

Distributed

Transaction

Architectures

Kyle Kingsbury

@aphyr

I Break
Databases!



jepseh.io

This work was
funded by

Fauna

Previous work funded by

- Mango
- Volt DB
- Cockroach Labs
- Yugabyte
- PingCAP

What

- do we -

Want?



Mutation is

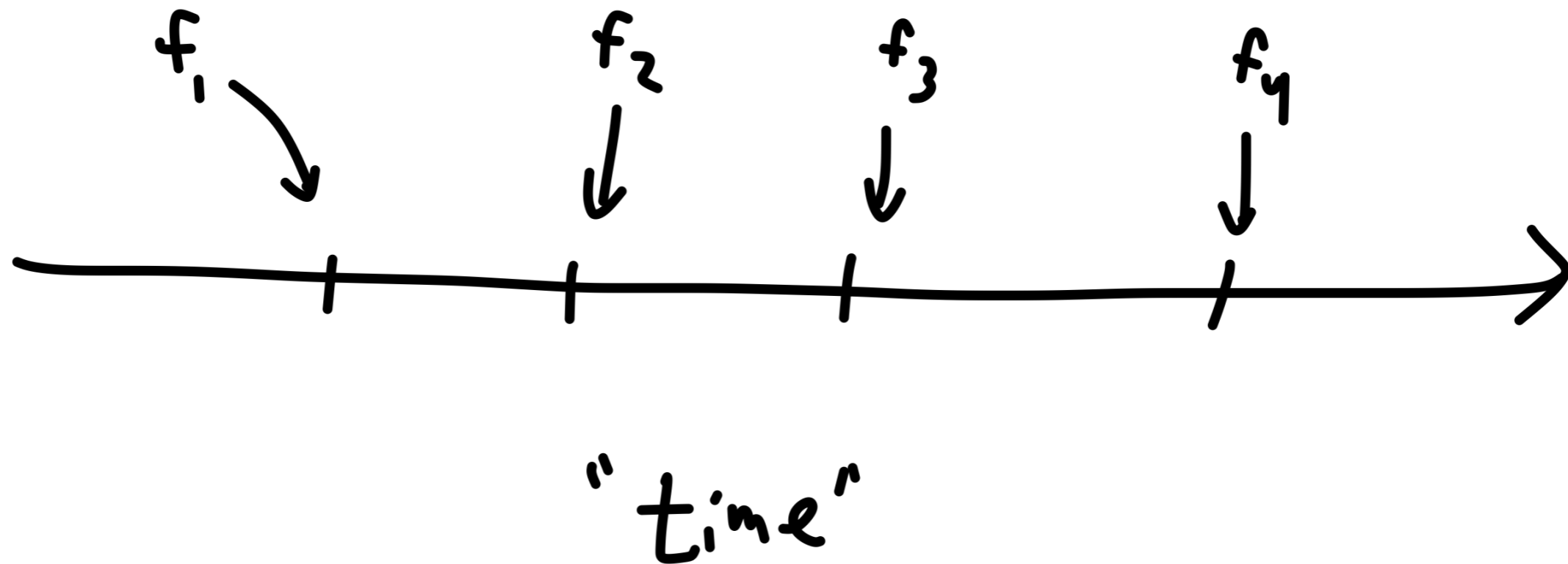
H A A R D

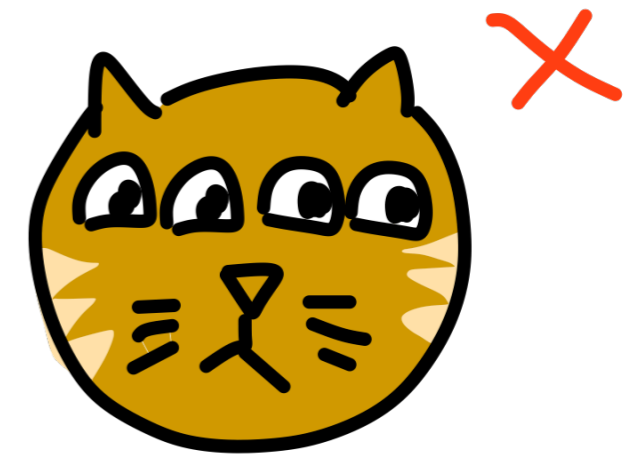
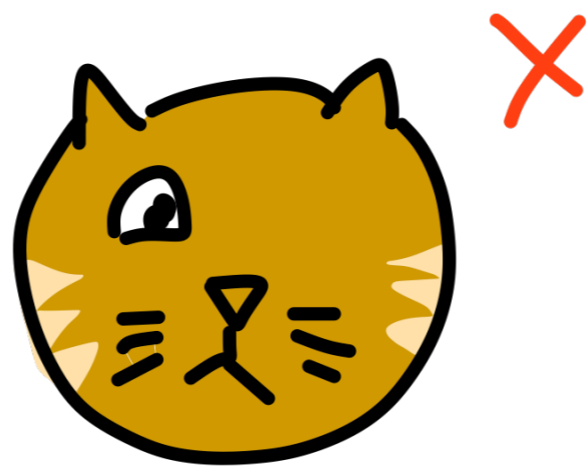
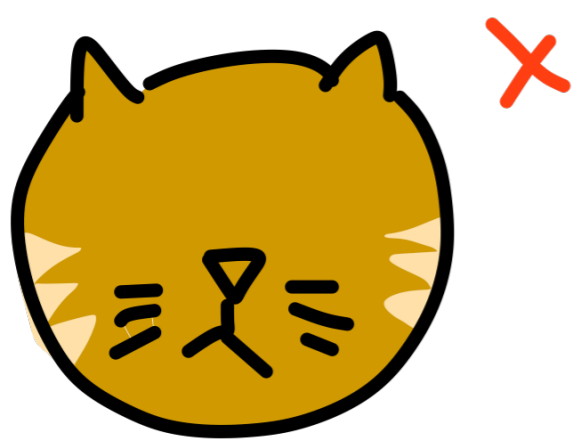
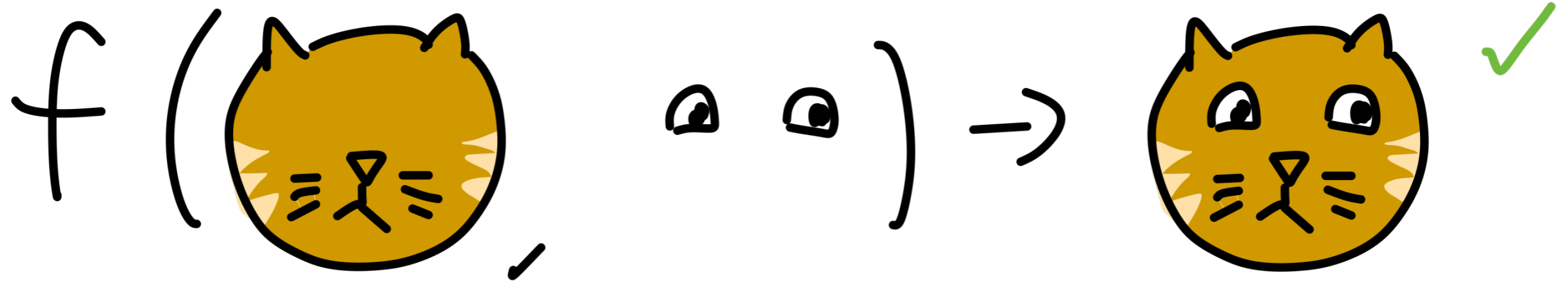
"

"

"

$$f(x, a) \rightarrow (x', r)$$





$$f(x, a) = f_3(f_2(f_1(x, a), a), a)$$

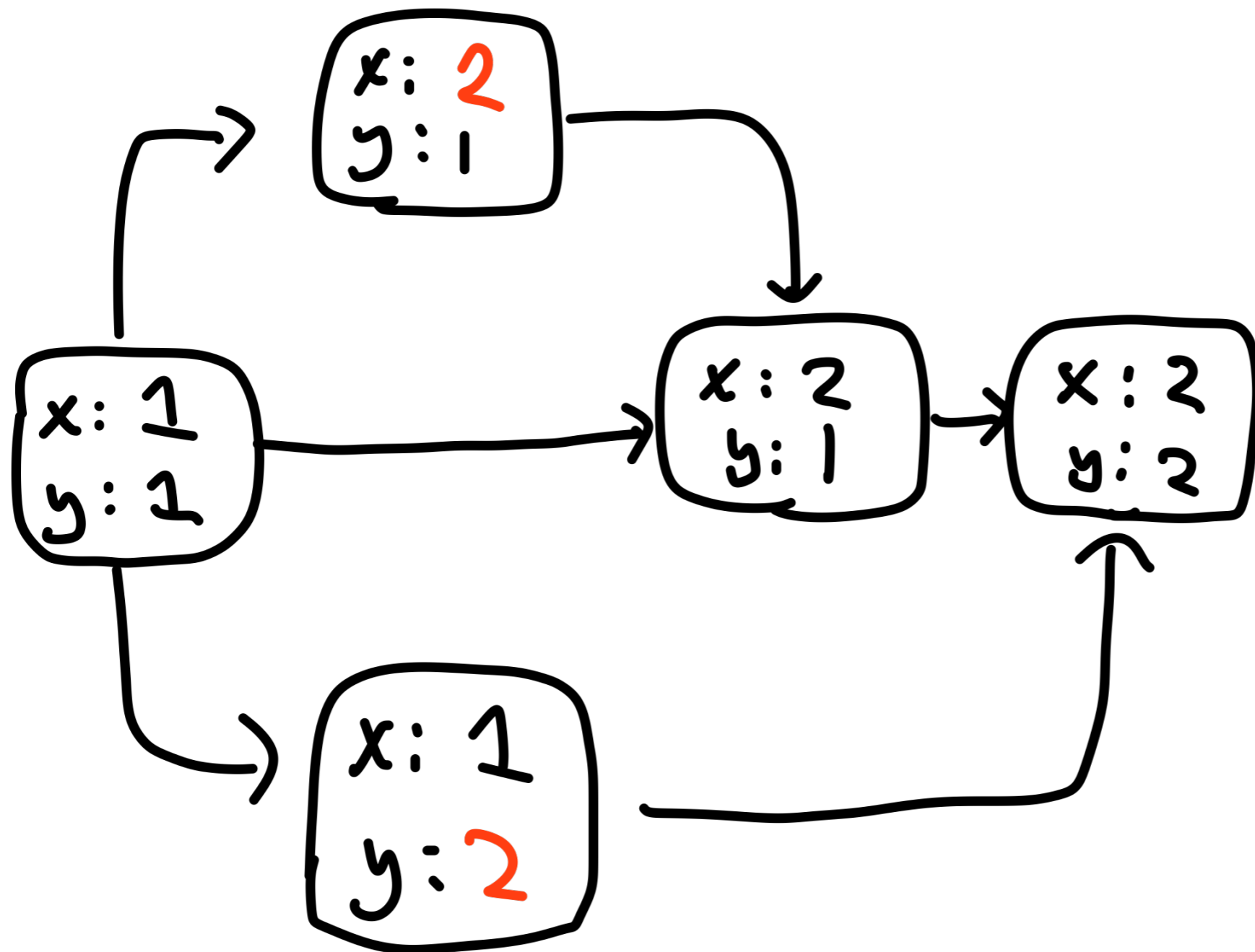
```
x = get ("ms fluffles")
```

```
y = get ("mx Cuddlekins")
```

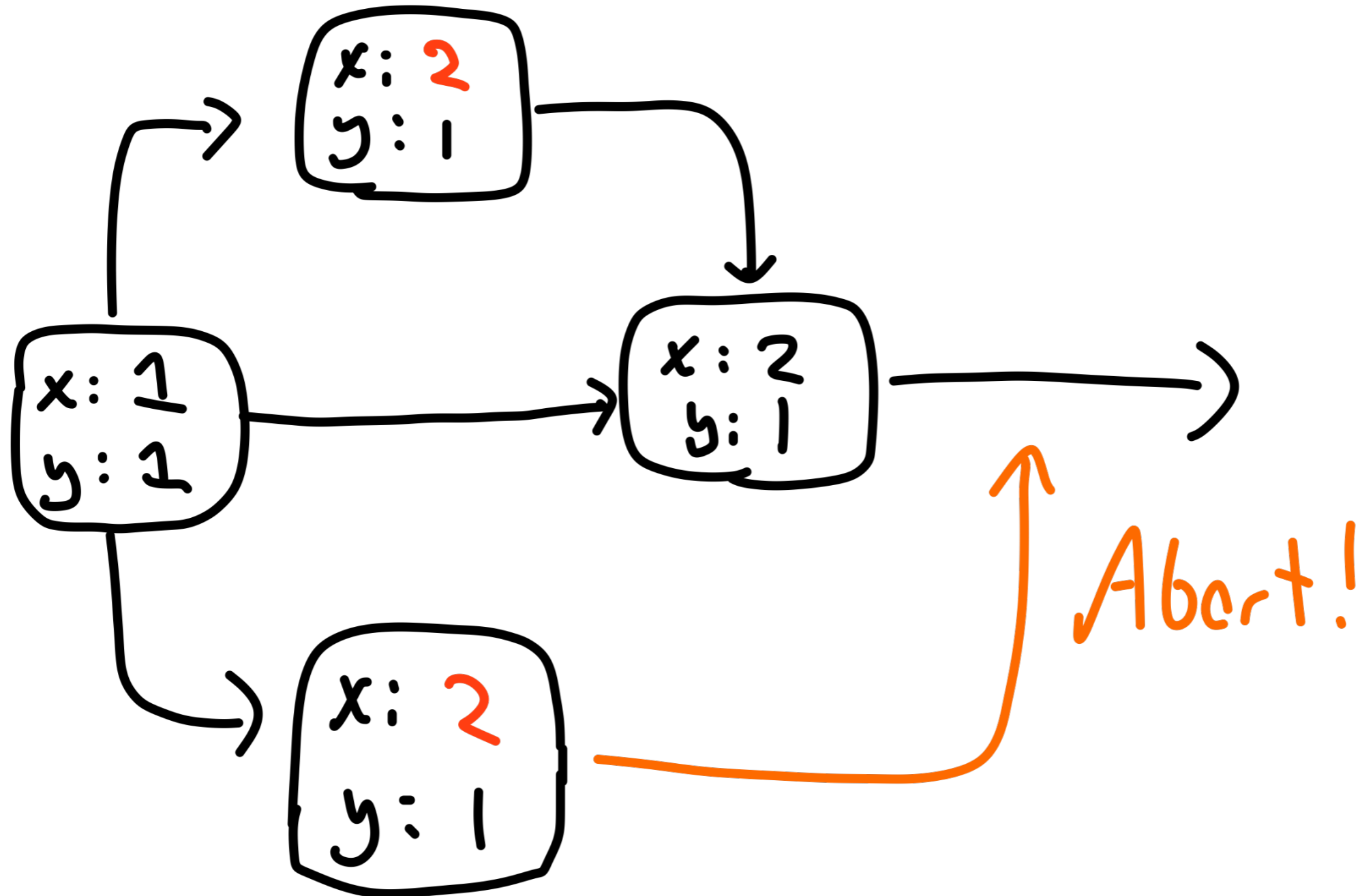
```
y.mice = x.mice.first
```

```
put (y)
```

