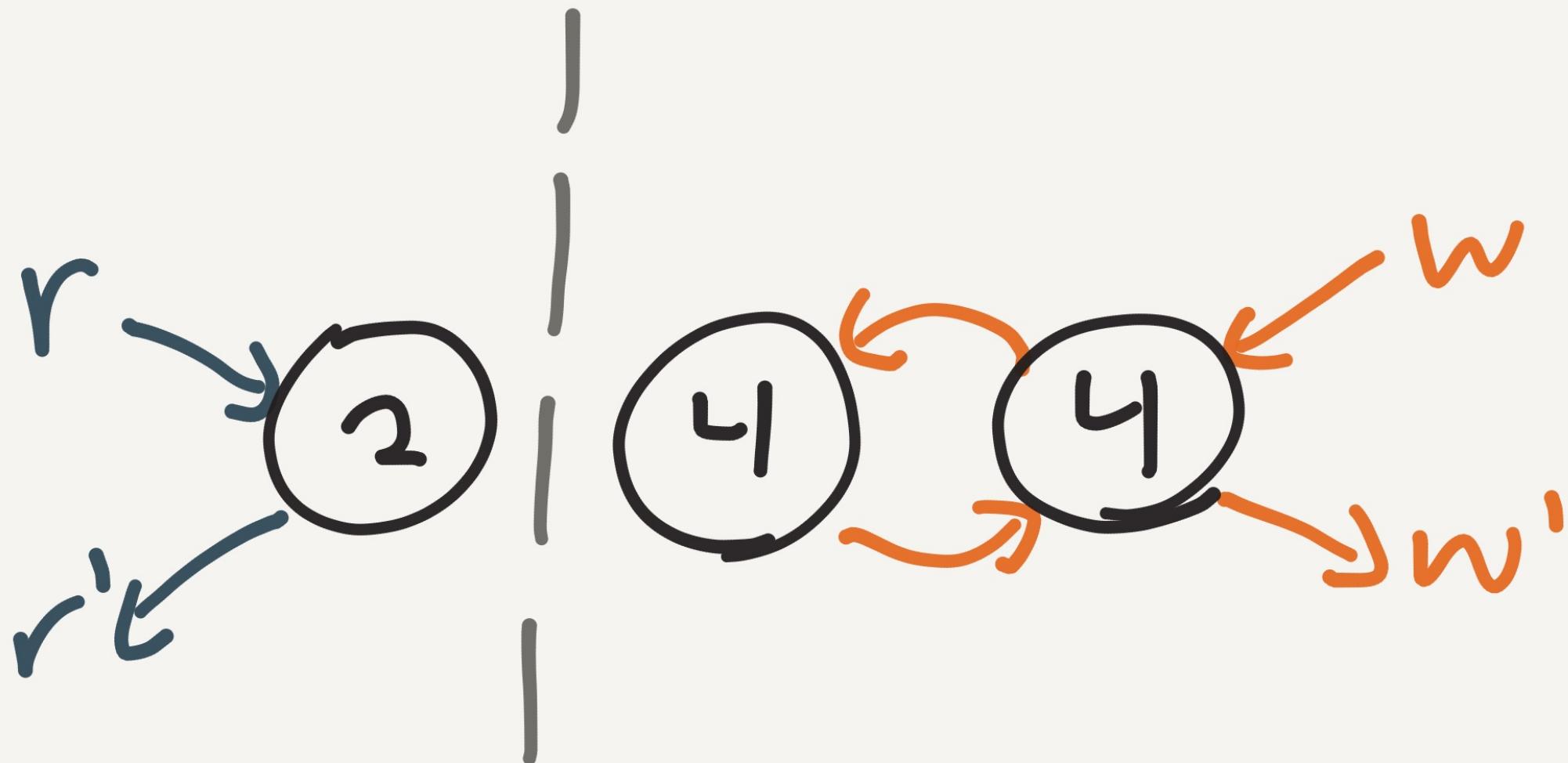
For a given partition



SPI orders updates

... but not reads





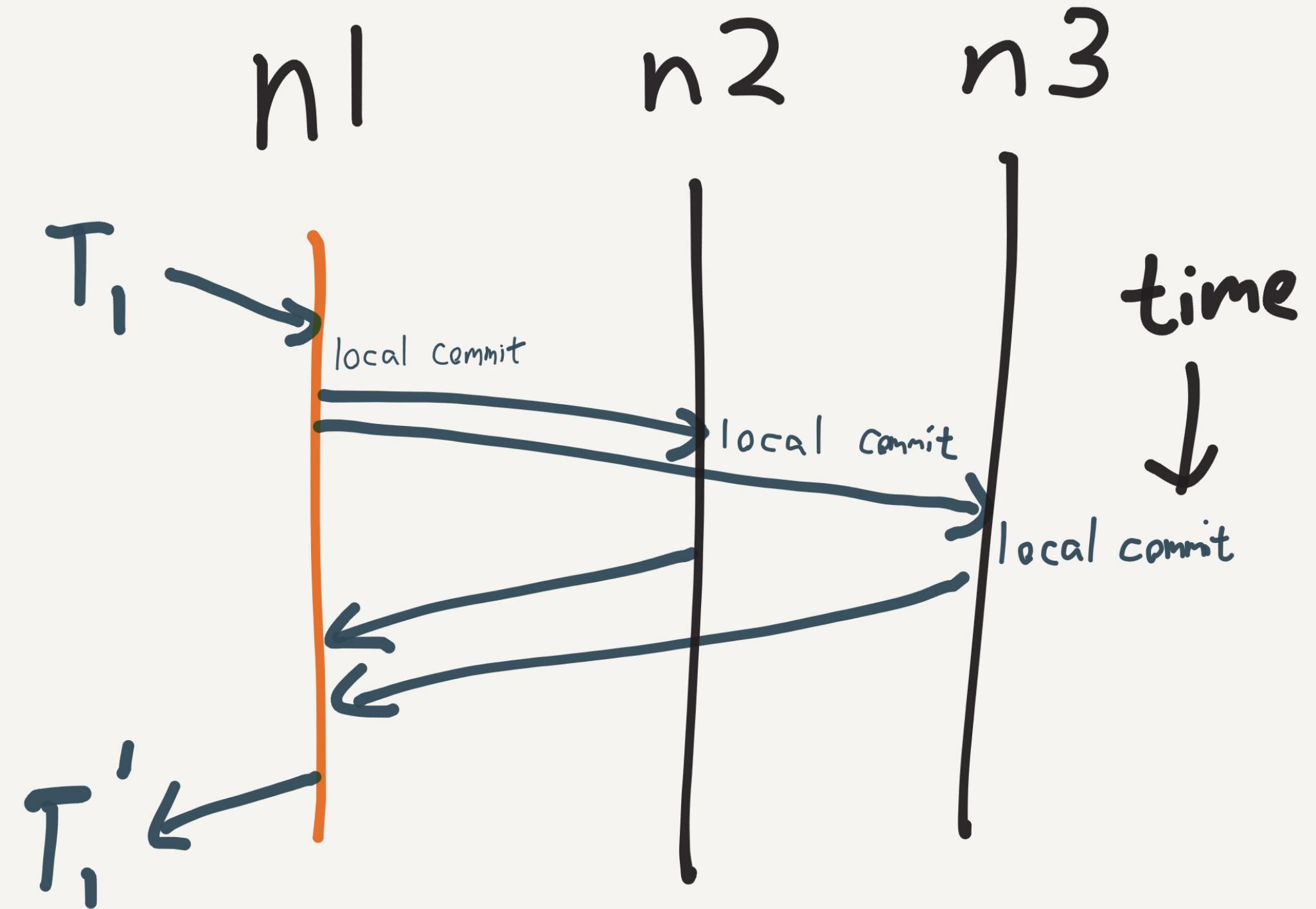
Solution:

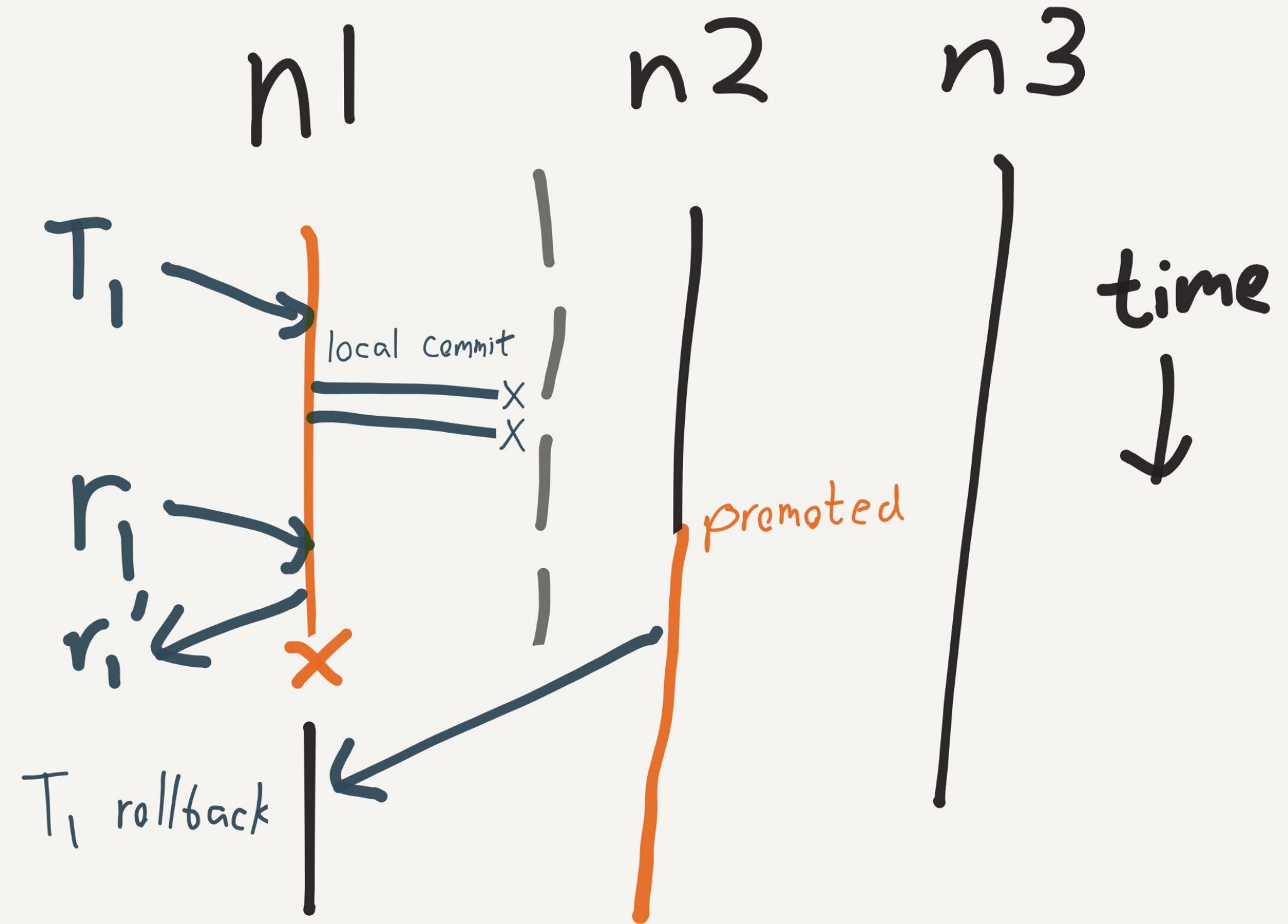
Push all txns
through the SPI
Ordering path

Dirty
Reads

T_1 w 3... abort

T_2 r3



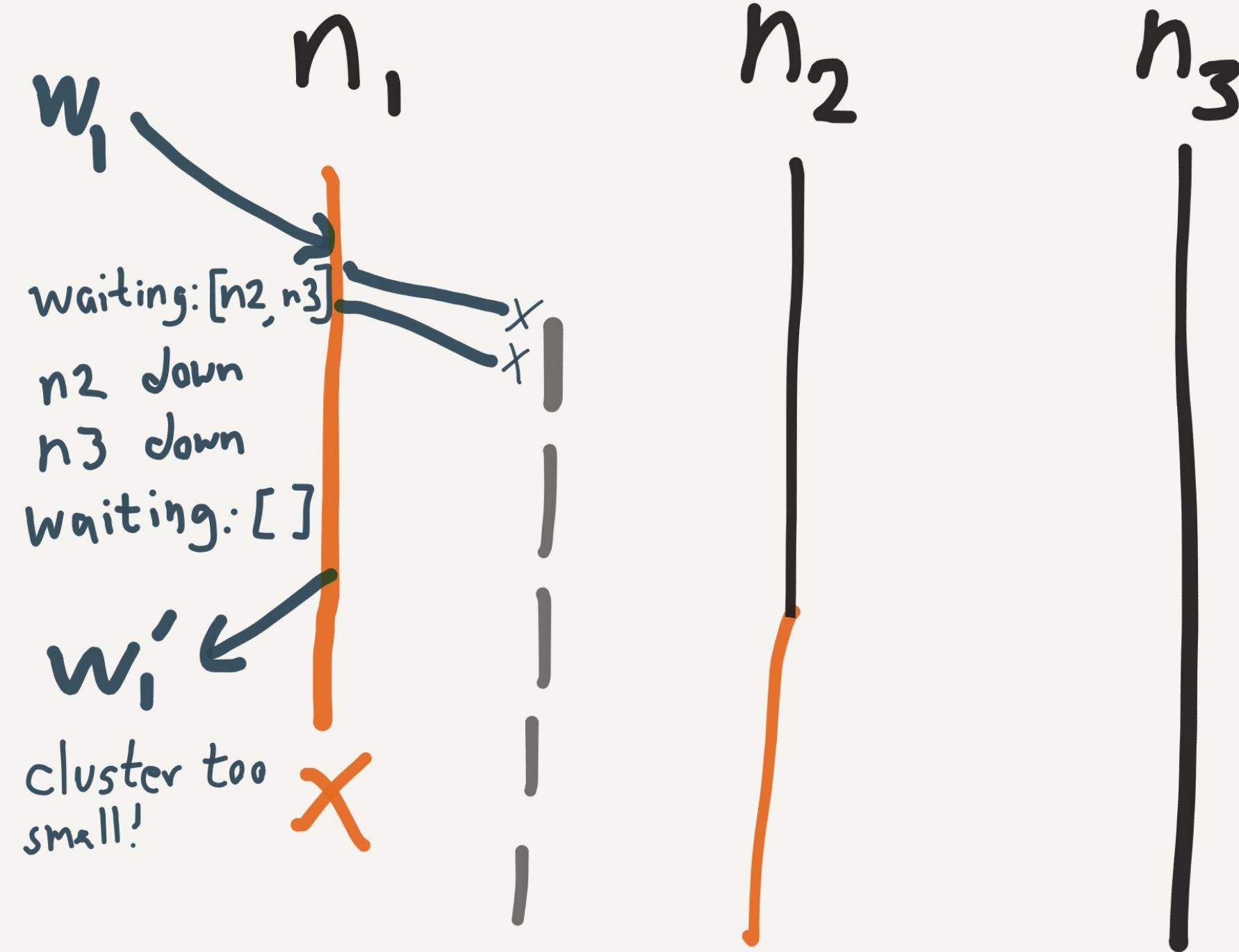


Lost

Updates

$T_1 : [w_1]$ ← lost!

$T_2 :$ [r_o]



"Zookeeper" watchos

are asynchronous;

race condition

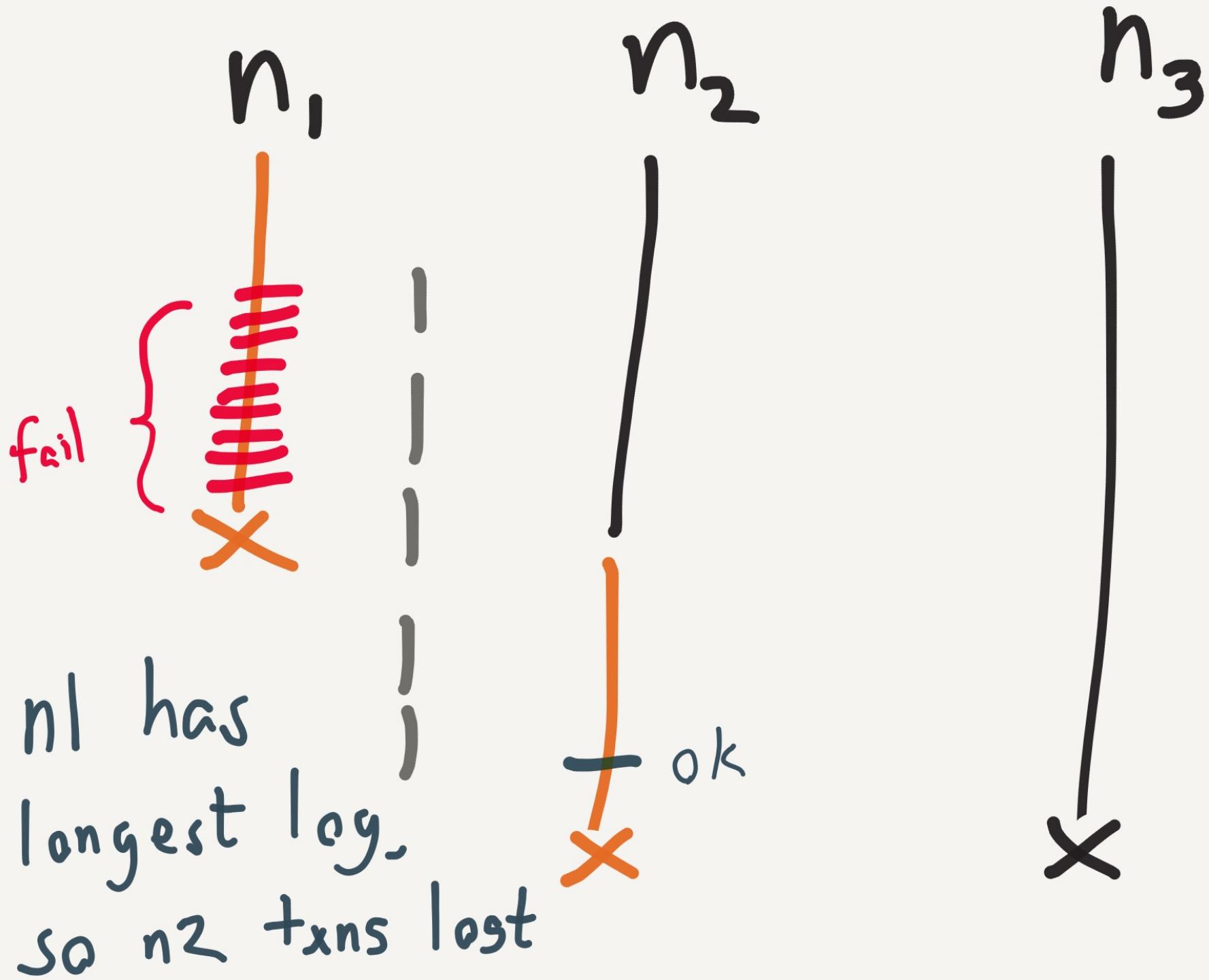
between ack & shutdown

Solution:

Make cluster shrink
↳ shutdown one
atomic step

Still loses
acknowledged
writes . . .

On Crash recovery,
Volt picks longest
log as authoritative.



Solution:

Reconstruct cluster
state from fault
logs

Multi-Partition

Transactions

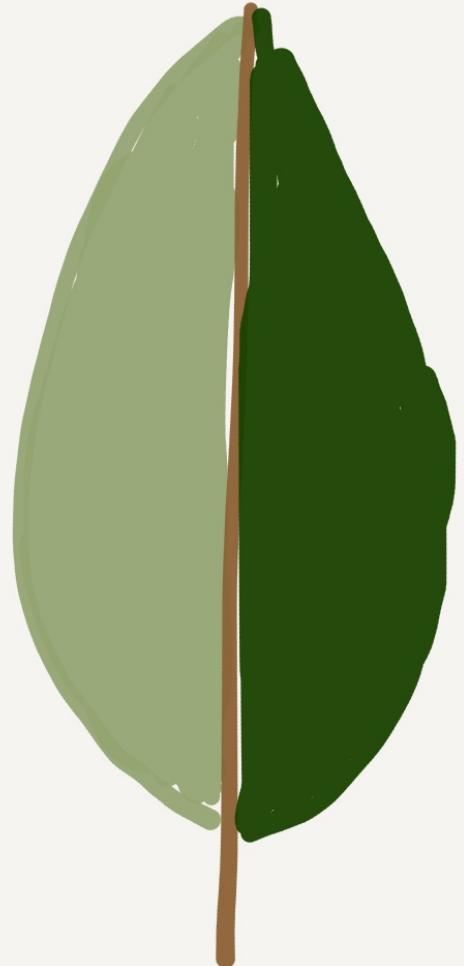
No faults found

MPI fully

Serializes read
and write txns,
may prevent divergence

VoltDB 6.4

Passes all Jepsen
tests for strict
serializability



mongoDB
3.4.0-rc3

- Document Store
- BSEN docs
- Custom replication
- Sharding layer

2013

Last writes

2015

Safe writes

Dirty/stale reads



3.0: Wired Tiger
Replication v1

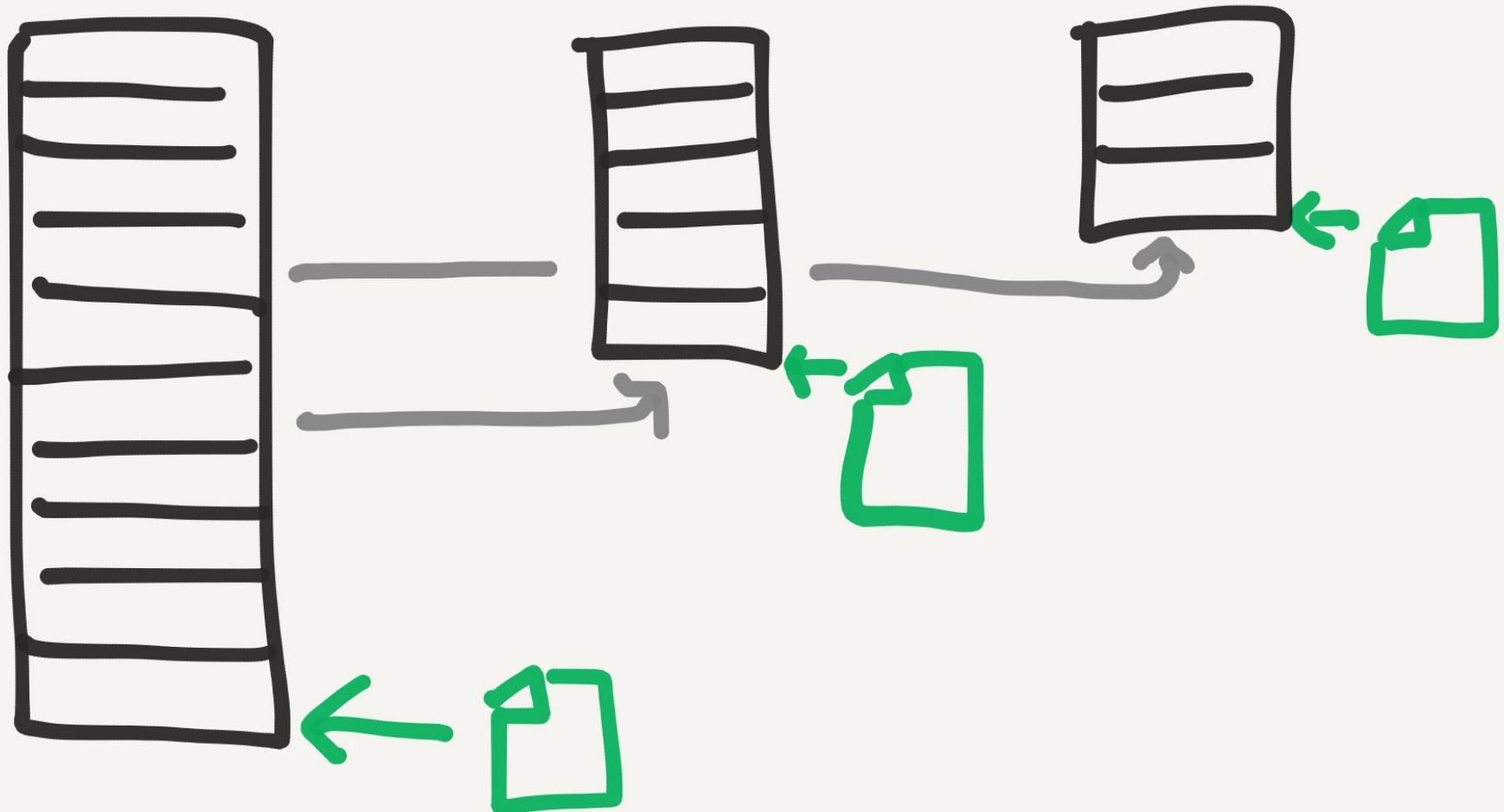
3.2: Majority reads

3.4: Linearizable reads

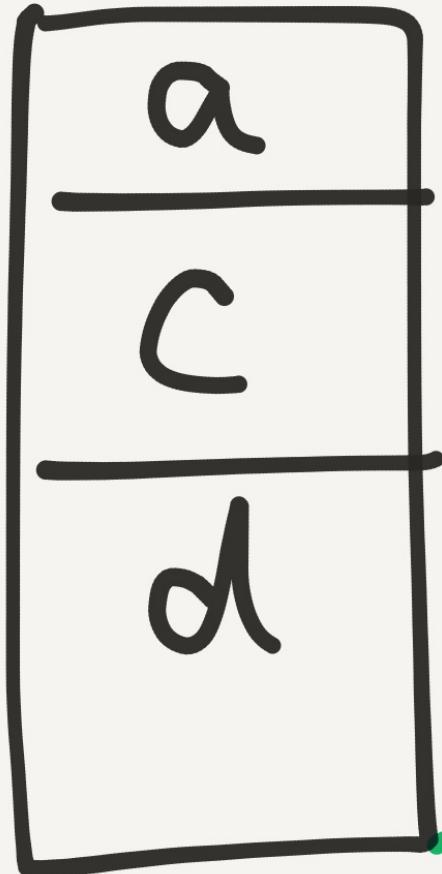
primary

sec.

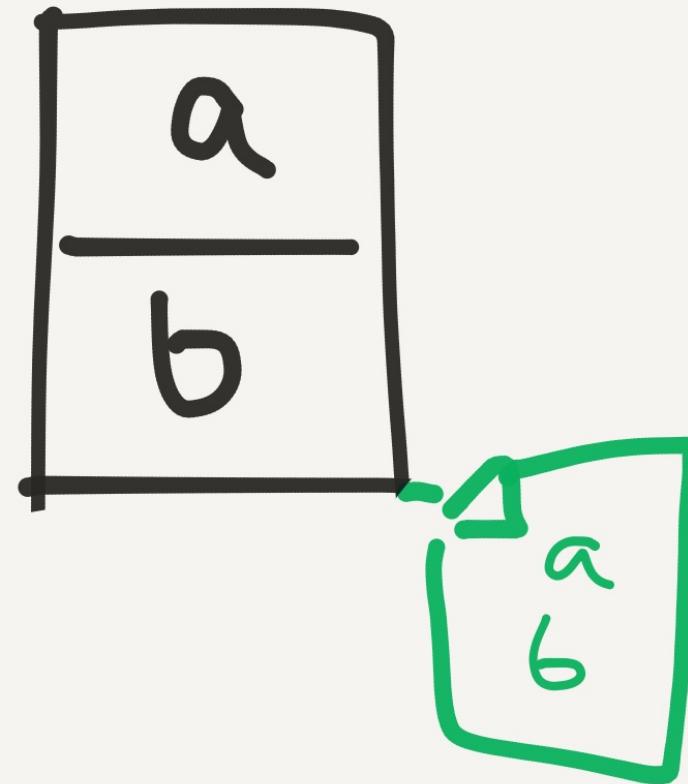
sec.



primary

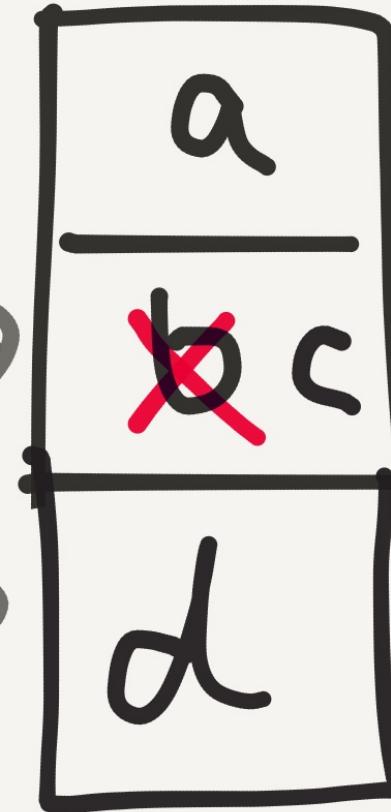
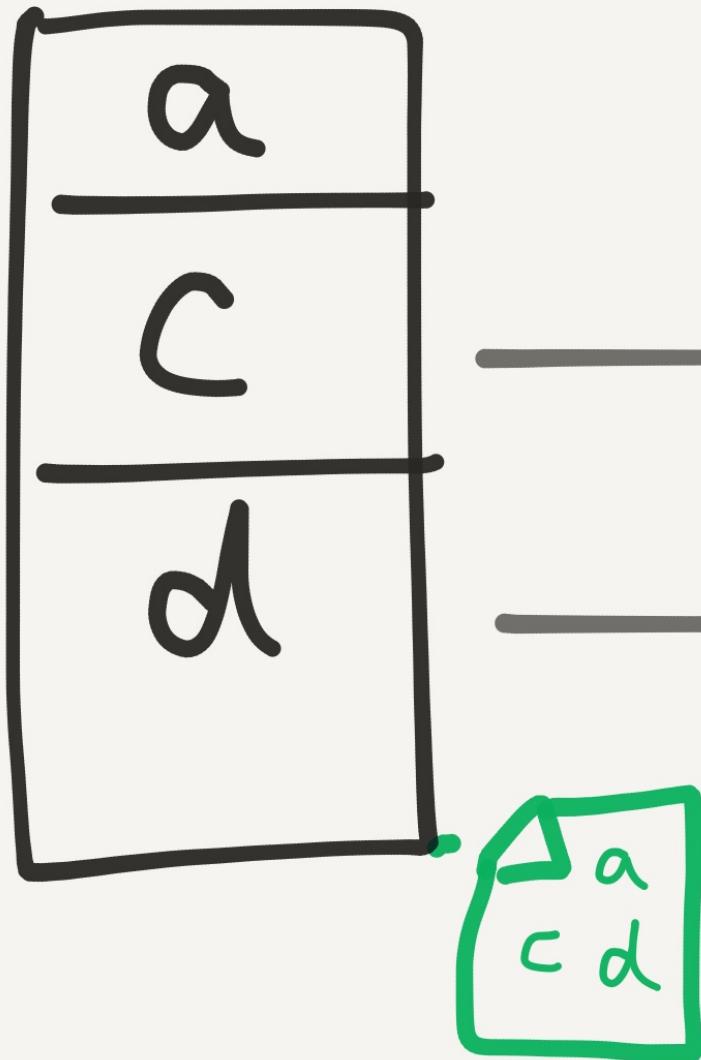


secondary

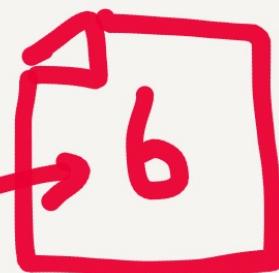


primary

secondary



rollback

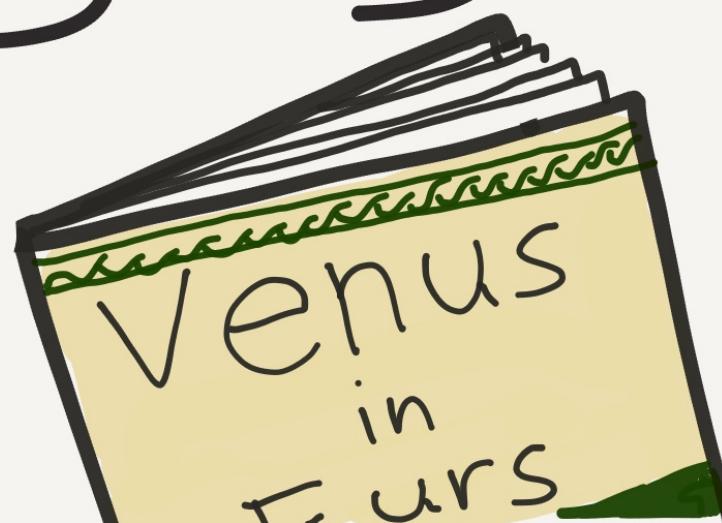


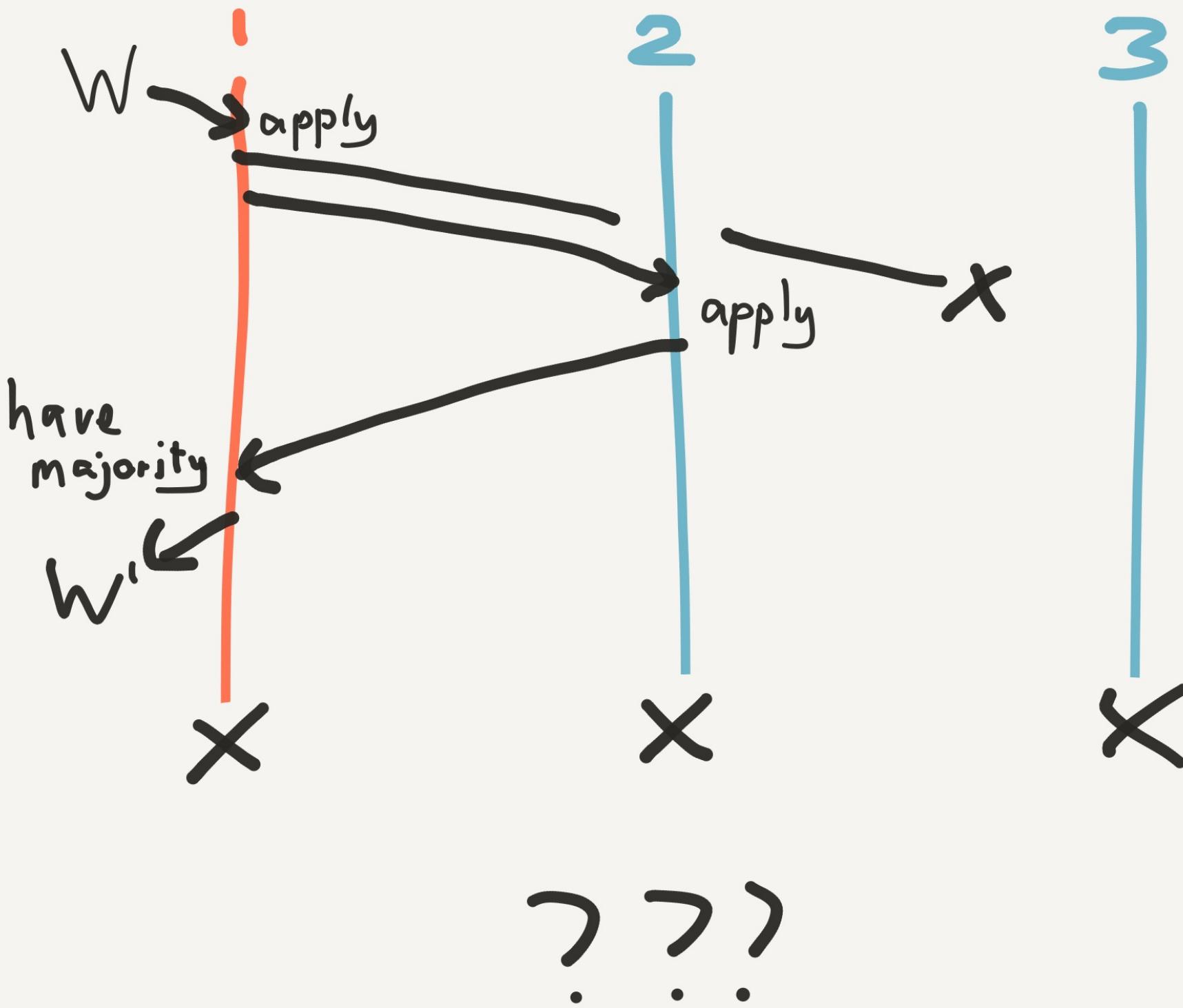
v0 protocol :

- Dirty reads
- Lost writes?