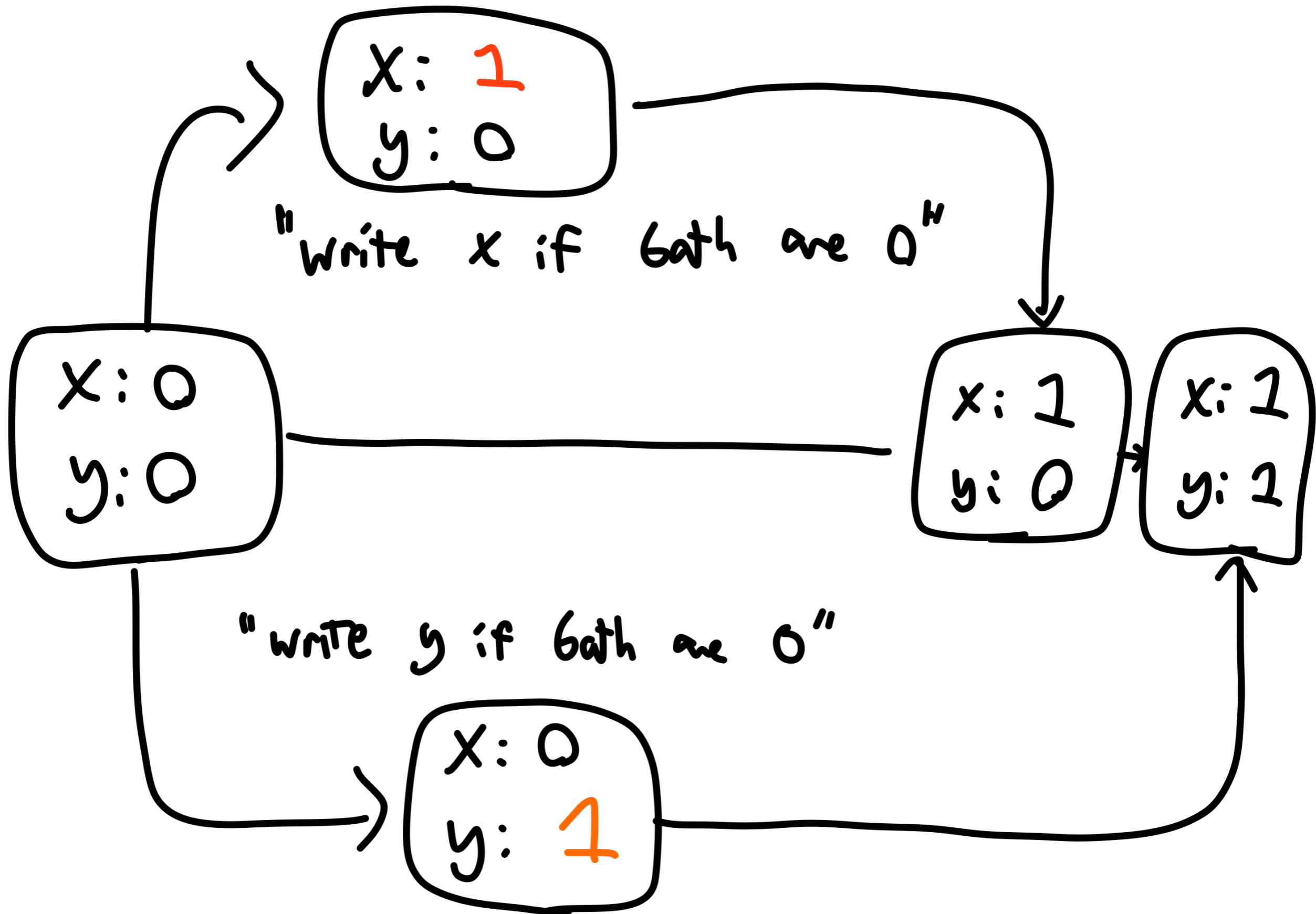
Snapshot Isolation



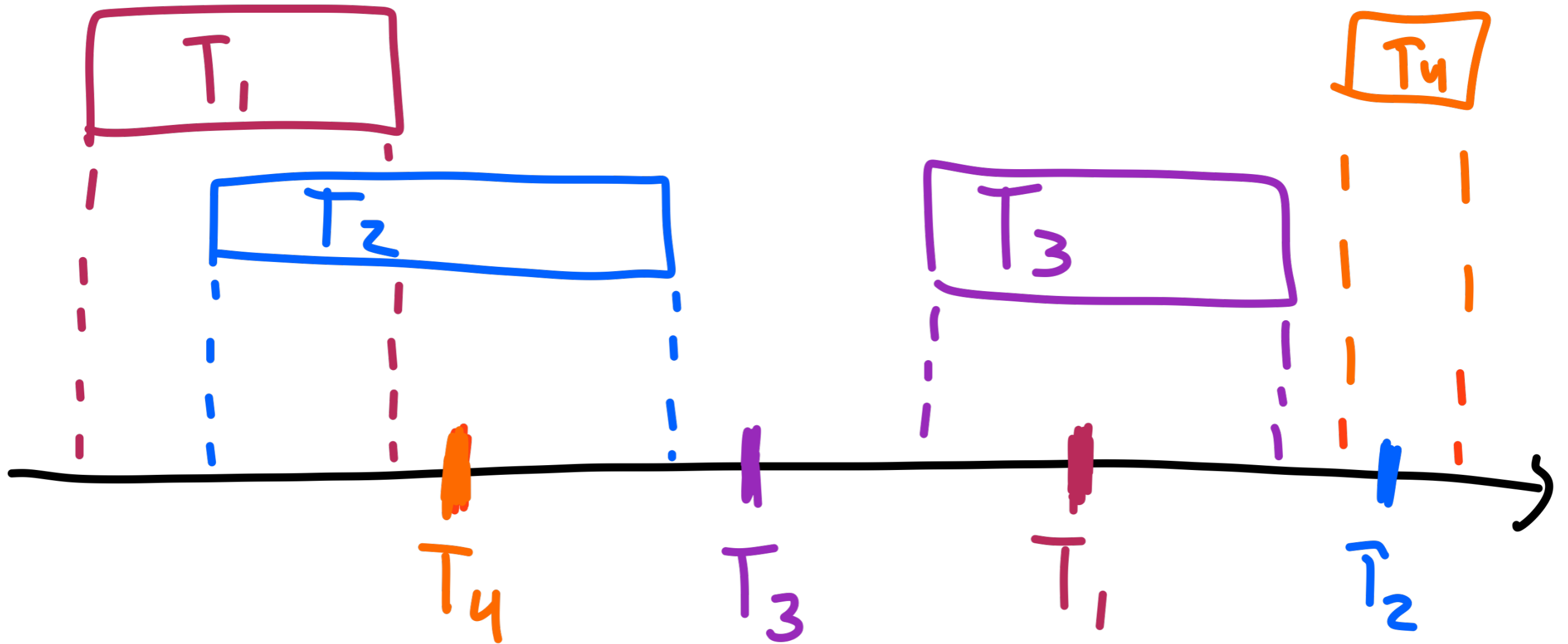
Snapshot Isolation



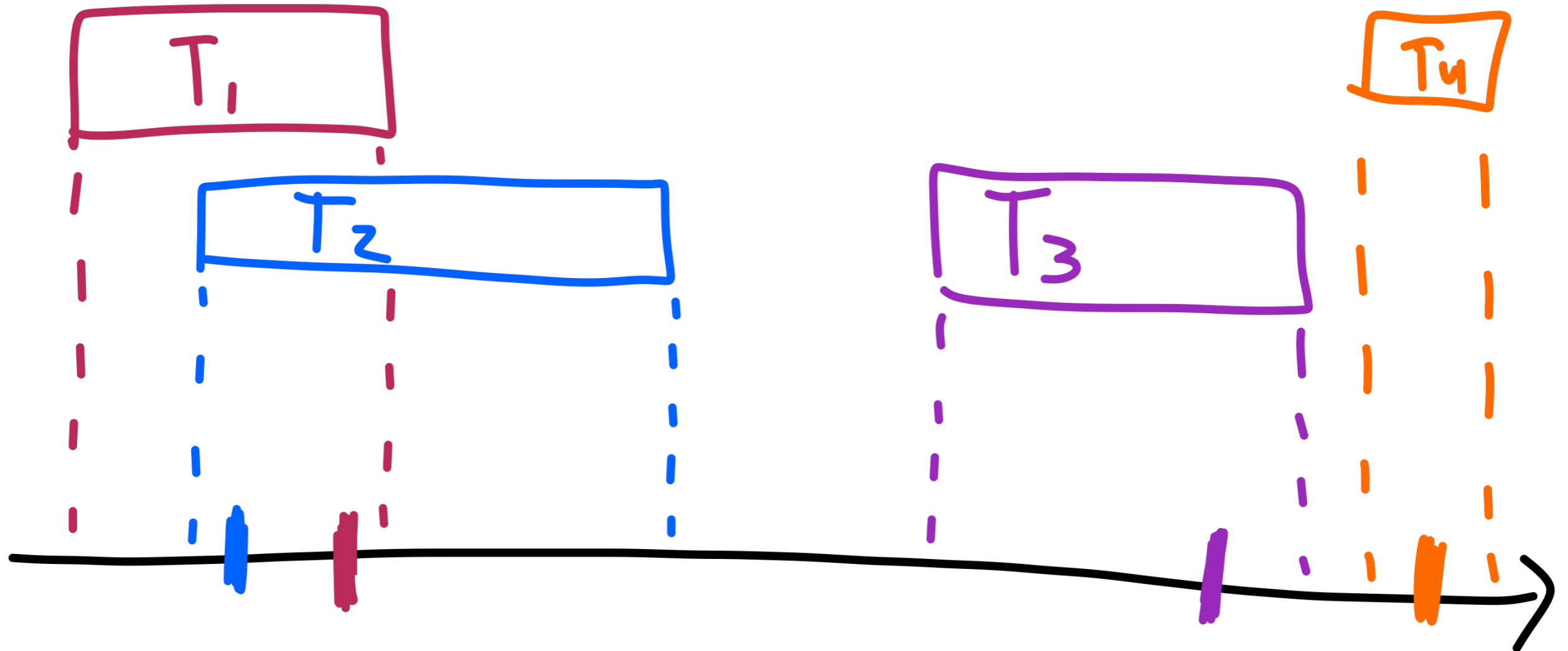
Write Skew



Serializable



Strict Serializable



↑ Where

— do we —

↑ Want it?

- Distributed

- Across Datacenters

- Worldwide

↑ When

— do we —

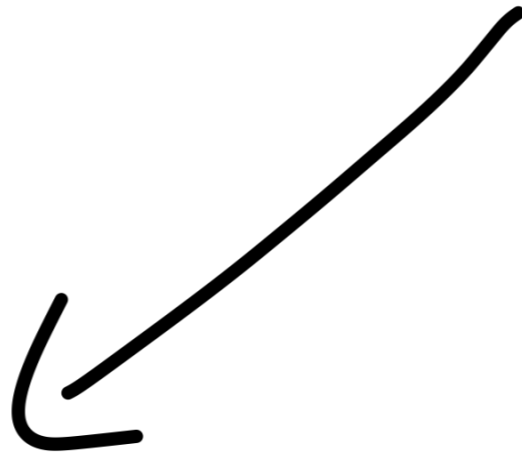
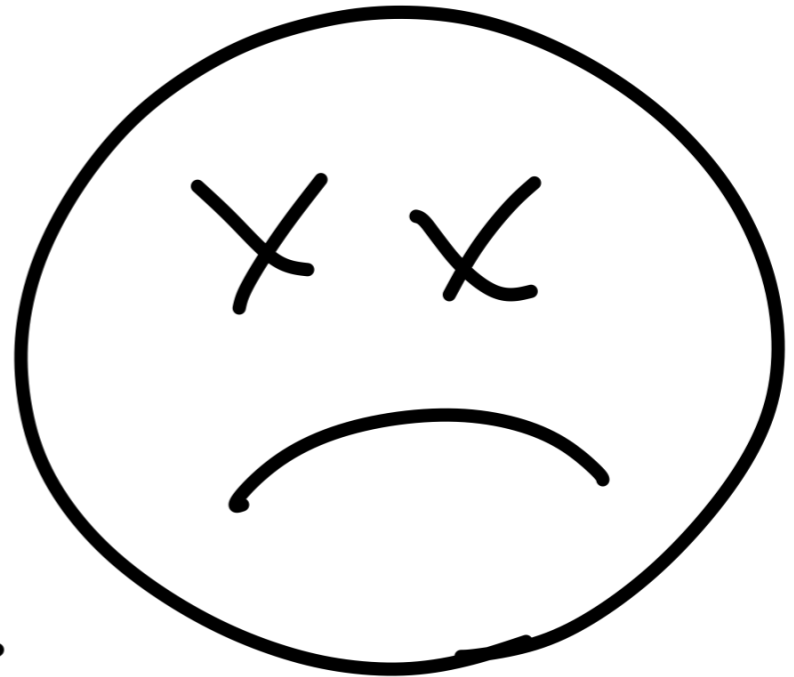
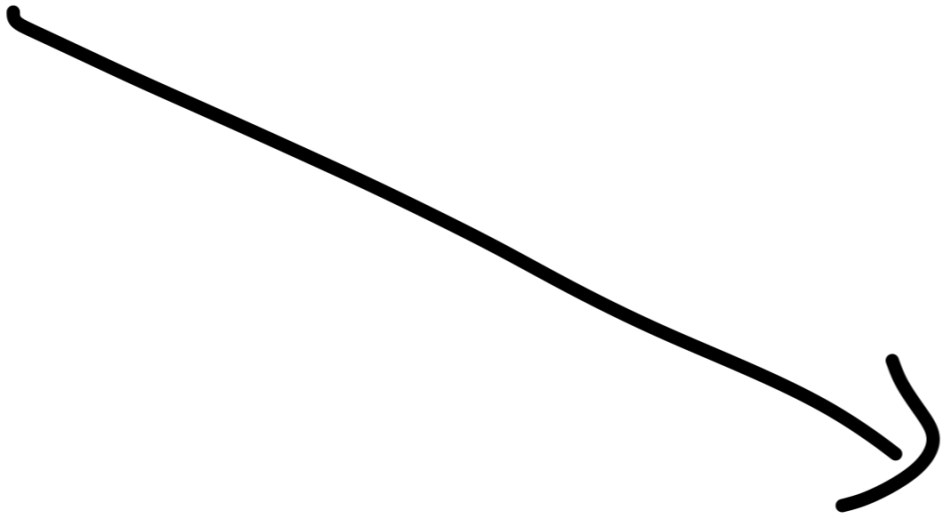
↑ Want It?.