QK but how bad
are dirty reads

anyway?



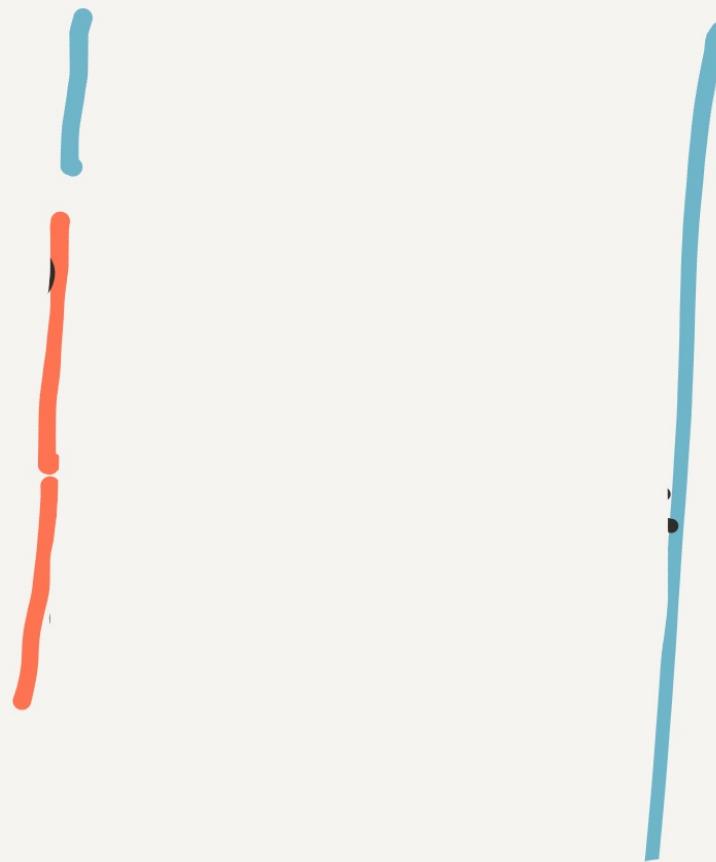
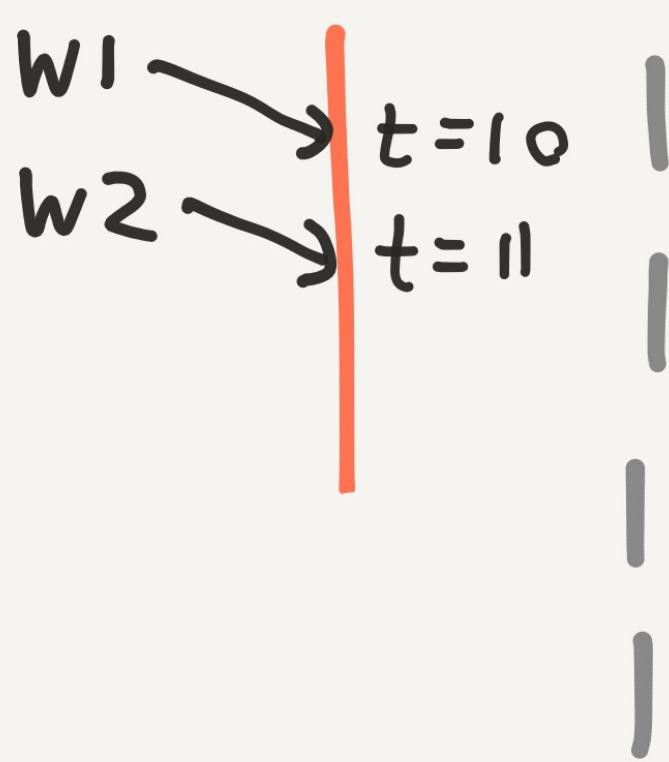


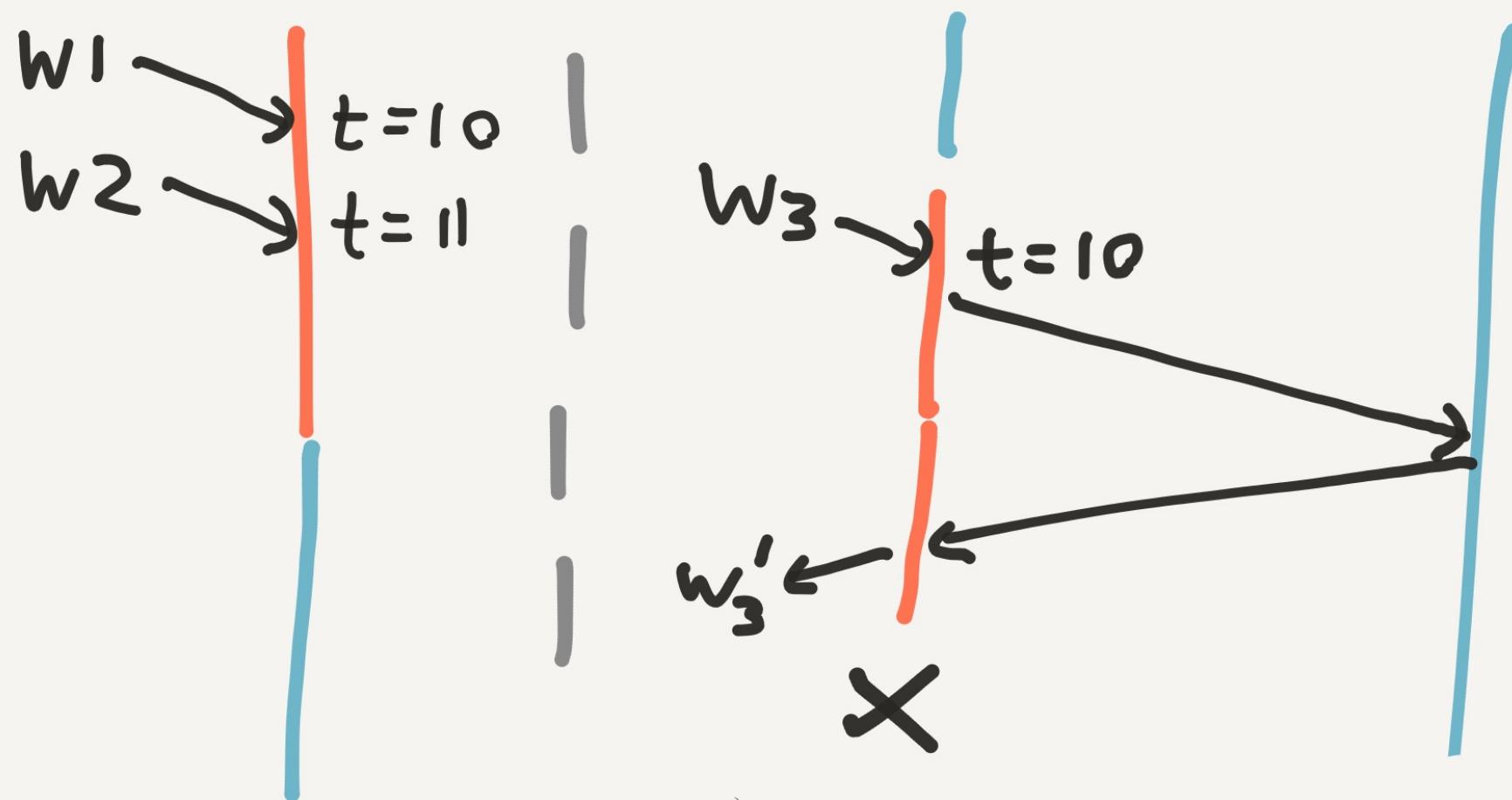
Optime:

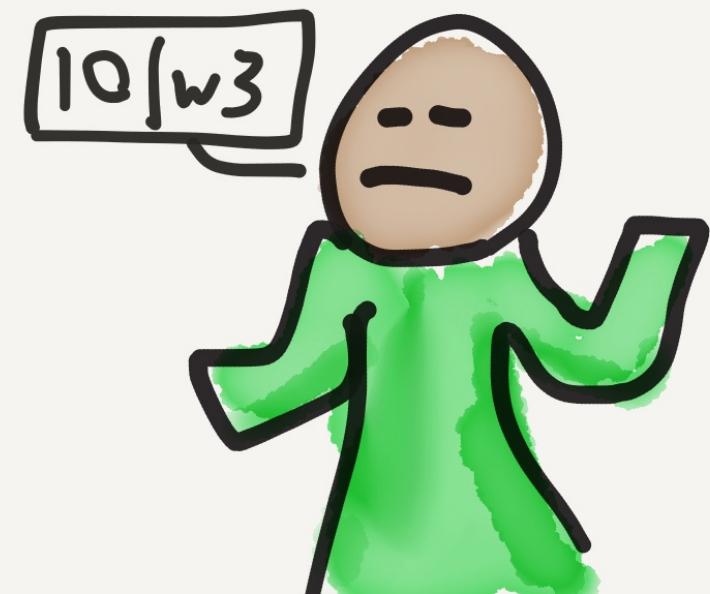
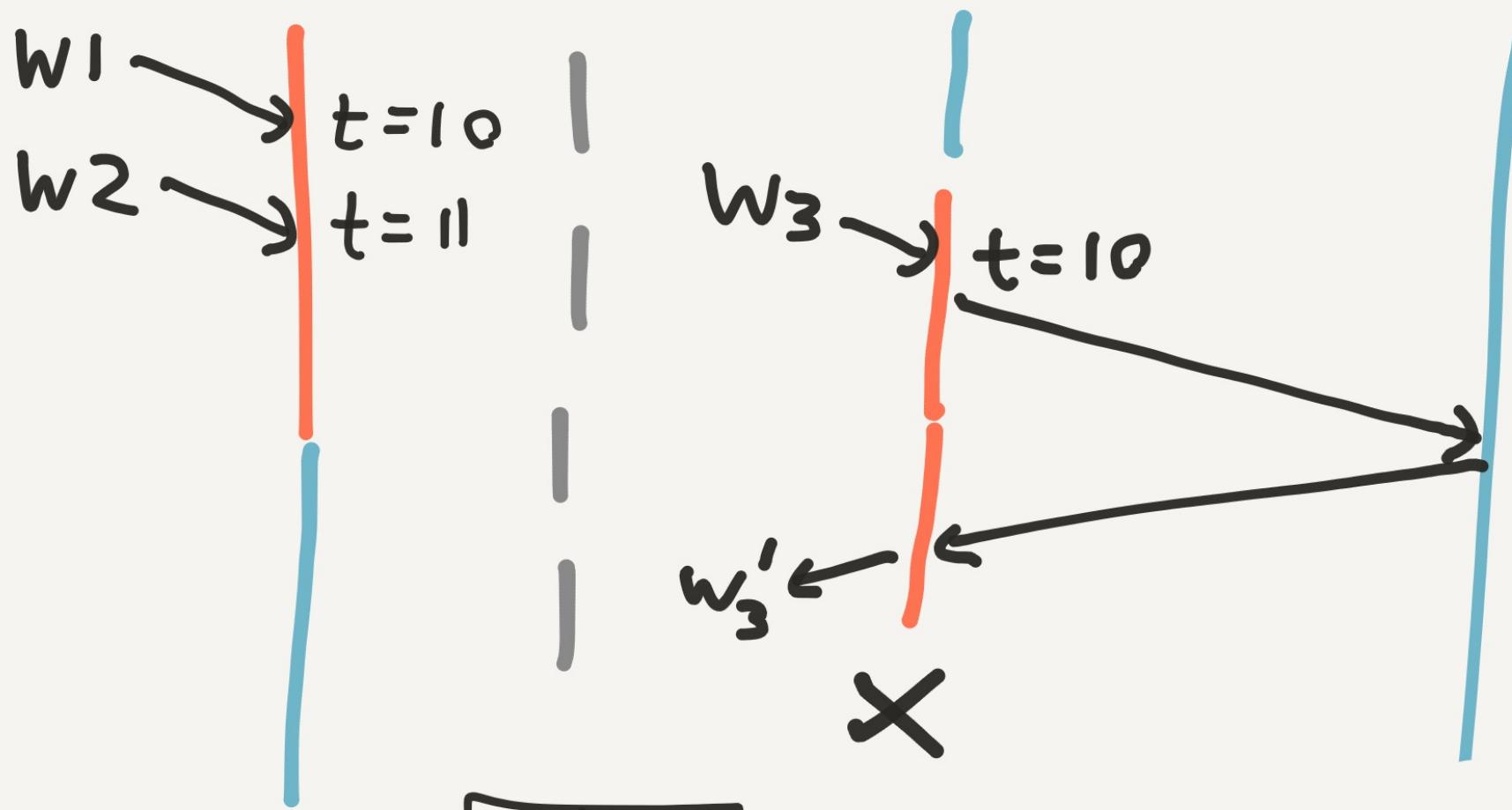
Max { last optime + 1
local wall clock

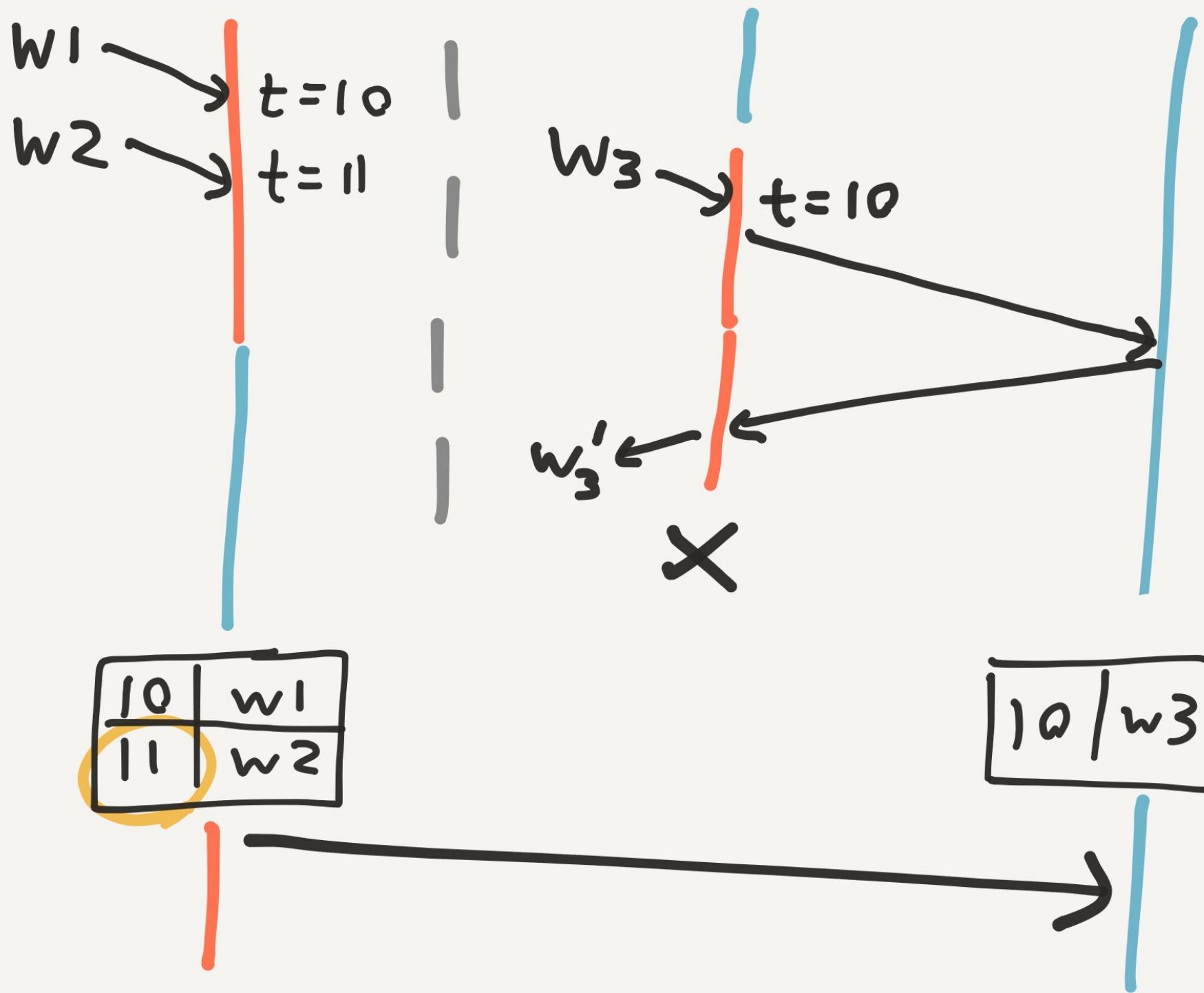
Optime	Operation
10	~~~
11	~~~
12	~~~
15	~~~
20	~~~

Nade
w/ highest
optime
wins









Requires precise timing

- But window grows
with clock skew!

$$\begin{array}{r} \text{ok} \\ 322 + 93 \\ \hline \text{lost} \\ 4525 \\ \text{attempted inserts} \end{array} \quad \left. \begin{array}{l} \text{lost} \\ \} \text{acked} \end{array} \right\}$$

vo protocol

is fundamentally

broken 😞

Protocol v1

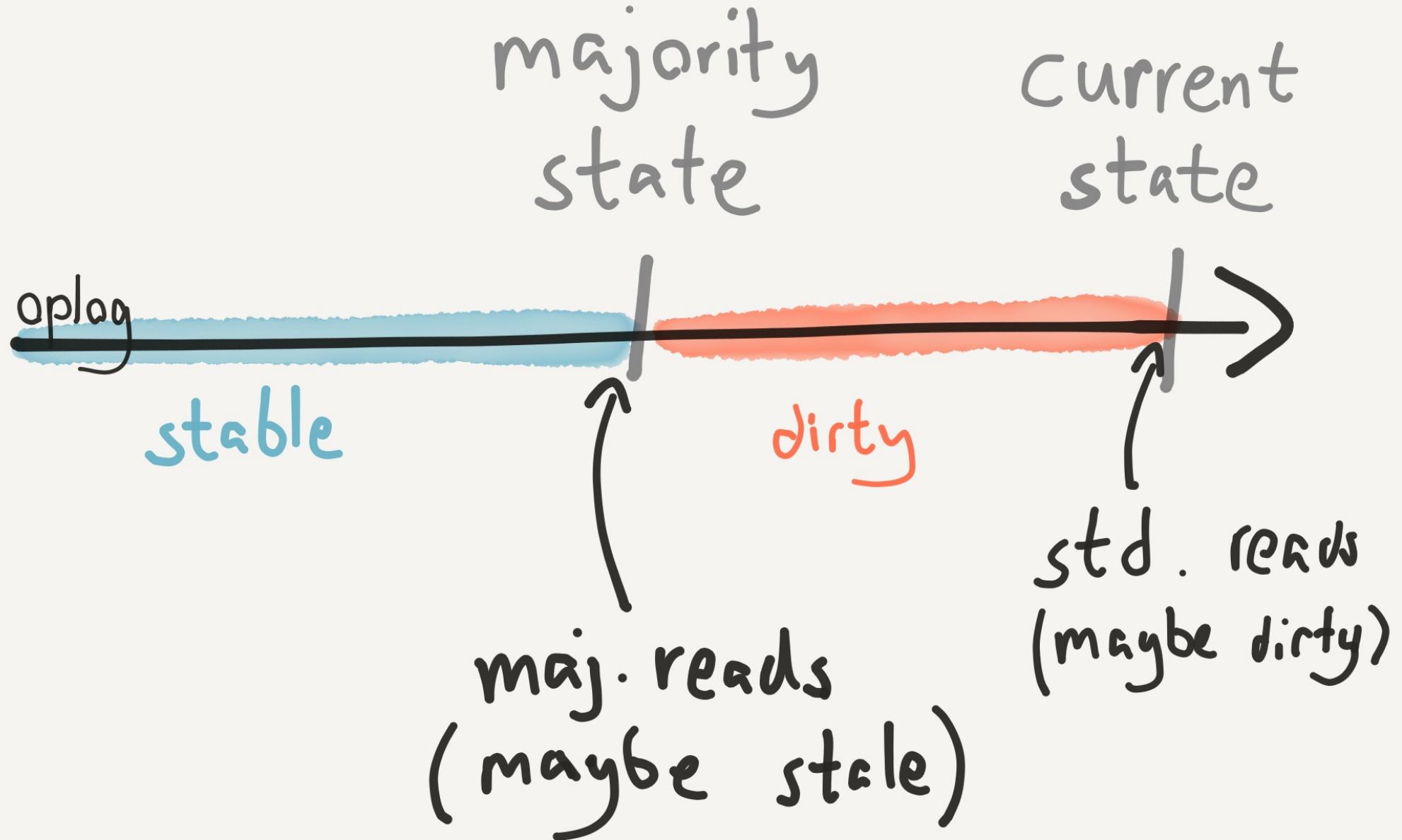
Default in 3.2
for new replsets

Optime

[term, time]