Now!

Always!

Forever!

T,

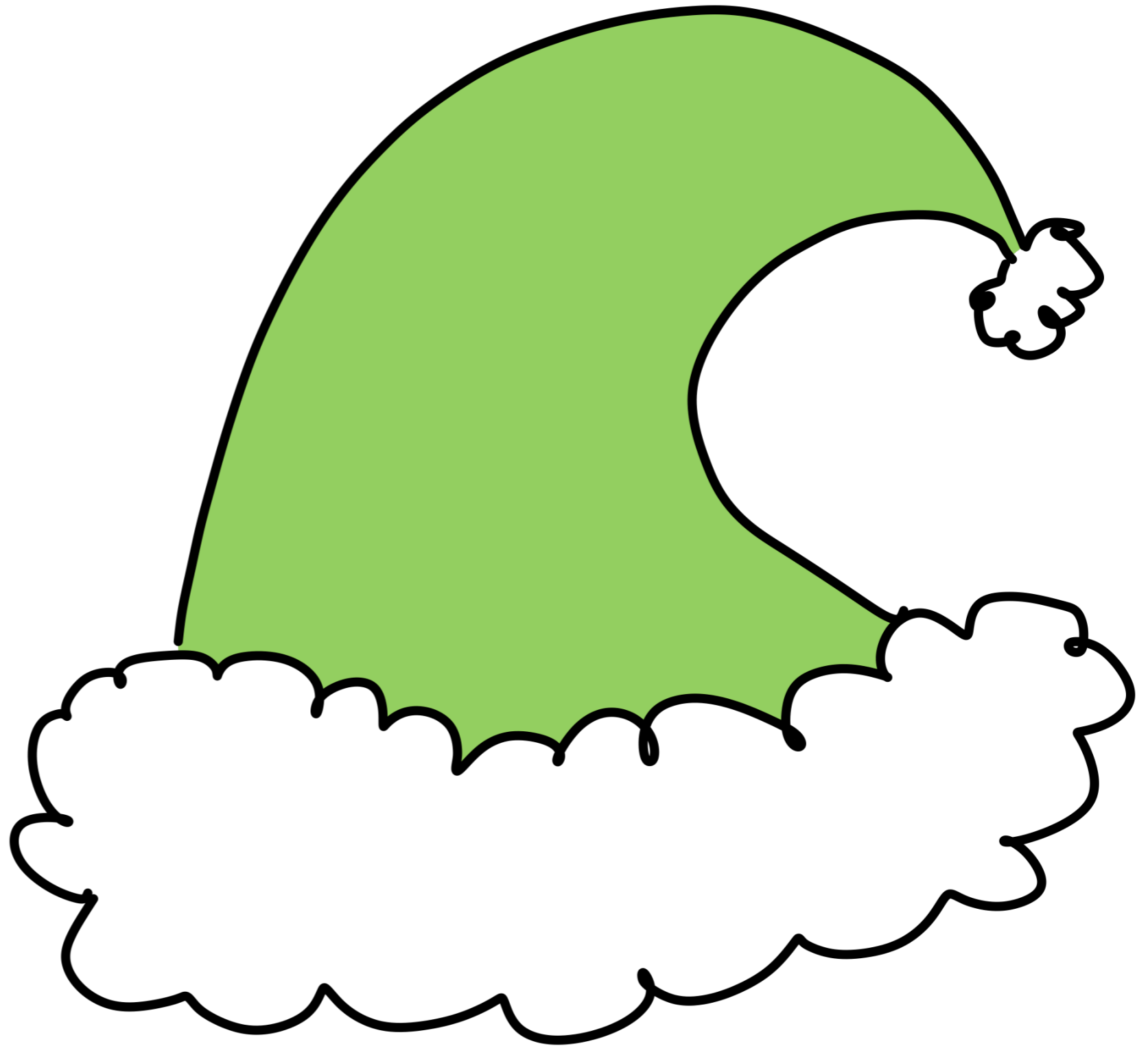


???

Total Availability

"Every request to a
non-faulty node succeeds"

the
CAP
theorem



Strict - 1SR

Serializable

Linearizable

Repeatable
Read

Snapshot
Isolation

Sequential

Cursor
stability

Monotonic
Atomic
View

Causal

Read

Committed

PRAM

Read

Uncommitted

WFR

MR

MW

RYW

If we can't have total
or sticky available... what

CAN we have for

strong consistency?

Consensus

Algorithms

- Nontriviality

- Safety

- Liveness

Lampart, 2002

- $n > 2f$

- Minimum 2 msg delays
(w/ conflicting requests)

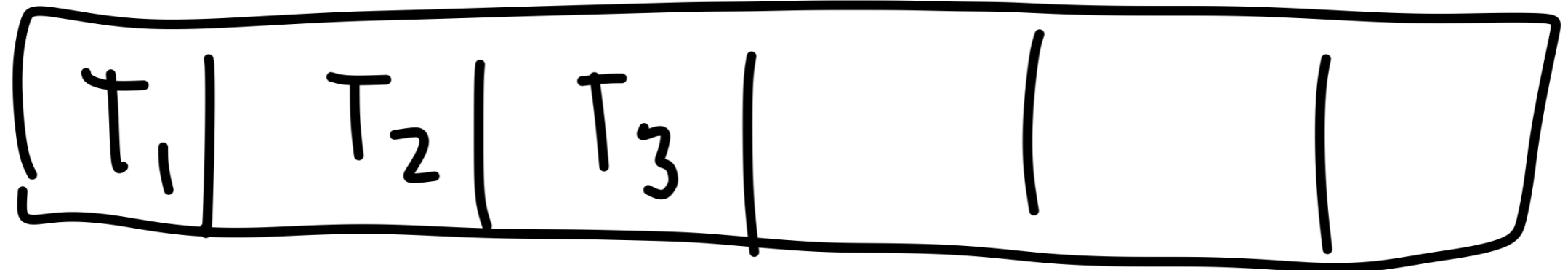
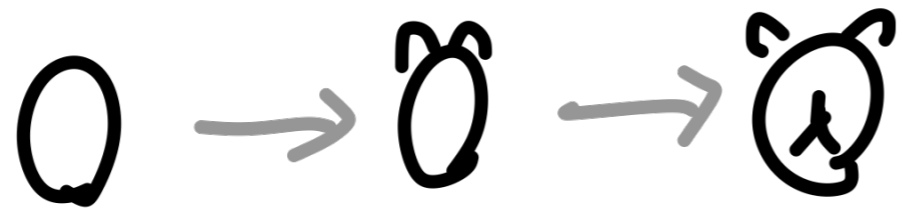
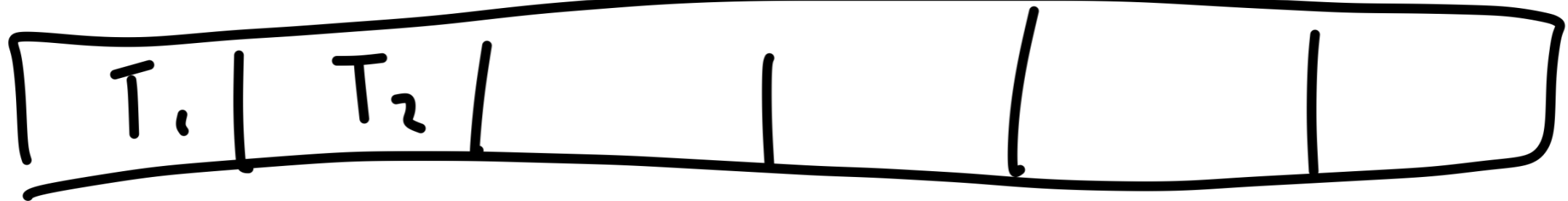
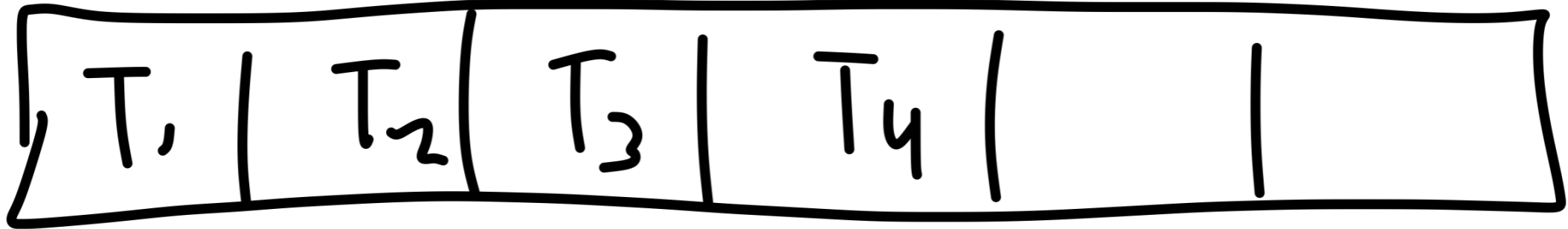
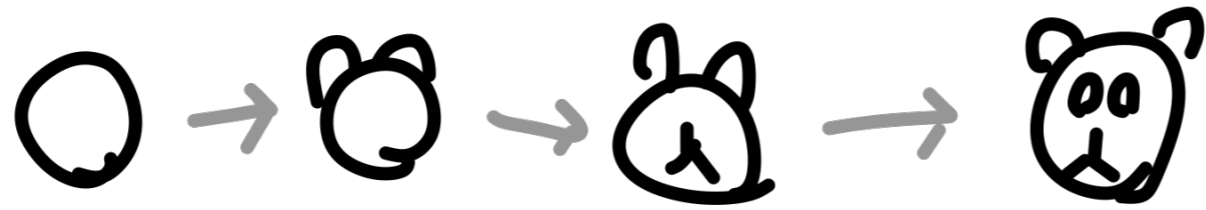
Practically Achievable

ZAB

Raft

Paxos

- Model system as a state machine
- Consensus on transitions
- And their order



- Linearizable in 2 hops
- Sequential reads if
desired, in O hops

- Mango DB

- Cassandra LWT

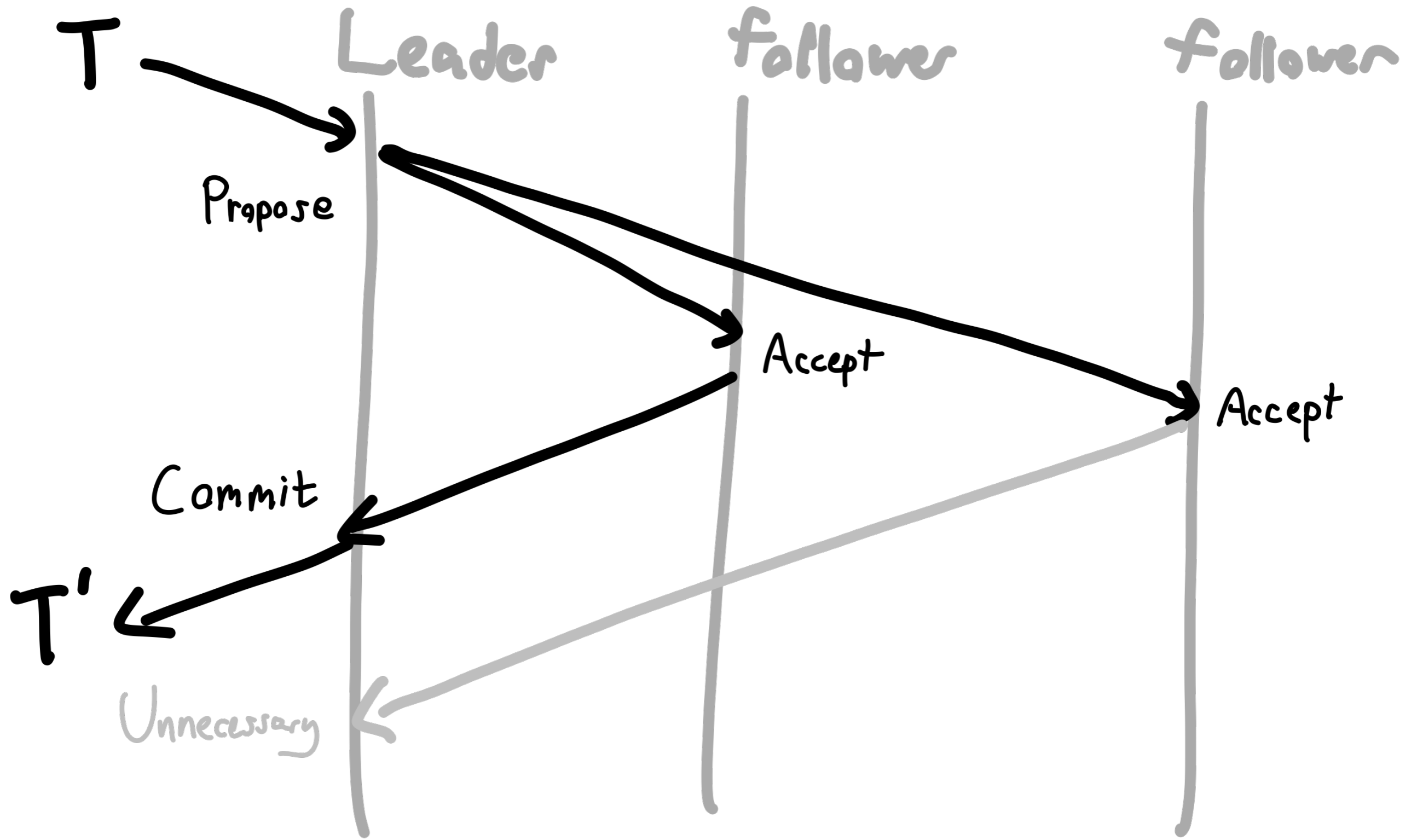
- Riak SC

- Aerospike

Majority Consensus

1 round trip

$\lfloor n/2 \rfloor + 1$ nodes to commit



The Consensus

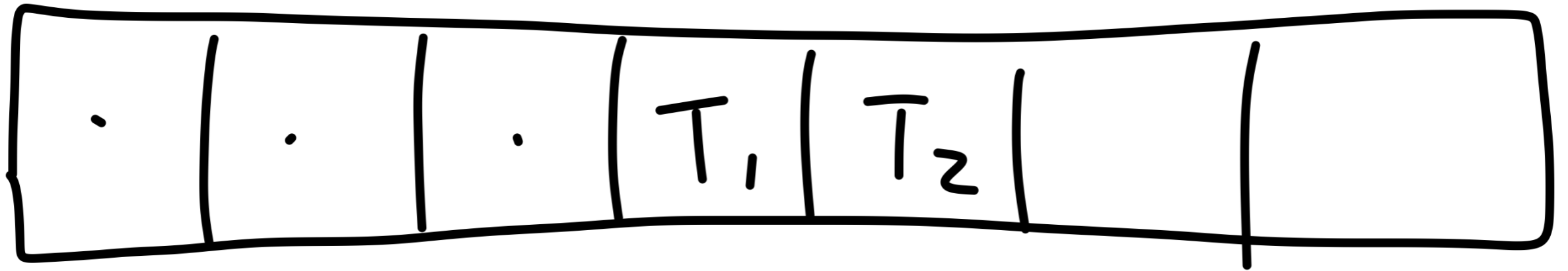


State machines are

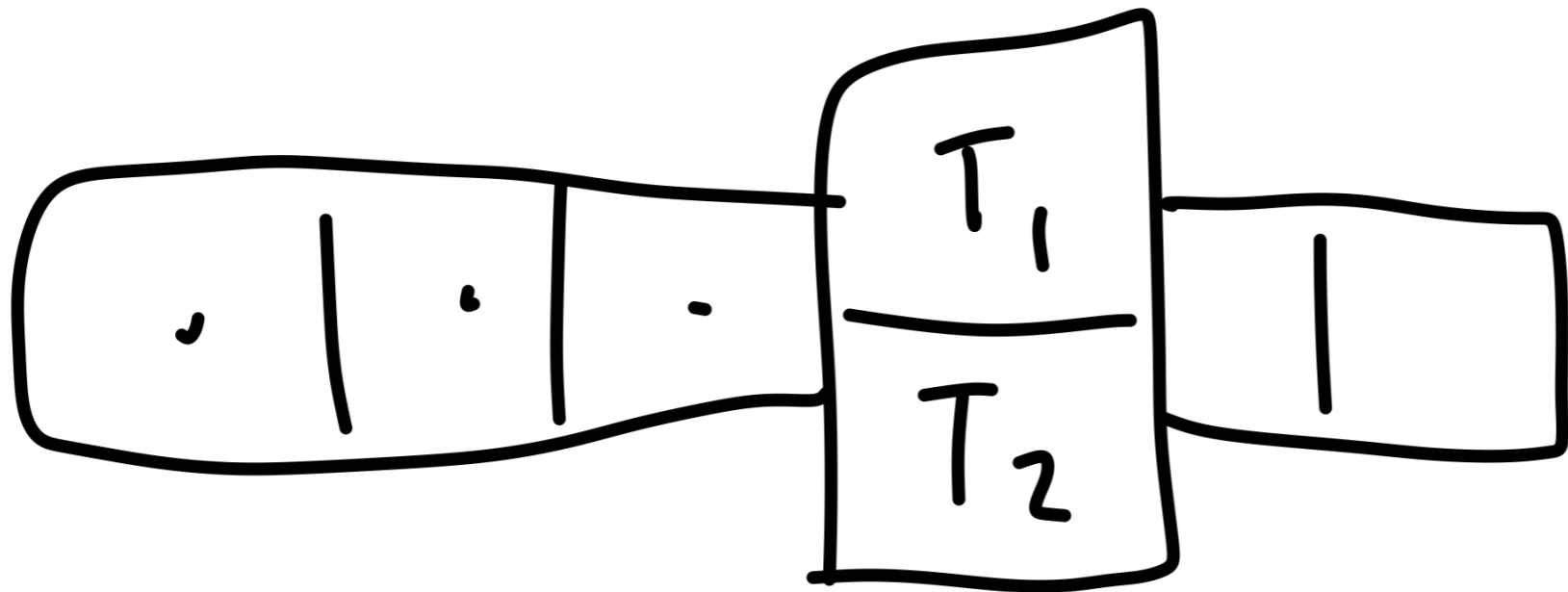
single-threaded, ... so

our DB is

single-threaded!



$$T_1 \circ T_2 = T_2 \circ T_1 ?$$



$$f(\cancel{X}) \rightarrow \cancel{X}$$

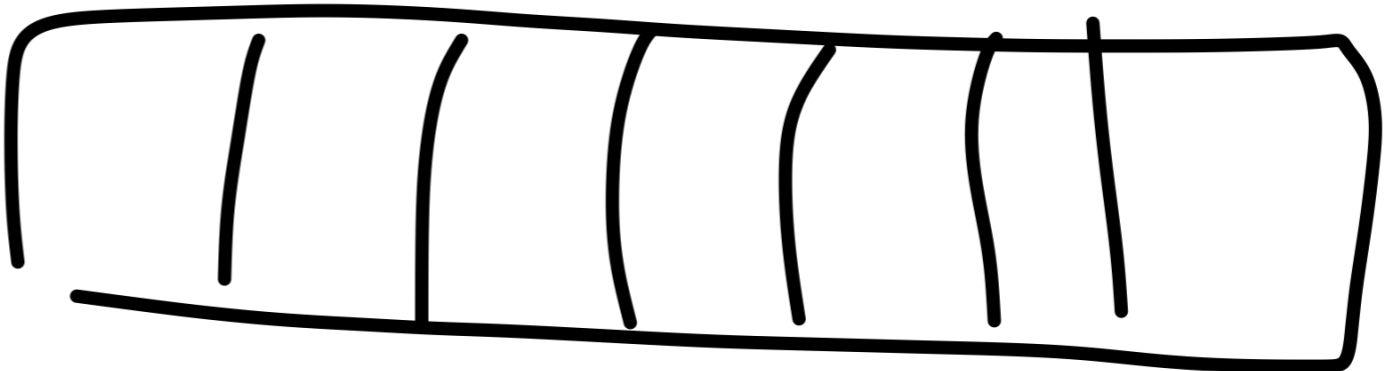
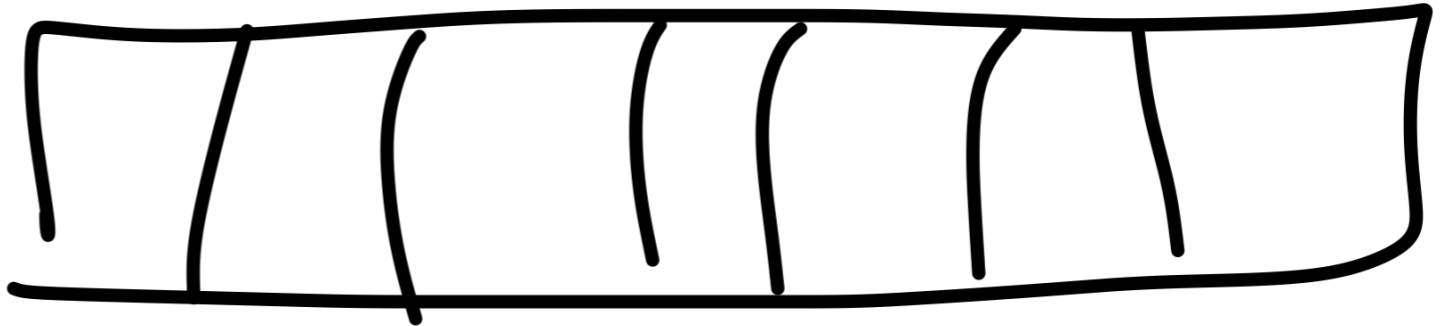
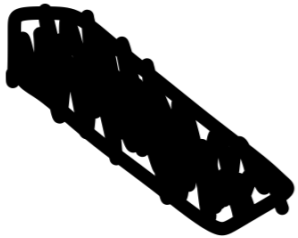
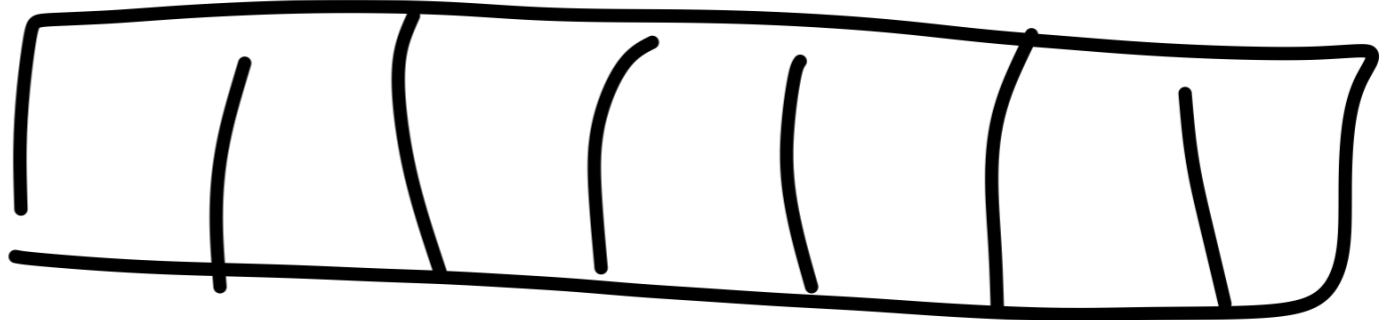
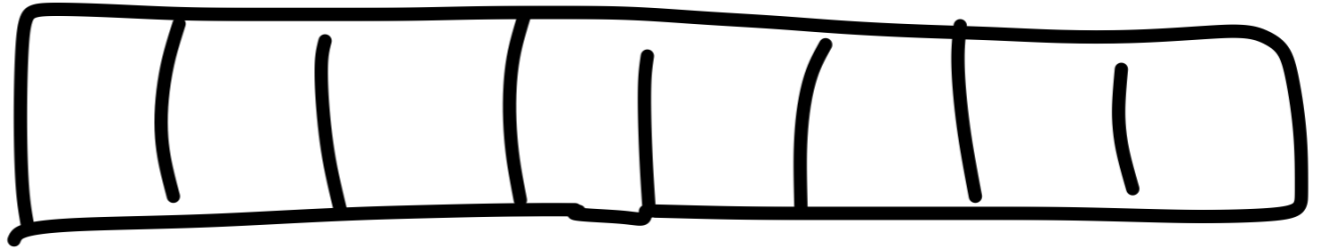


$f(\text{[scribble]}) \rightarrow \text{[orange scribble]}$

[scribble]

[scribble]

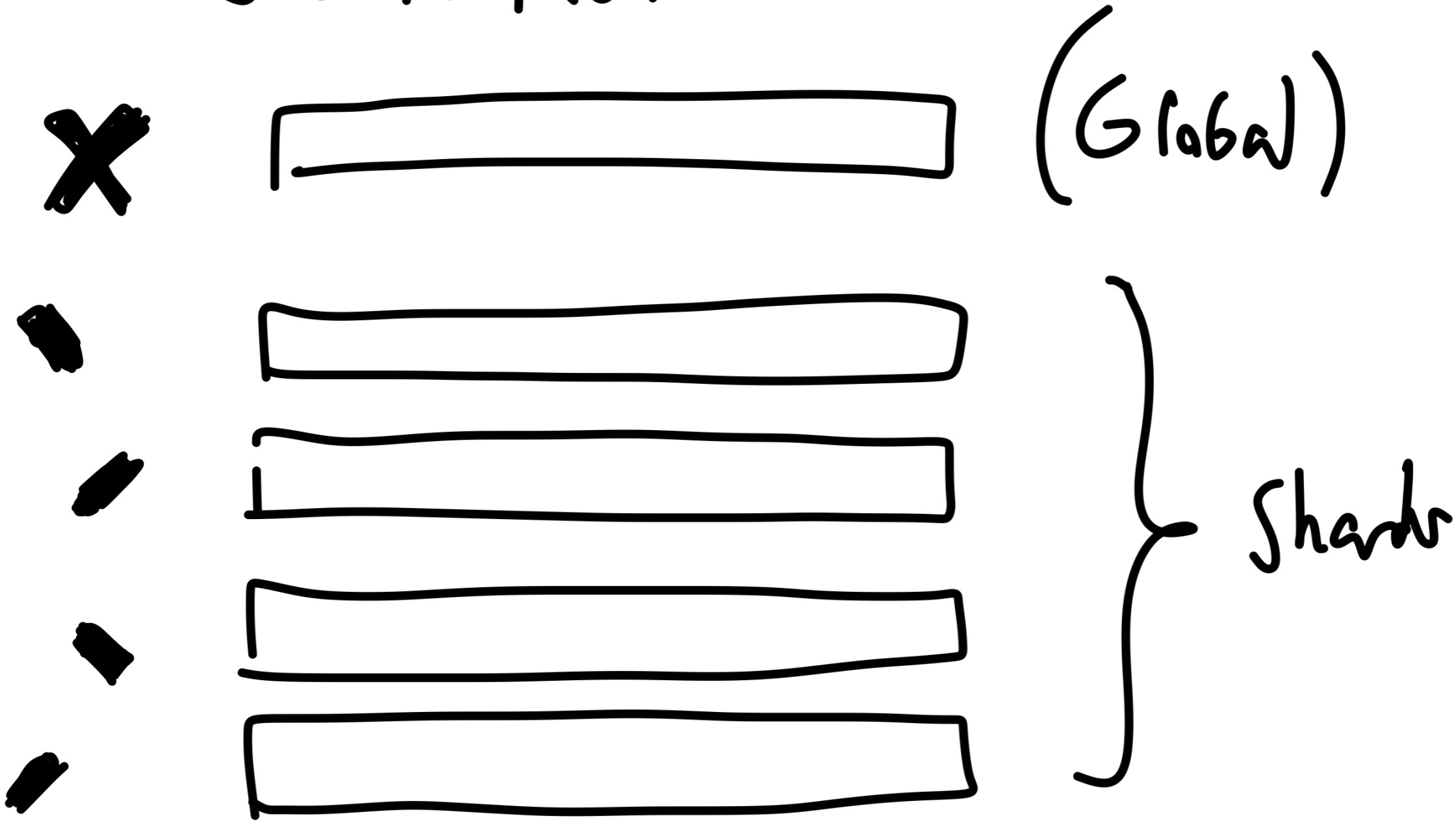
[scribble]



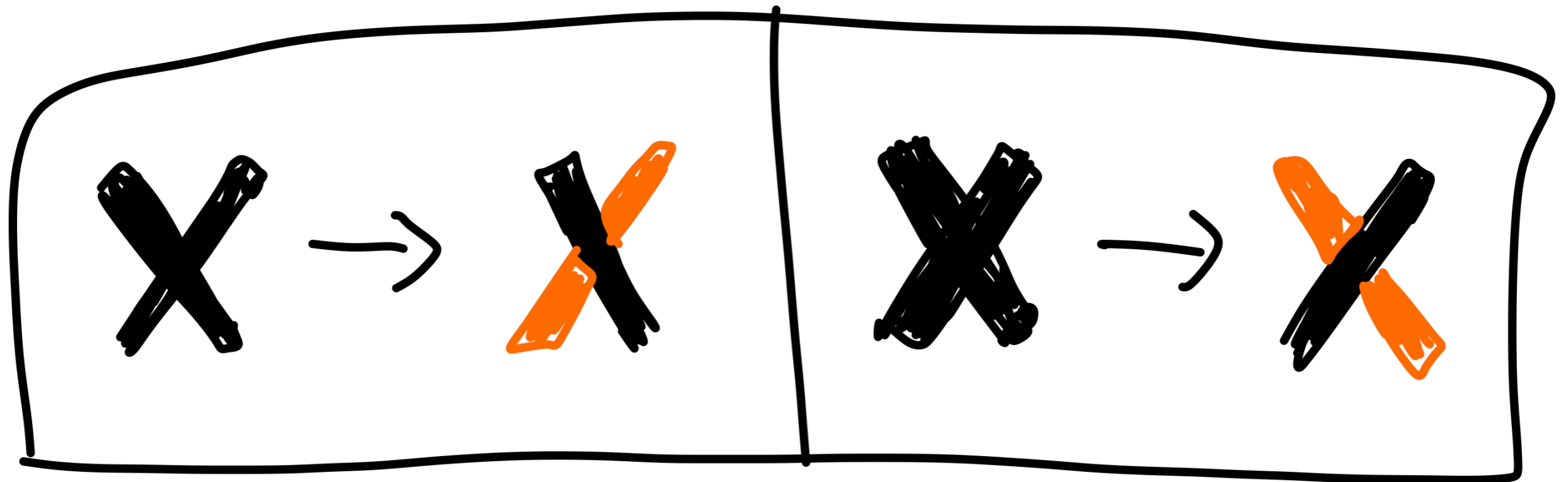
Sharding

(e.g. Val+DB)

What about cross-shard transactions?