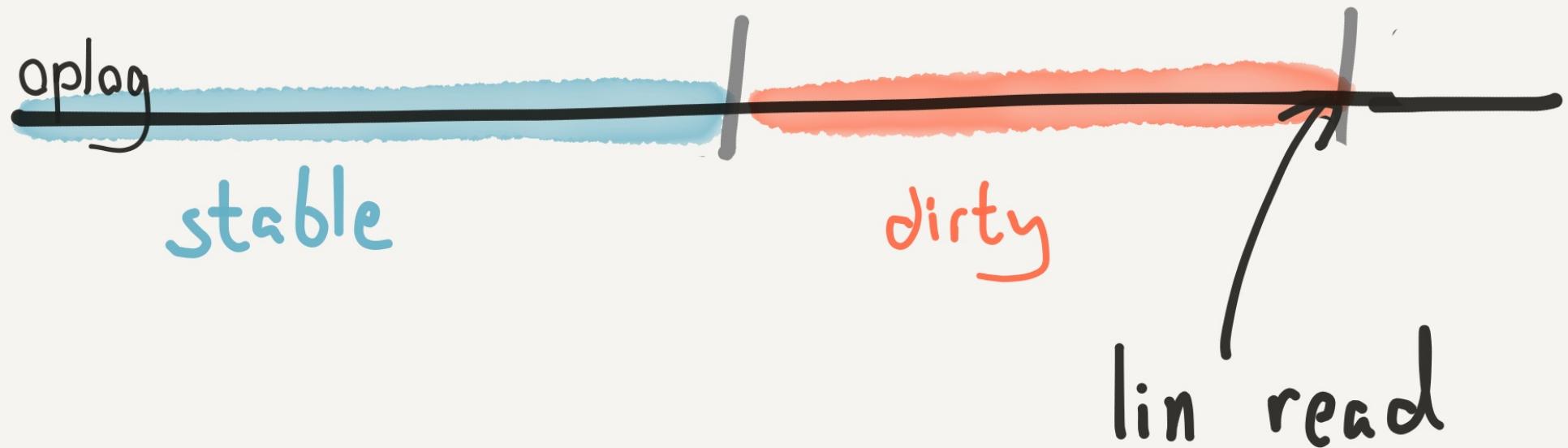
logical clock

old optime



majority
state

Current
state



majority
state



OK

lost

5|5 + 4|7

} acked

6078

attempted inserts

Server - 27053

primaries could ack
writes from prior
terms

| -Heartbeats received

by outdated primary

could advance Commit
index

3.4.0 - rc4

Server - 27149

Secondaries ignored
terms þ took ops
from old leaders

THE KINGDOM

Fixed in

3.2.12

3.4.0

3.5.1

Recommendations

- Get off v0
- Don't use arbiters
- Upgrade to 3.4.*

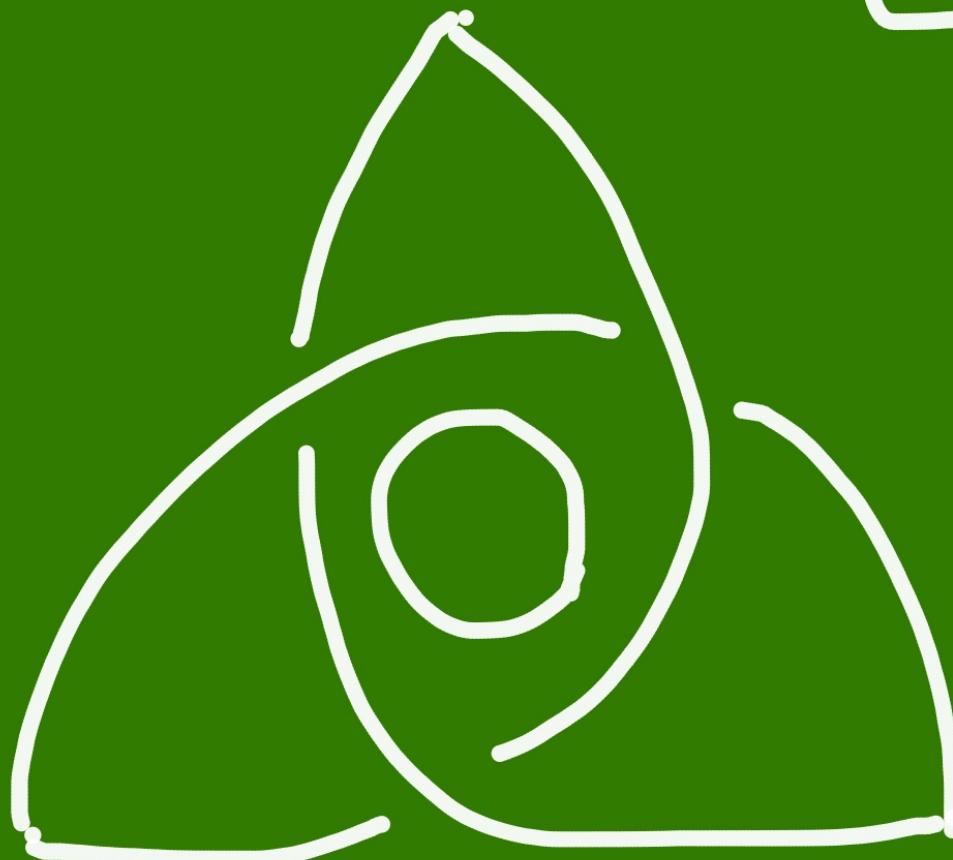
Use Write Concern

MAJORITY

Read Concern

MAJ/LIN

0.10.2



Tendermint

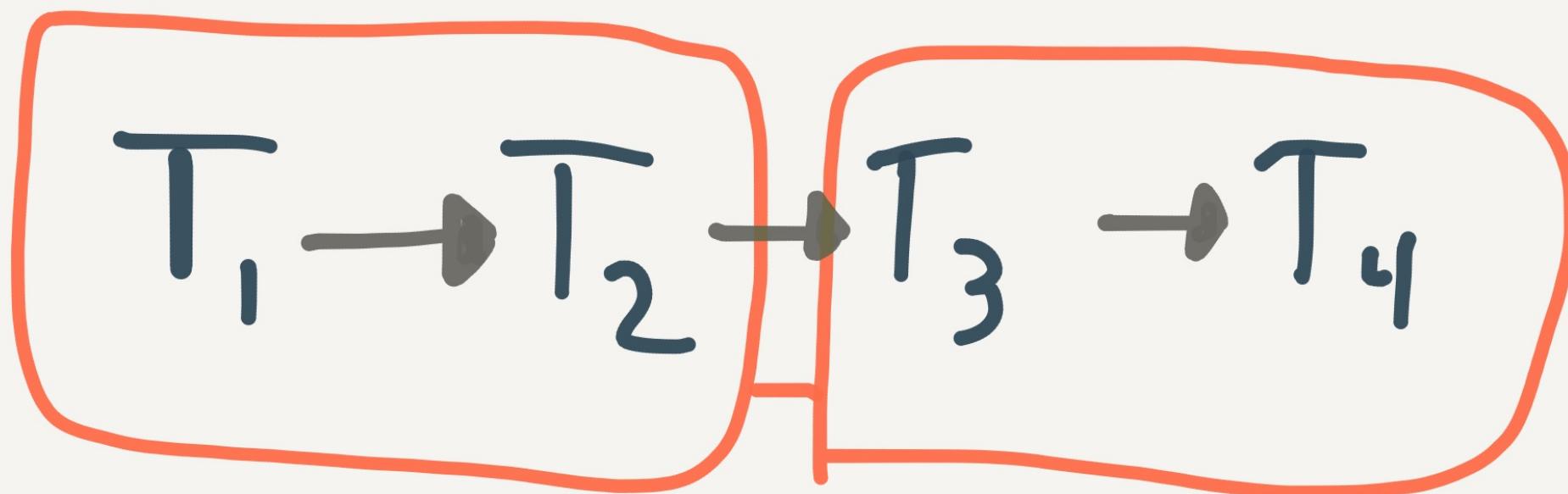
- Transaction replication
- Pluggable FSMs
- Byzantine fault tolerant

"Raft Meets
Blockchain"

Like Raft, provides a linearizable order of txns

$$T_1 \rightarrow T_2 \rightarrow T_3 \rightarrow T_4$$

Like Raft, provides a linearizable order of txns



blockchain!

Byzantine

Fault
Tolerant!

shun!



shun!



shun!



shun!



Pluggable
State
Machines

validator

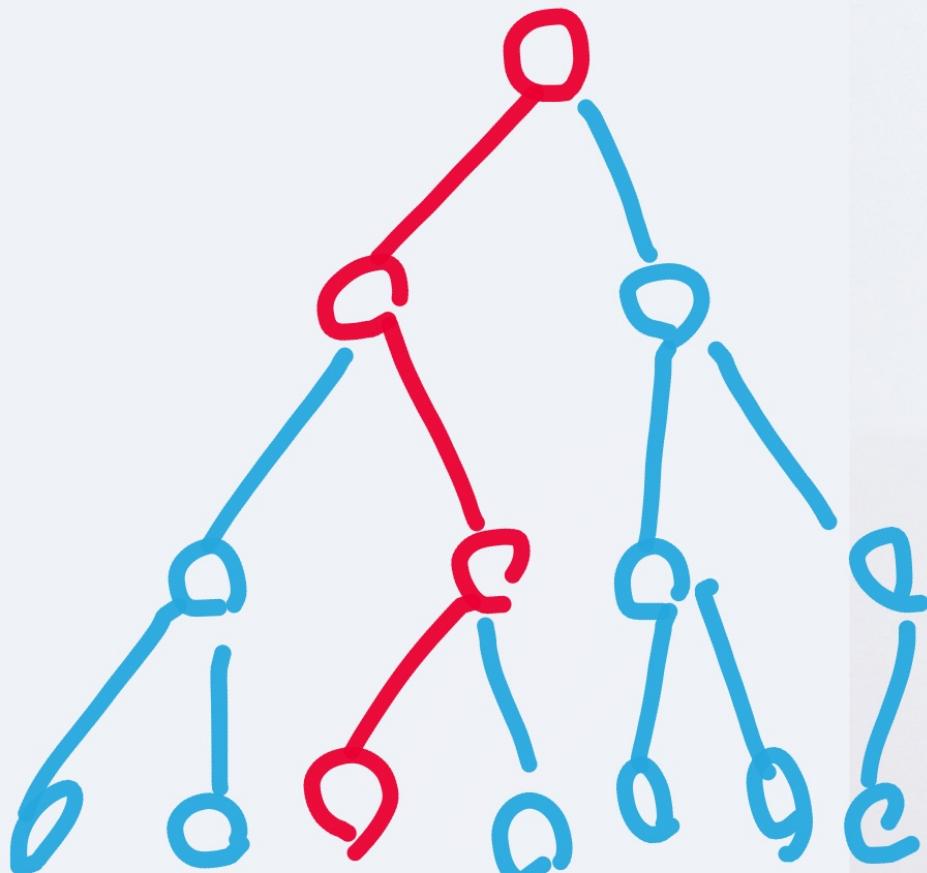
validator

|

State

|

State



merkle
eyes



CAS-register

- read (2)
- write (5)
- cas (1, 4)

set

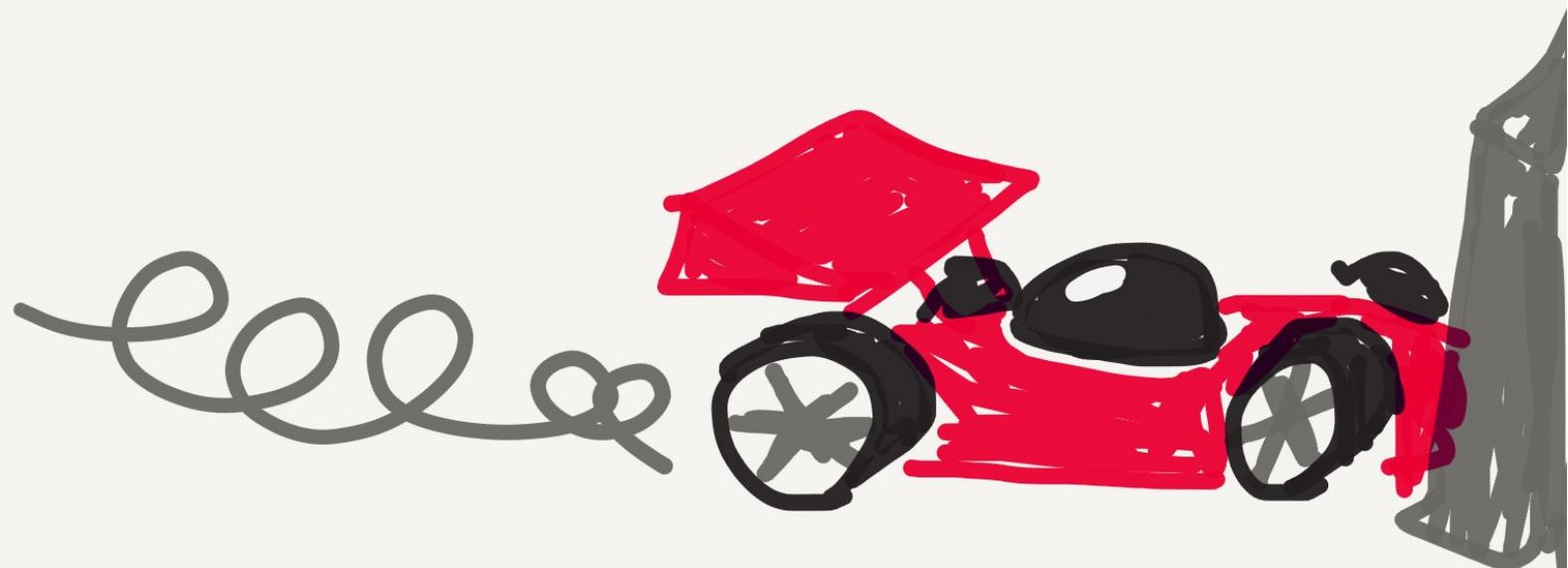
- add(6)

read({1,2,3})