Leaves a single
point of contention
for cross-shard transactions!



Let's Talk

WORST

Case

- Cross-Shard Transactions

- Strong (SI - Strict-ISR)

- Globally Distributed

- Site redundancy

- Access everywhere

You CAN do better!

- Stay inside shard

- Relax consistency

- Relax durability

- Colocate clients w/data

- Colocate shard leaders

Percolator

Google 2010

Peng & Dabeek

Alternative to MapReduce

for txnal index building

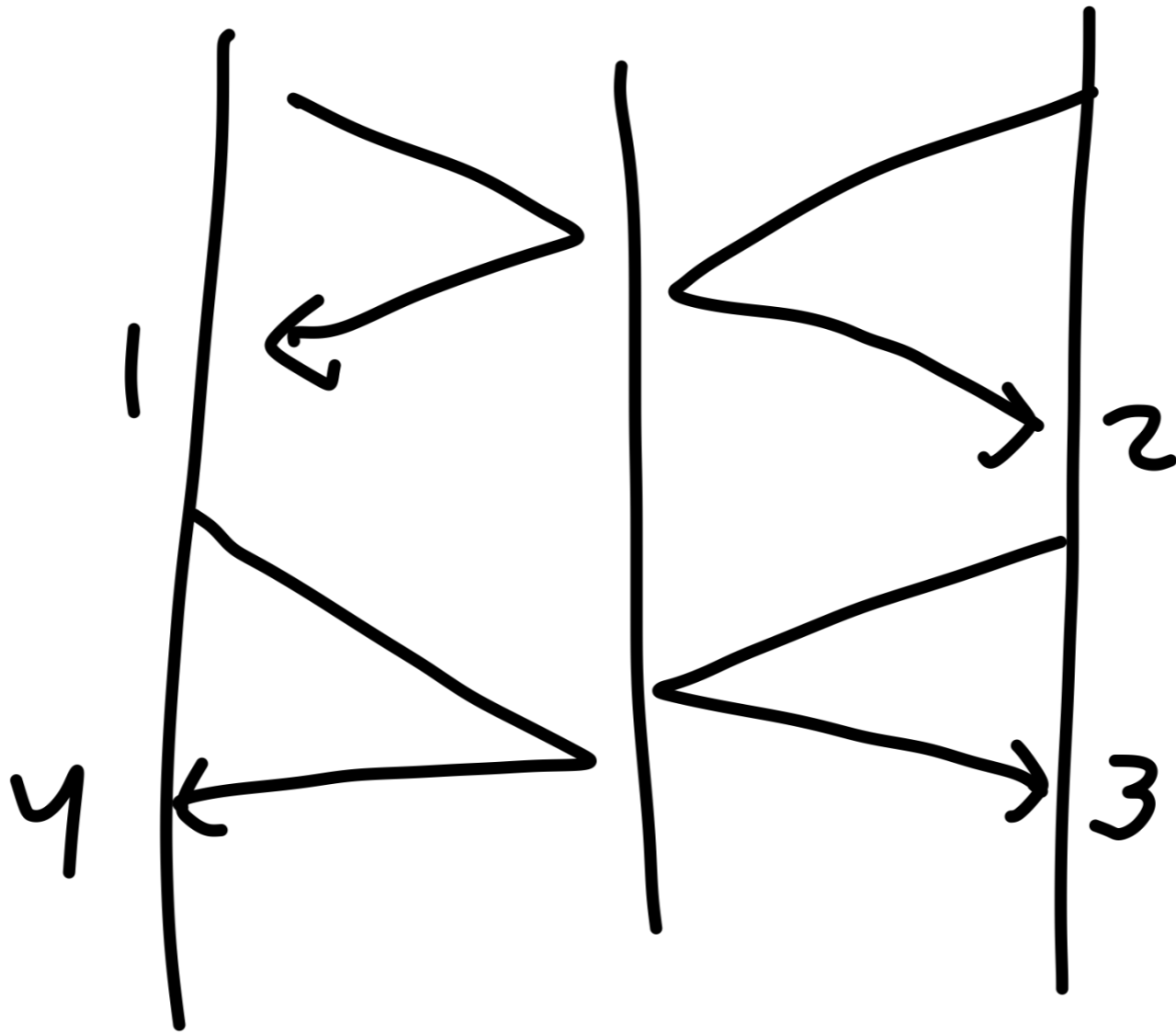
- "ACID Txns" (really SI)
- Fast read only txns
- 2PC for updates
- Relaxed latency

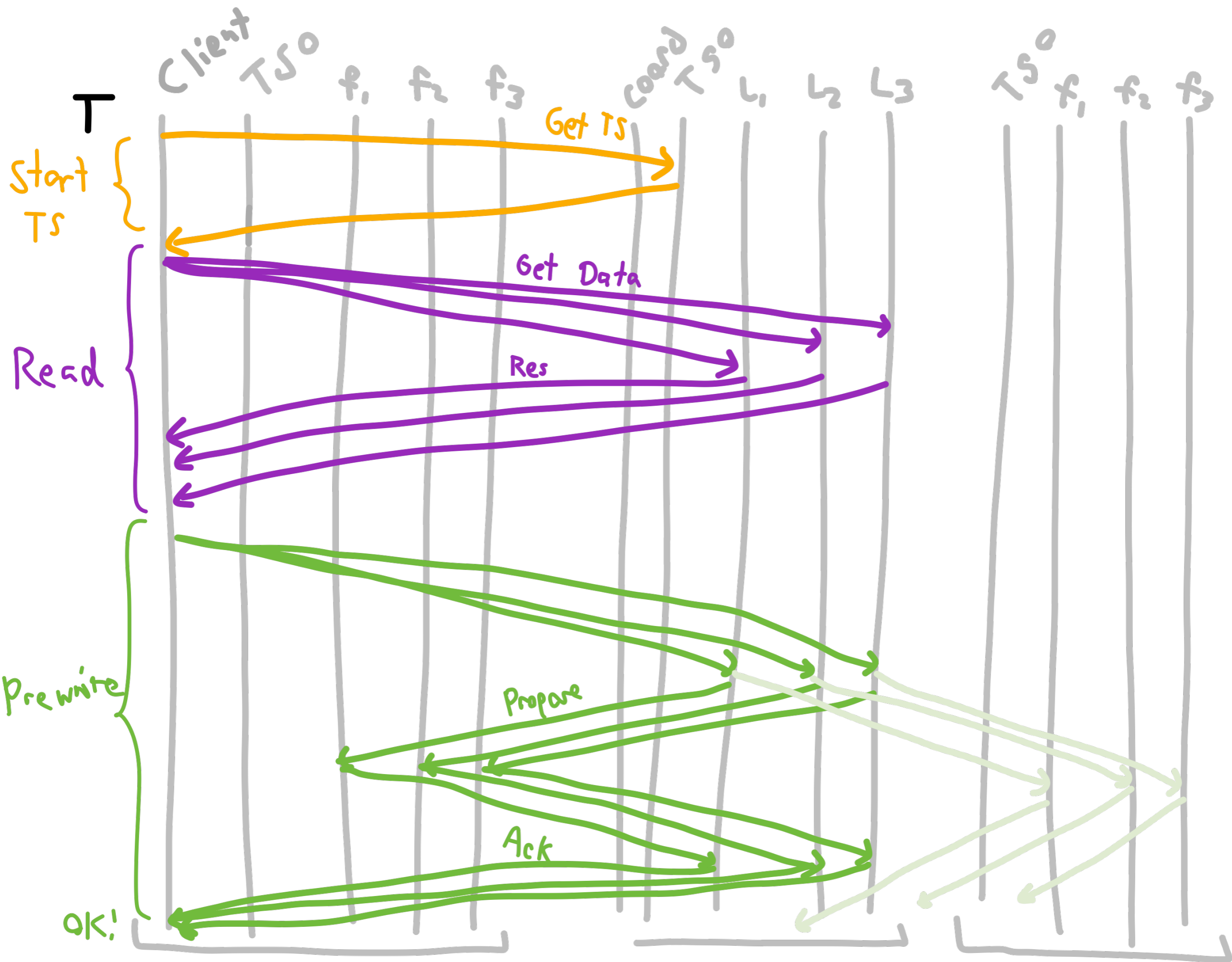
Built on BigTable

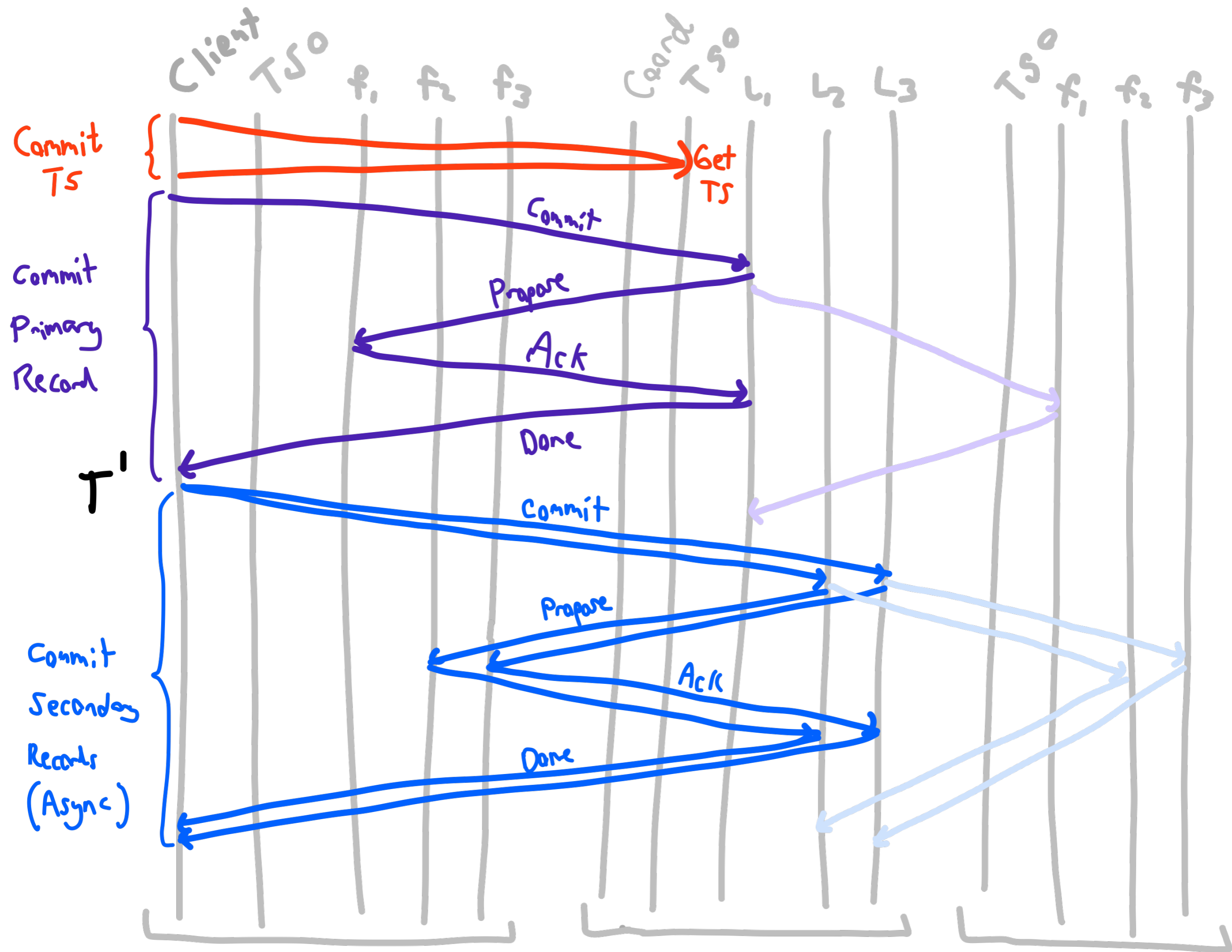
- k/v store
- Atomic read - modify - write
- No multi-key transactions
- Percolator is a client lib!

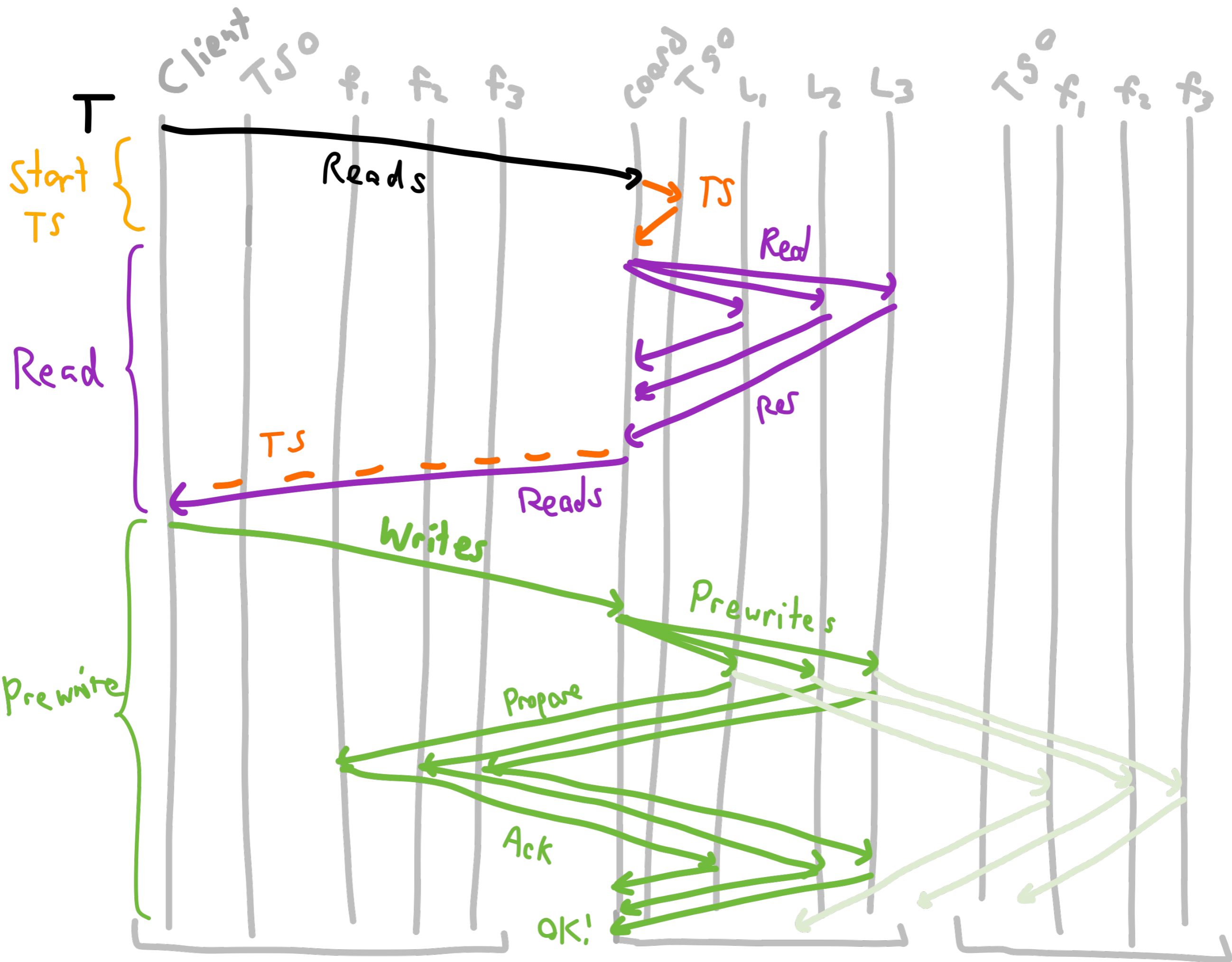
Time Stamp Oracle

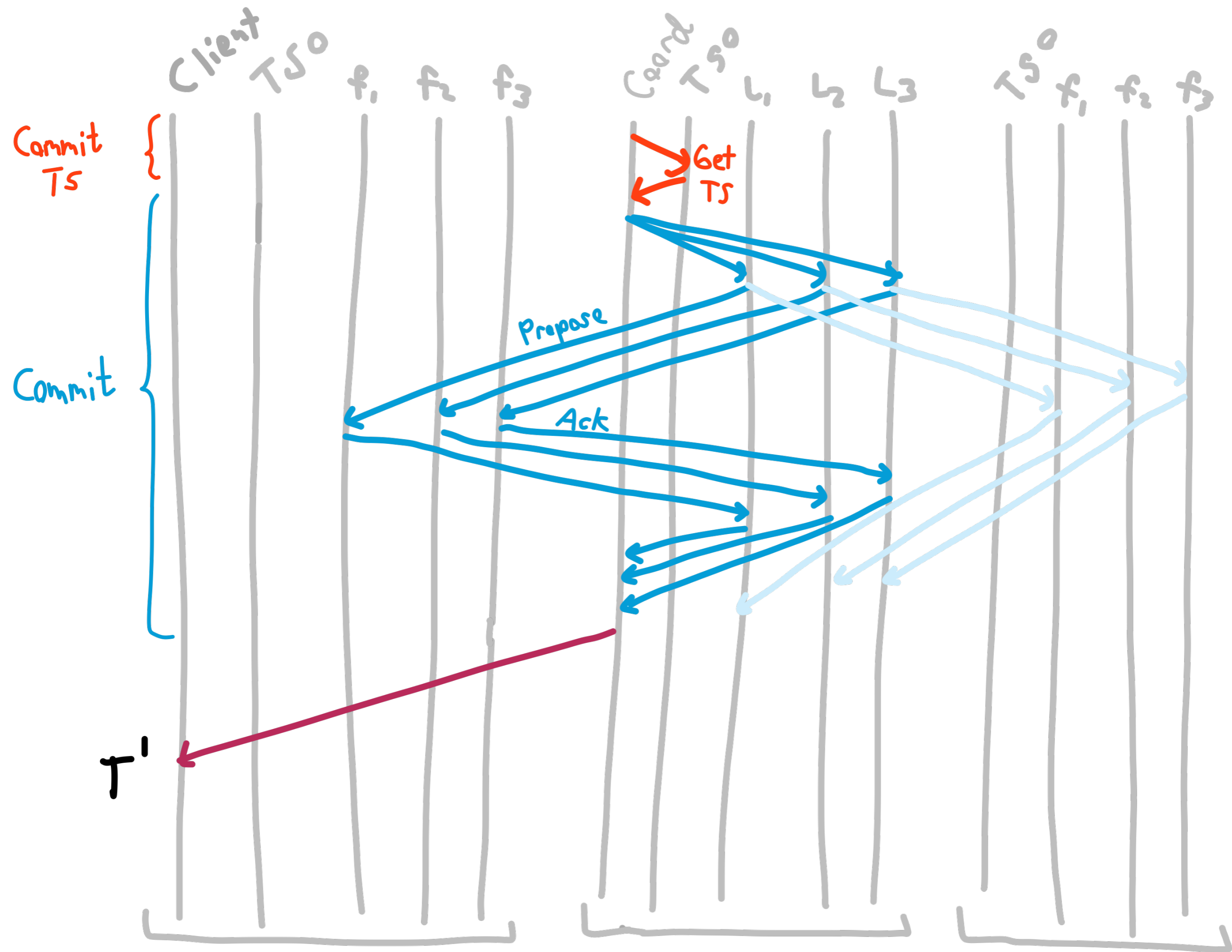
TSO











R/w Txns

14 Cross-DC haps

0 local haps

T; DB

- Commercial/OSS implementation
- SI K/V store + MySQL layer
- Raft for each shard
- Timestamp via Raft + leader leases

Spanner

Google 2012

- Carbett, Dean, Epstein, ...
- Globally replicated SQL-style

OLTP DB

- Versioned

- Fast snapshot reads

- "External Consistency"

(Strict - 1SR?)

TSO \rightarrow True Time

- Service

- Paxos

- latency insensitive?

- library

- magic clocks

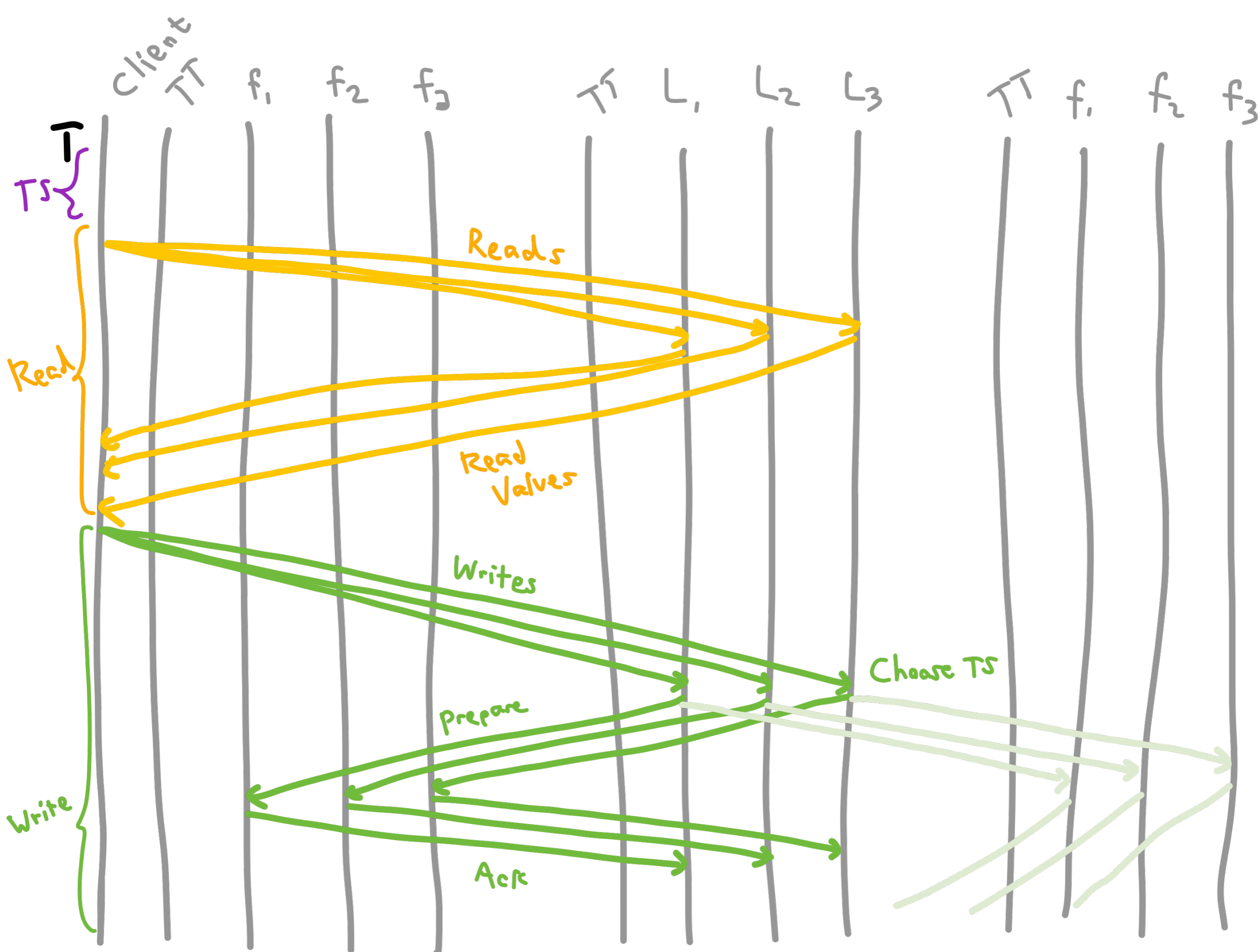
- latency sensitive?

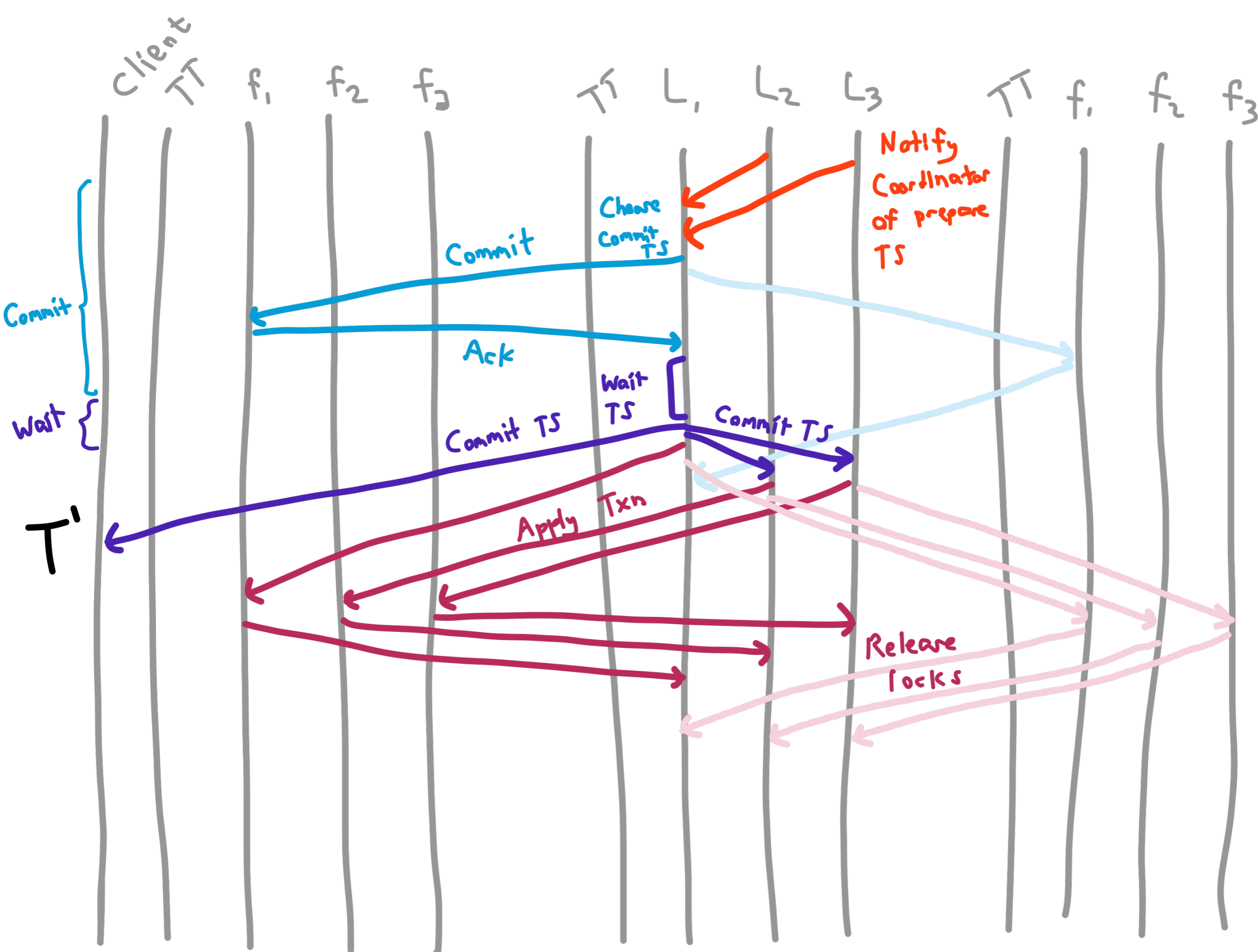
True Time is Unusual

- GPS + atomic clocks
- Voting out bad nodes
- Clock sync algo
- Dedicated network & nodes
 - Instrument GC / IO pauses?