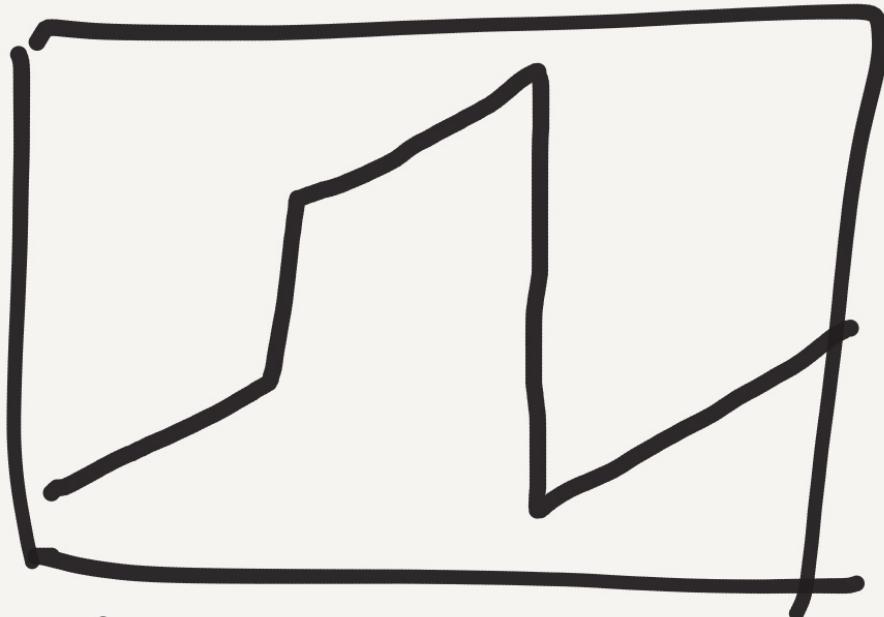
[cas({1,2,3}, {1,2,3,6})

- read({1,2,3, ...})