Basically 2PC over Paxos

- TS \Rightarrow lock intervals
- Locks taken on Paxos leader
- Coordinator Paxos group commits
- Block to ensure TS passed
- Clean up locks async





R/w Txns

8 Cross-DC hops

1 local hop

R/O Txns

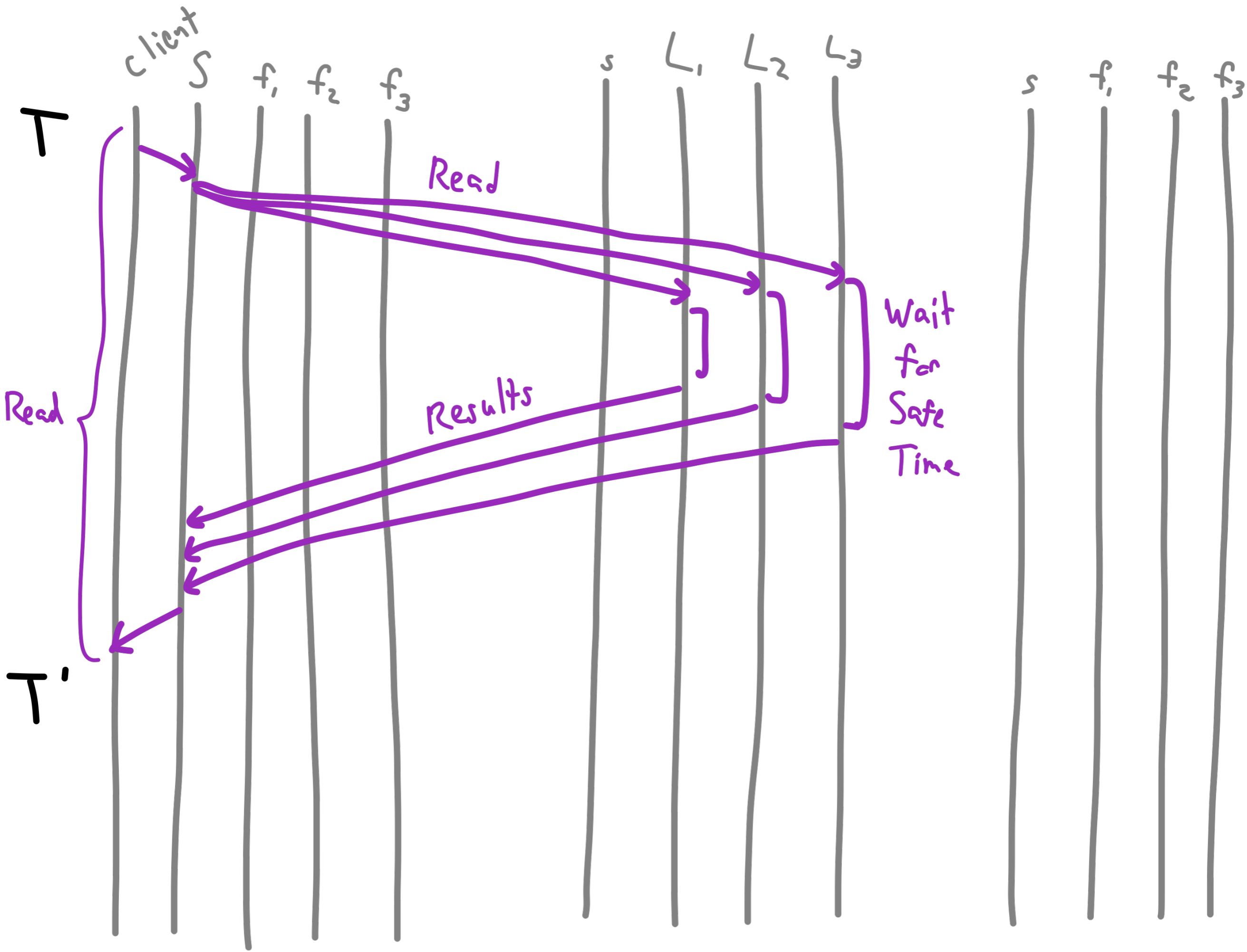
- Pick timestamp locally, read locally
- Linearizable in O round

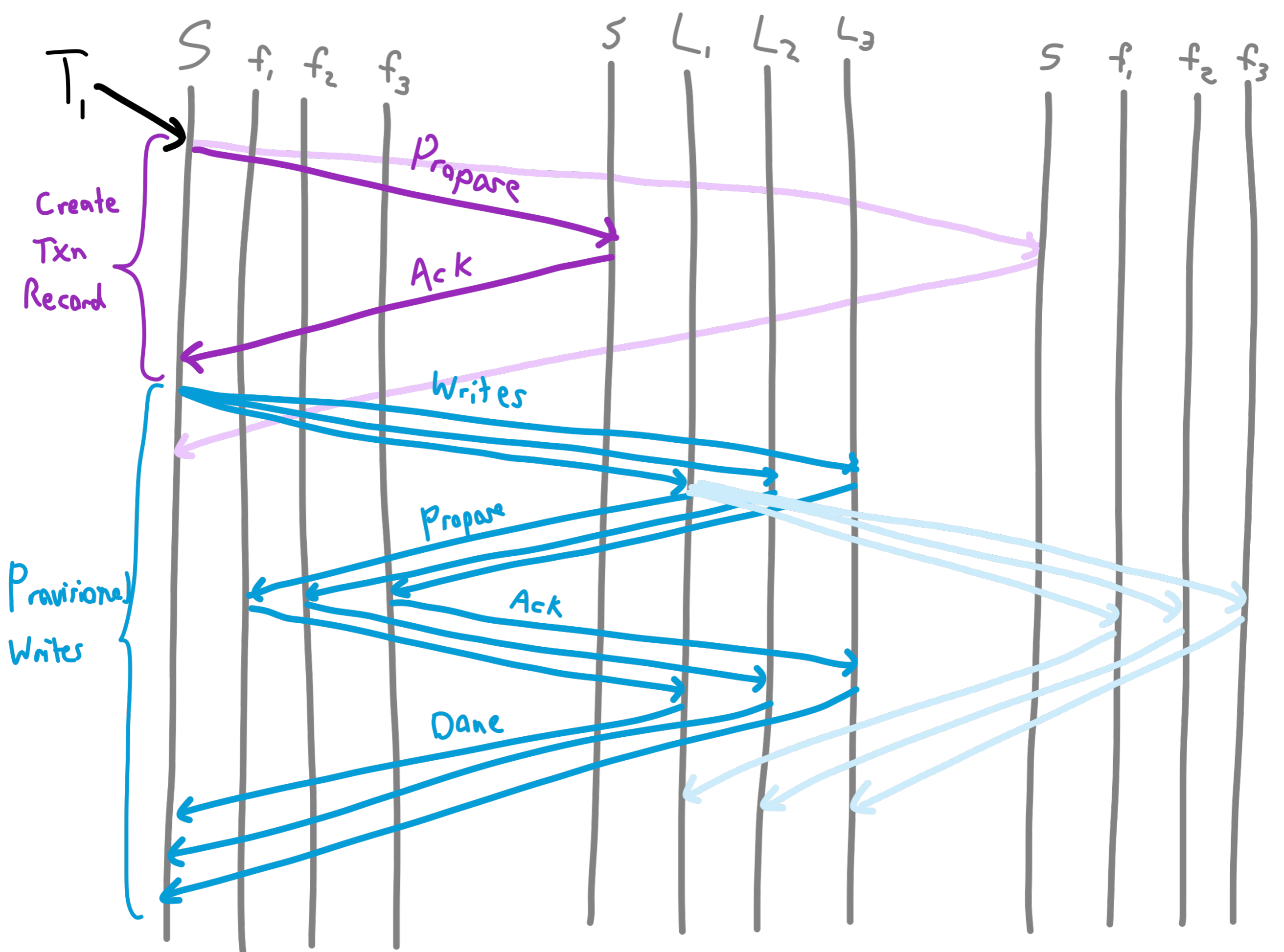
tips?

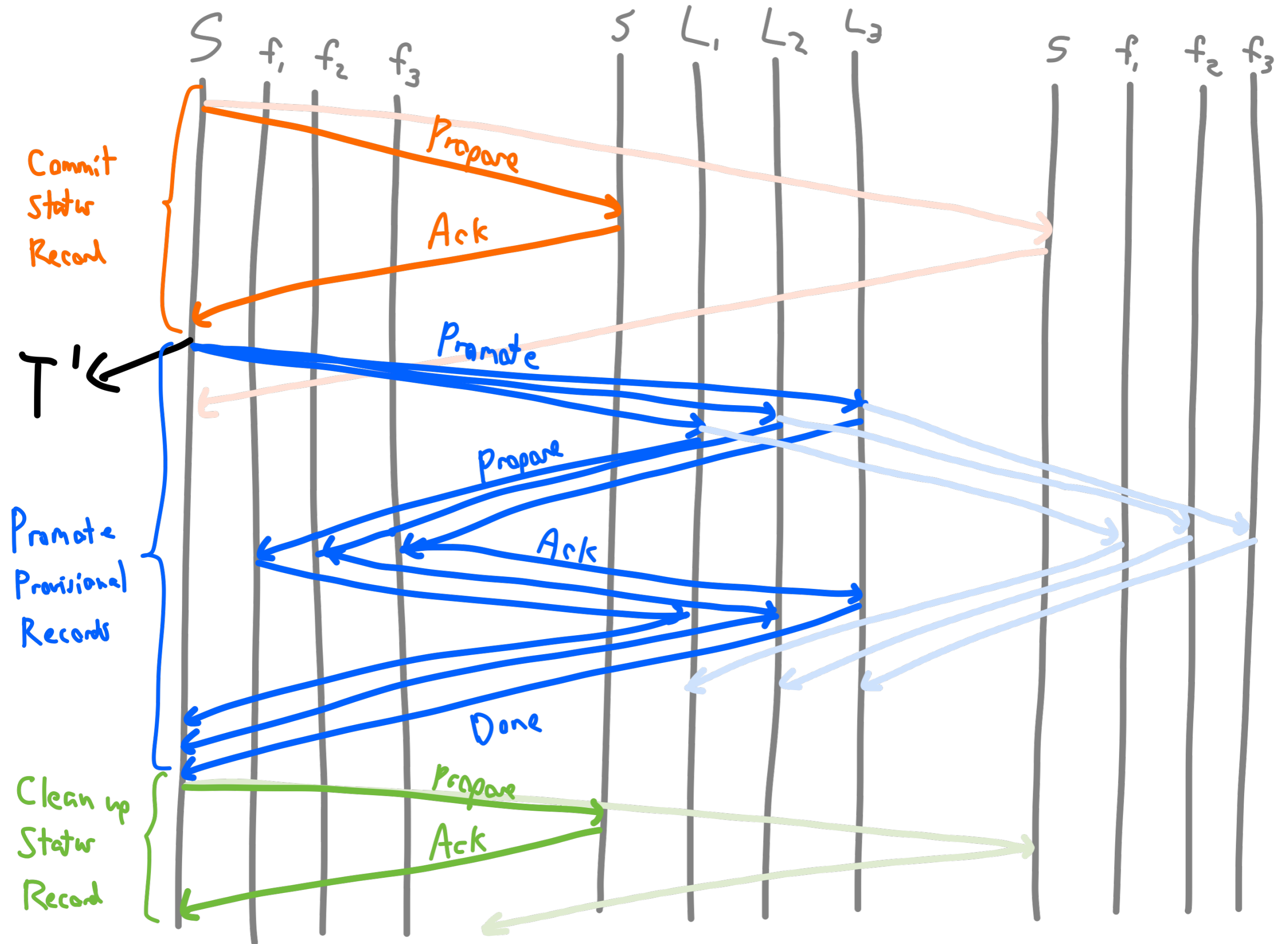
Yuggabyte DB

- Another Commercial/ASS impl
- Sort of Spanner - inspired
- HLCs coupled to Raft logs
- Leader leaves
- "Safe times"
- 2PC, txn record, async cleanup

- SI CQL
 - Serializable & SQL coming
- Linearizable single-key ops
- No general txns
 - One read
 - Batch write







R/W Transactions

2 + 8 cross-DC haps

2 + 2 local haps

Cockroach DB

- Commercial/OSS Spanner
- HLCs
- Serializable (lin. on keys?)
- Raft per shard
 - Leader leases

- No TrueTime

- Rely on your clocks!

- And your network

- 500 ms clock skew tolerance

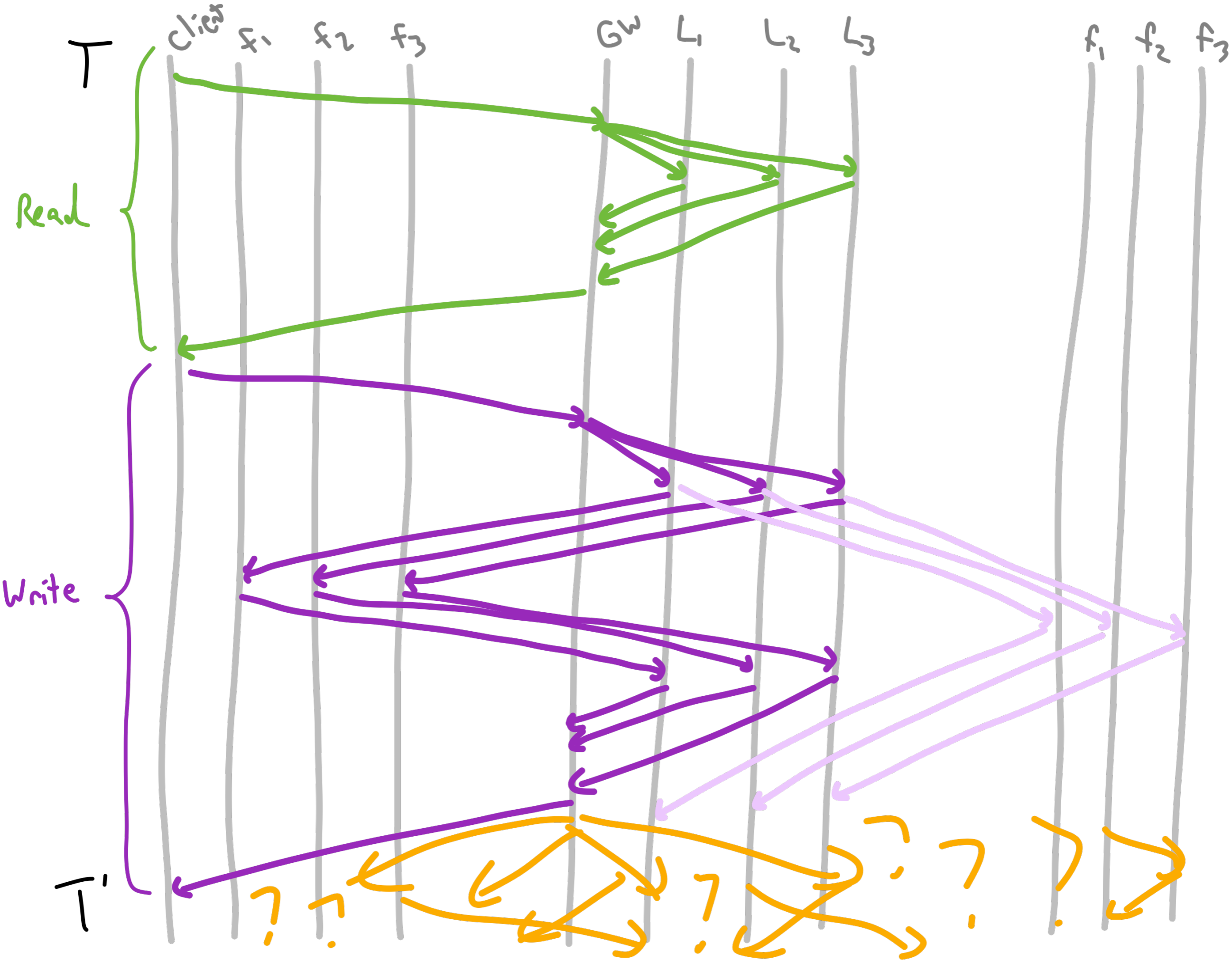
Spanner blocks \sim 7 ms

on every write

CockroachDB blocks contended

reads up to 500 ms

LPCC???