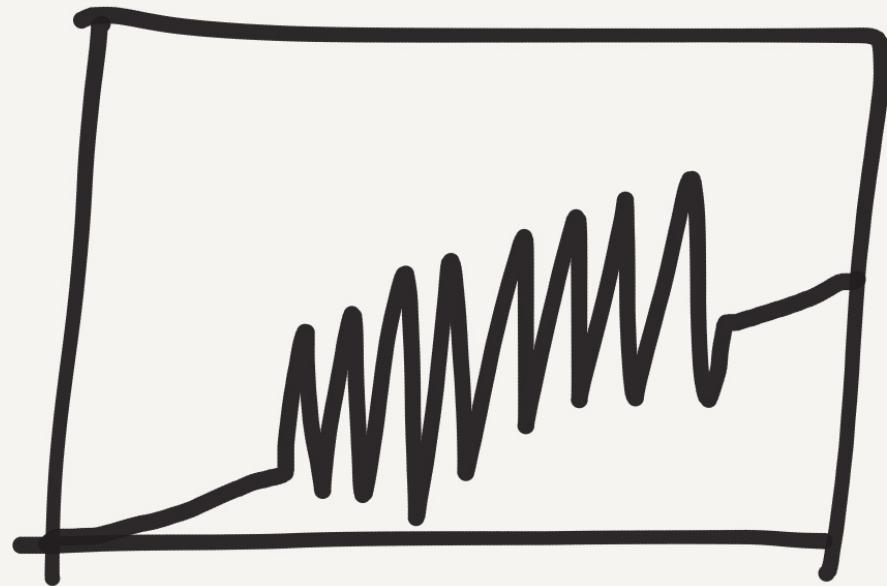
Merkleeyes Race



clock skew



bump



strobe

Crashes

$n1$ —

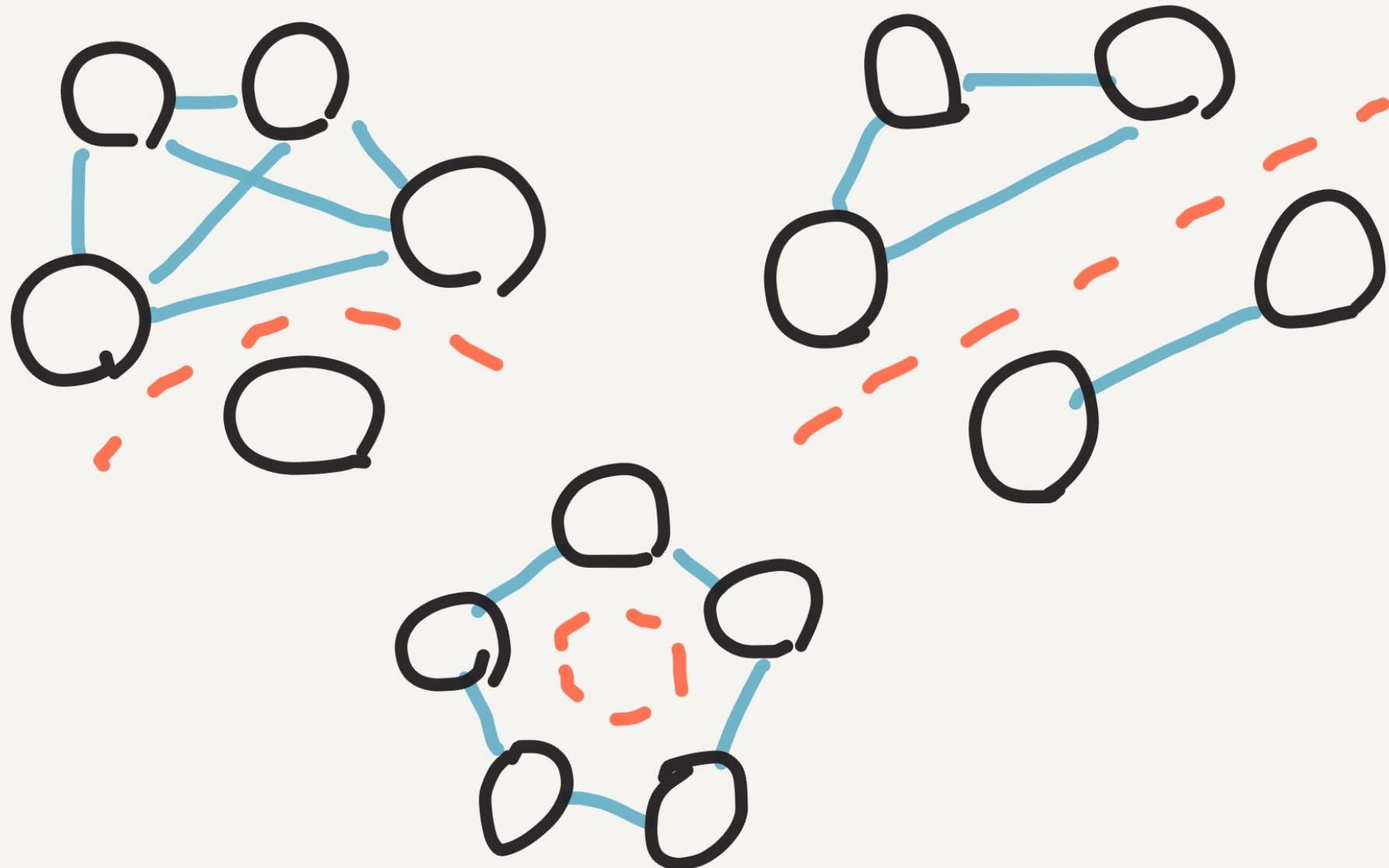
$n2$ —

$n3$ —

$n4$ —

$n5$ —

Partitions



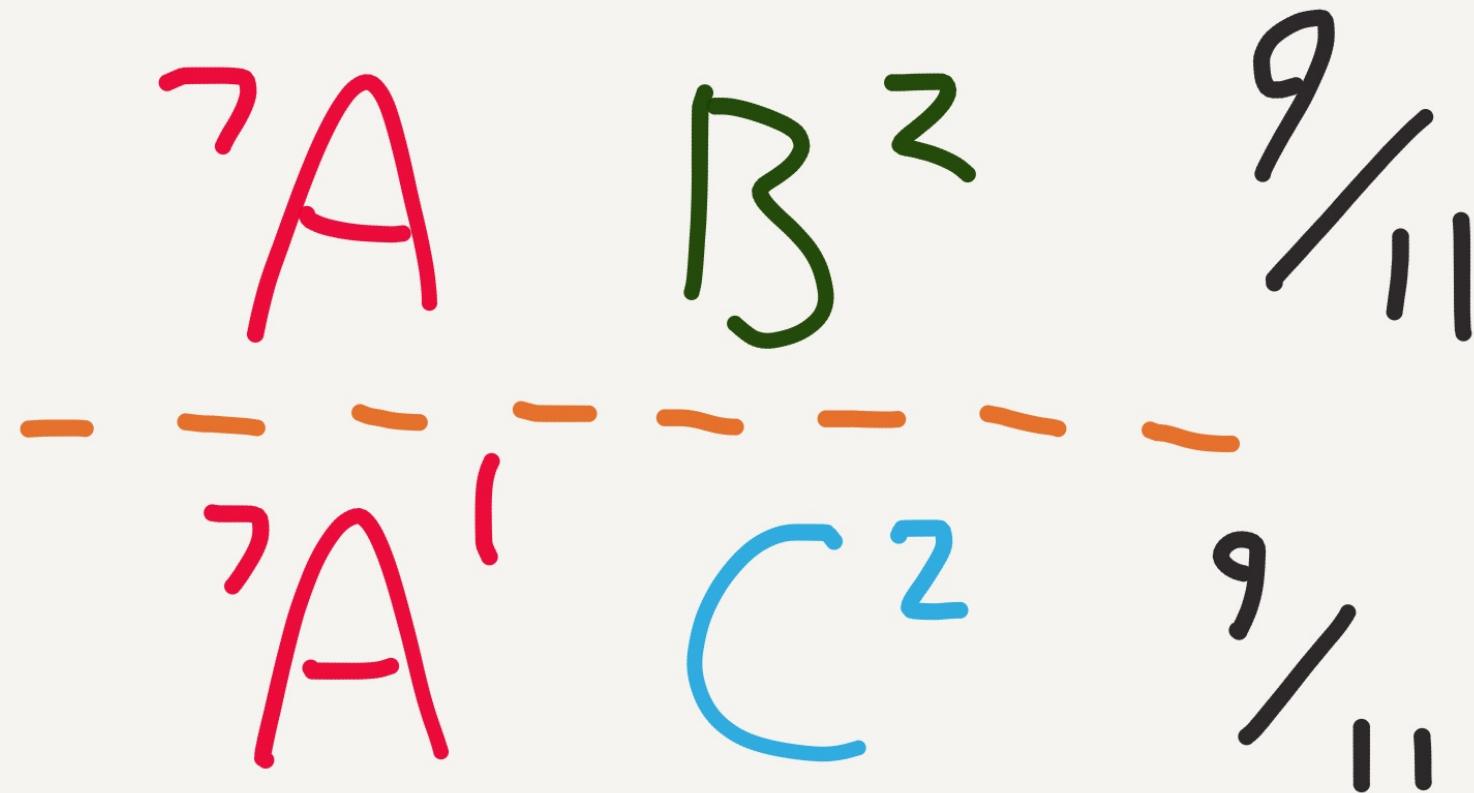
Duplicate Validators

?A B²

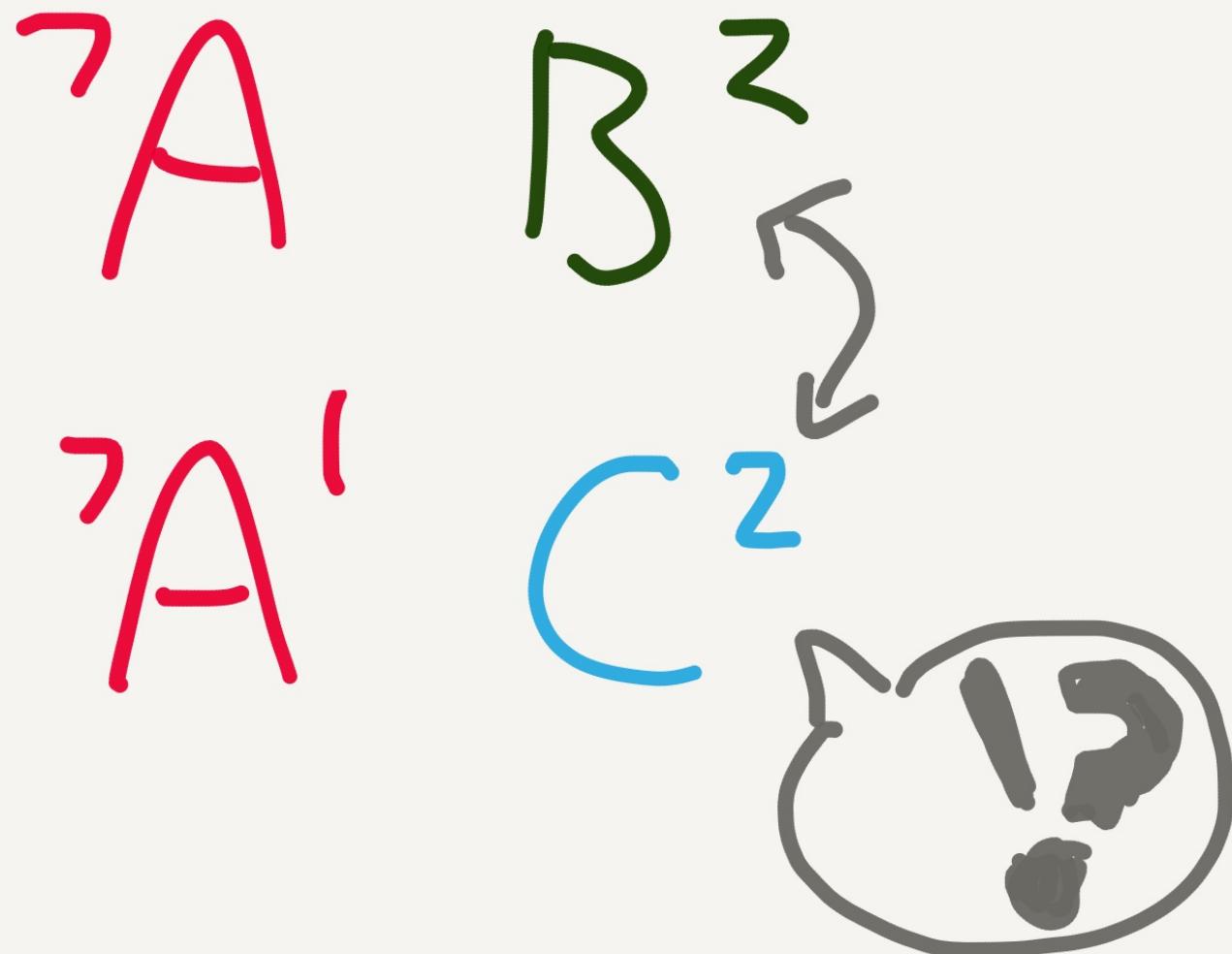
?A' C²

$$7 + 2 + 2 = 11$$

Duplicate Validators



Duplicate Validators



lost

Documents



iff



> 1/3

of votes

file truncation



Merkleeyes crash!

(goleveldb bug)

Tendermint

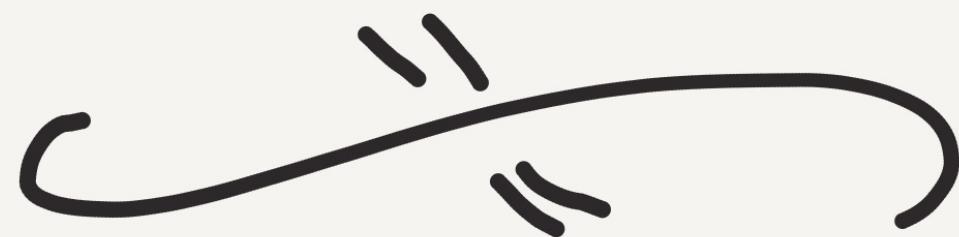
Doesn't fsync to
disk!

Dynamic Reconfig



Data Corruption &
Crashes are to be
expected — Tendermint
is still in beta

Team working to
fix these issues



Recap



Read the docs!



CAREFULLY

Then test it

— for —

YOURSELF

"Strict"

"ACID"

"Strong"



Be Formal

Be Specific.

Figure out the

invariants

your system needs

Consider
your
failure modes

Consider
your
failure modes

Process Crash

```
#kill -9 1234
```

Node failure

- AWS terminate
- Physical power switch

Clock Skew

```
# date 1028000  
# fake time ...
```

GC/JQ Pause

killall -s STOP foo

killall -s CONT foo

Network Partition

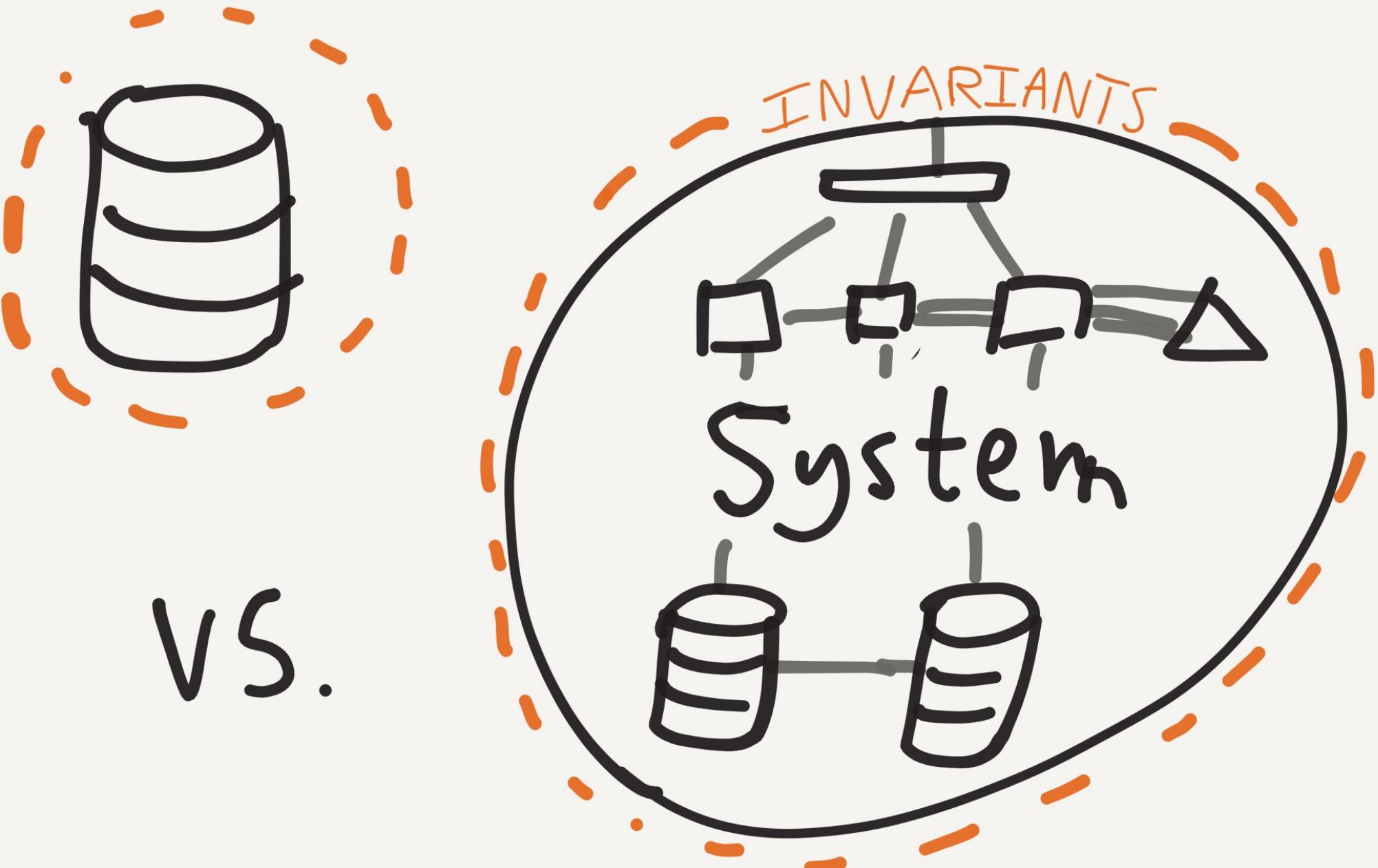
iptables -j DReP

tc qdisc ... delay...
drop...

Test your

Systems

| end ————— to —————> end |



If you look for

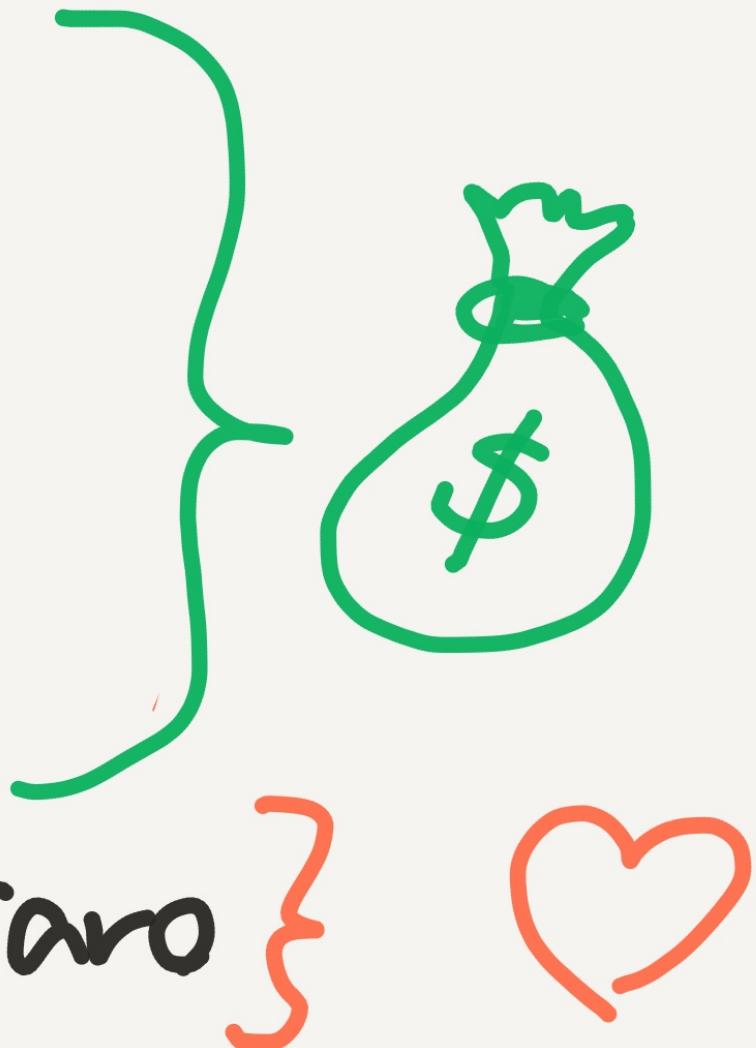
* perfection *

you will never be
content

- Property testing
- High-level invariants
- With distsys failure modes

Thanks

- VoltDB
- MongoDB
- Tendermint
- Peter Alvaro



<http://jepsen.ie>