6 cross-DC heaps

4 local heaps

Calvin

Yale 2012

Ren, Shao, Abadi

Globally distributed

serializable transactions

Key Insights

- Locks/OCC \Rightarrow Predetermined order

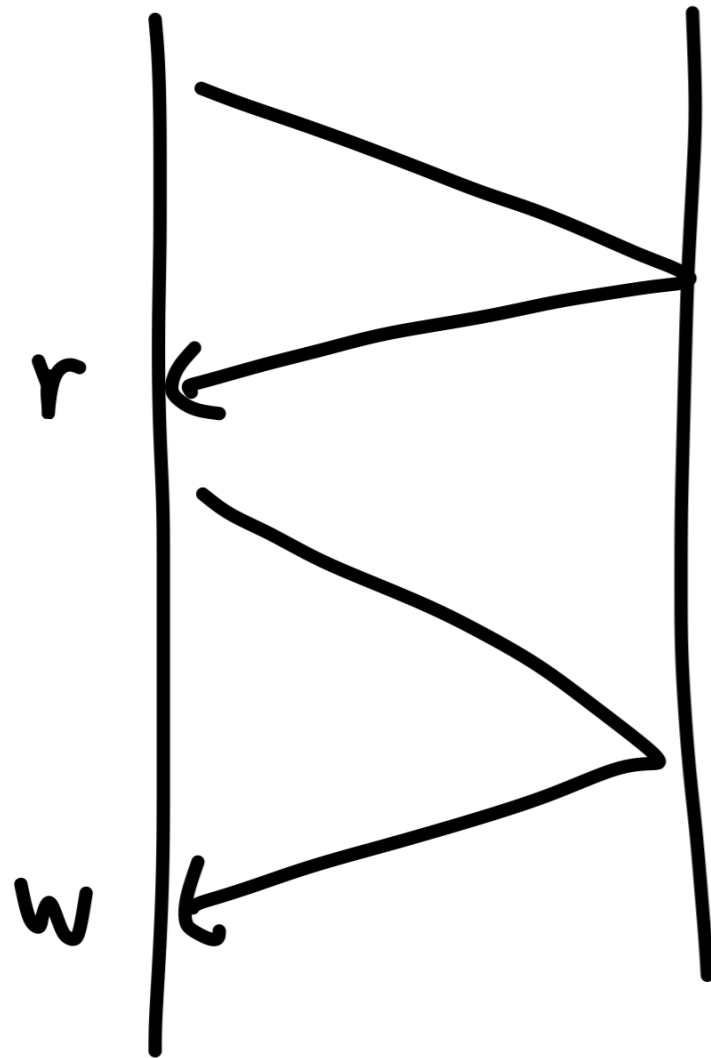
- Txns must be pure

- Txns known up front

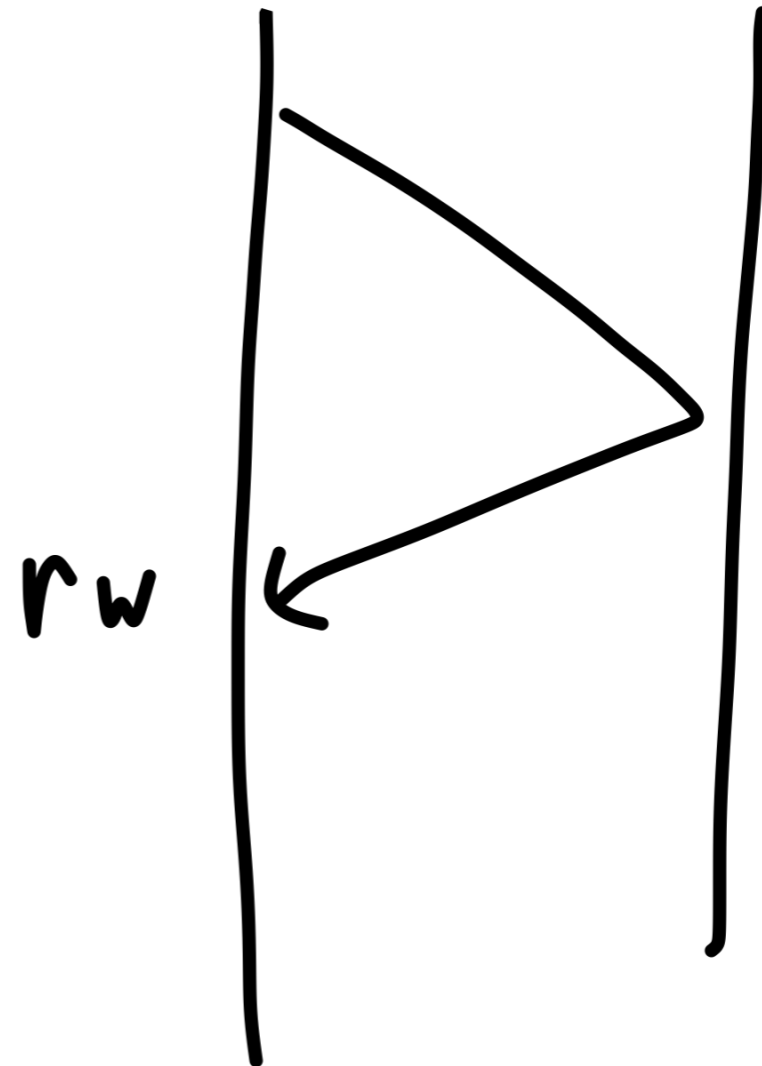
- No interactive sessions!

Key Insights

Client DB

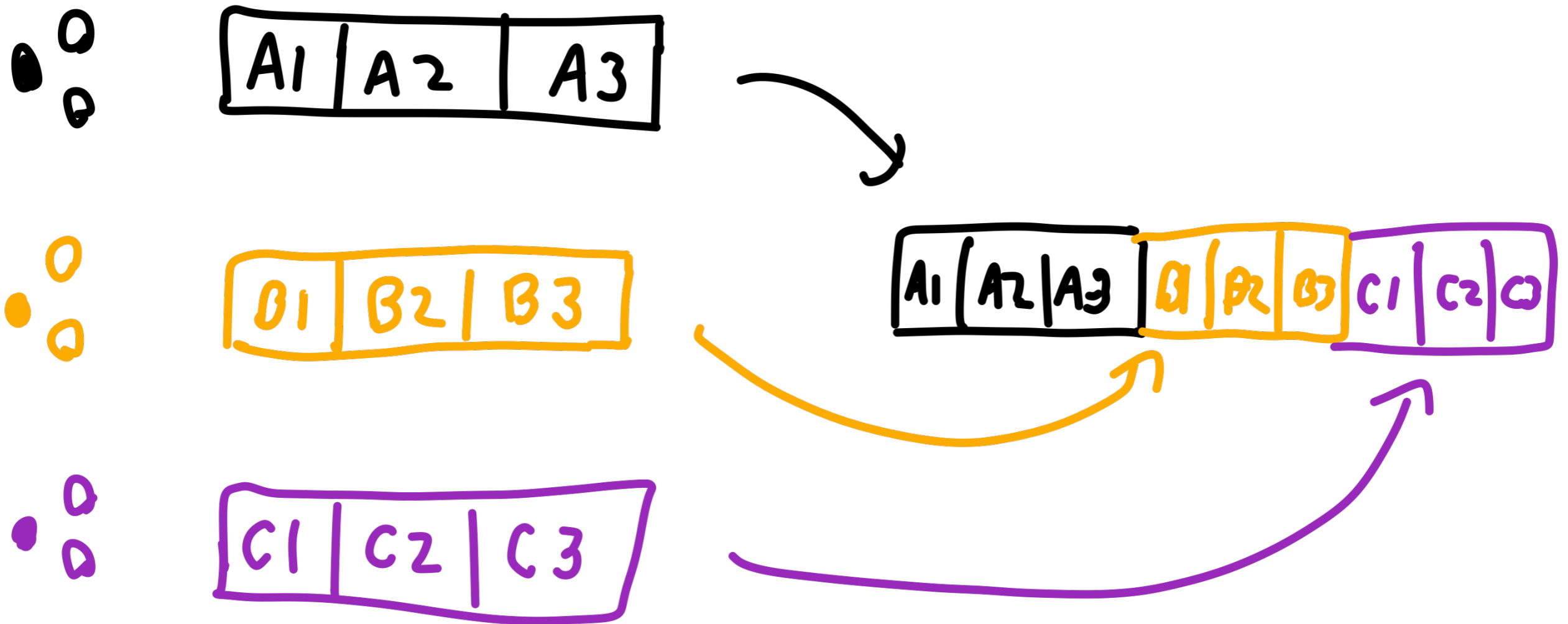


Client DB

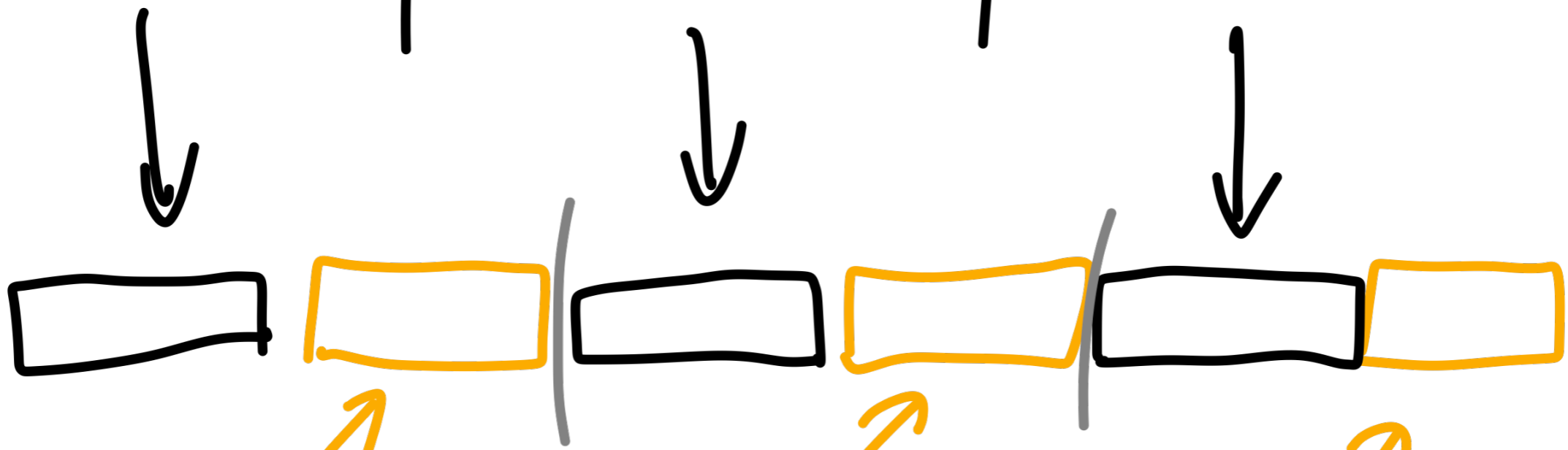


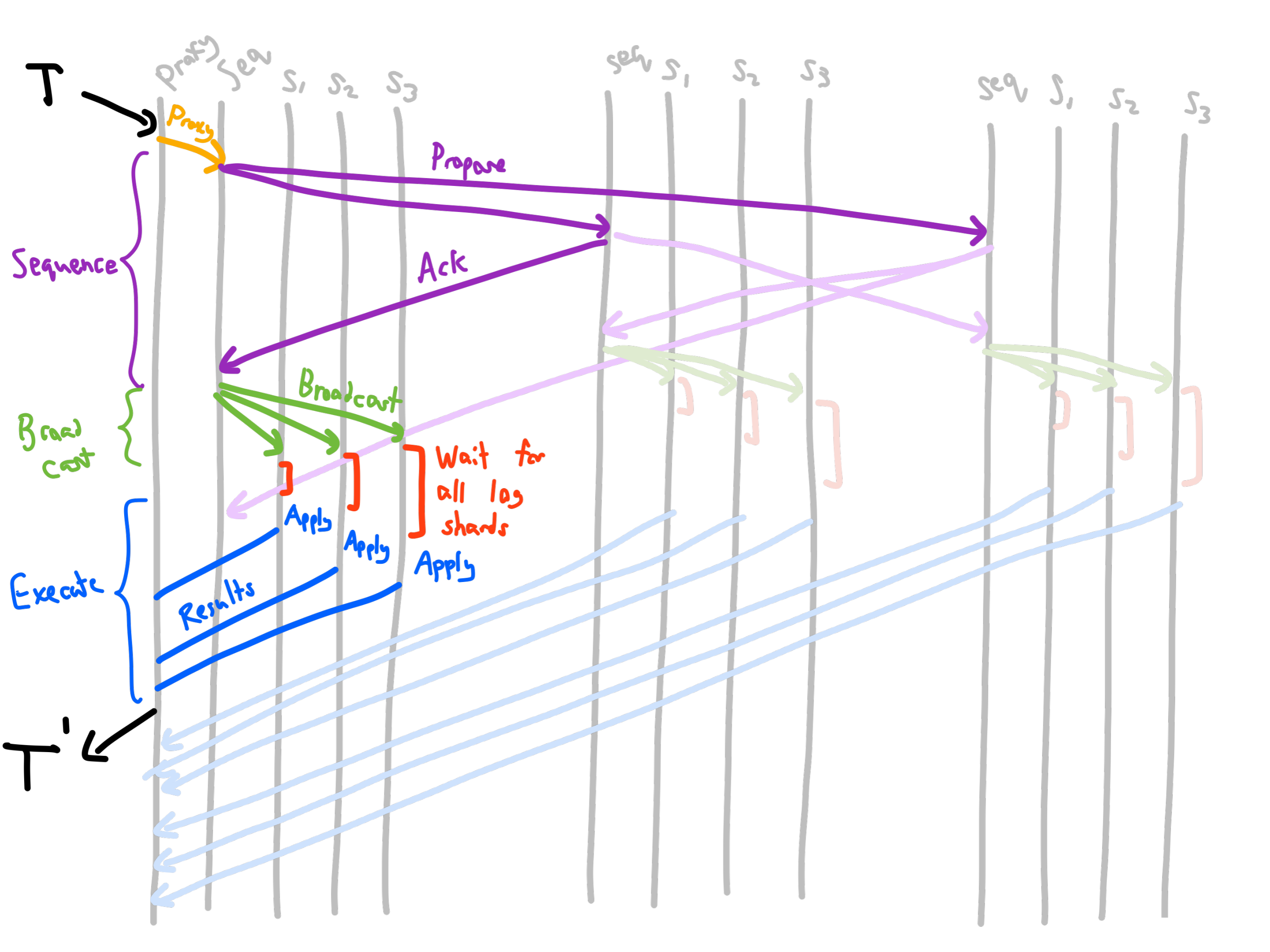
Key Insights

- You don't NEED Paxos for order...



epoch





2

Cross-DC

hops

5

local

hops

No clock dependency
(for safety)

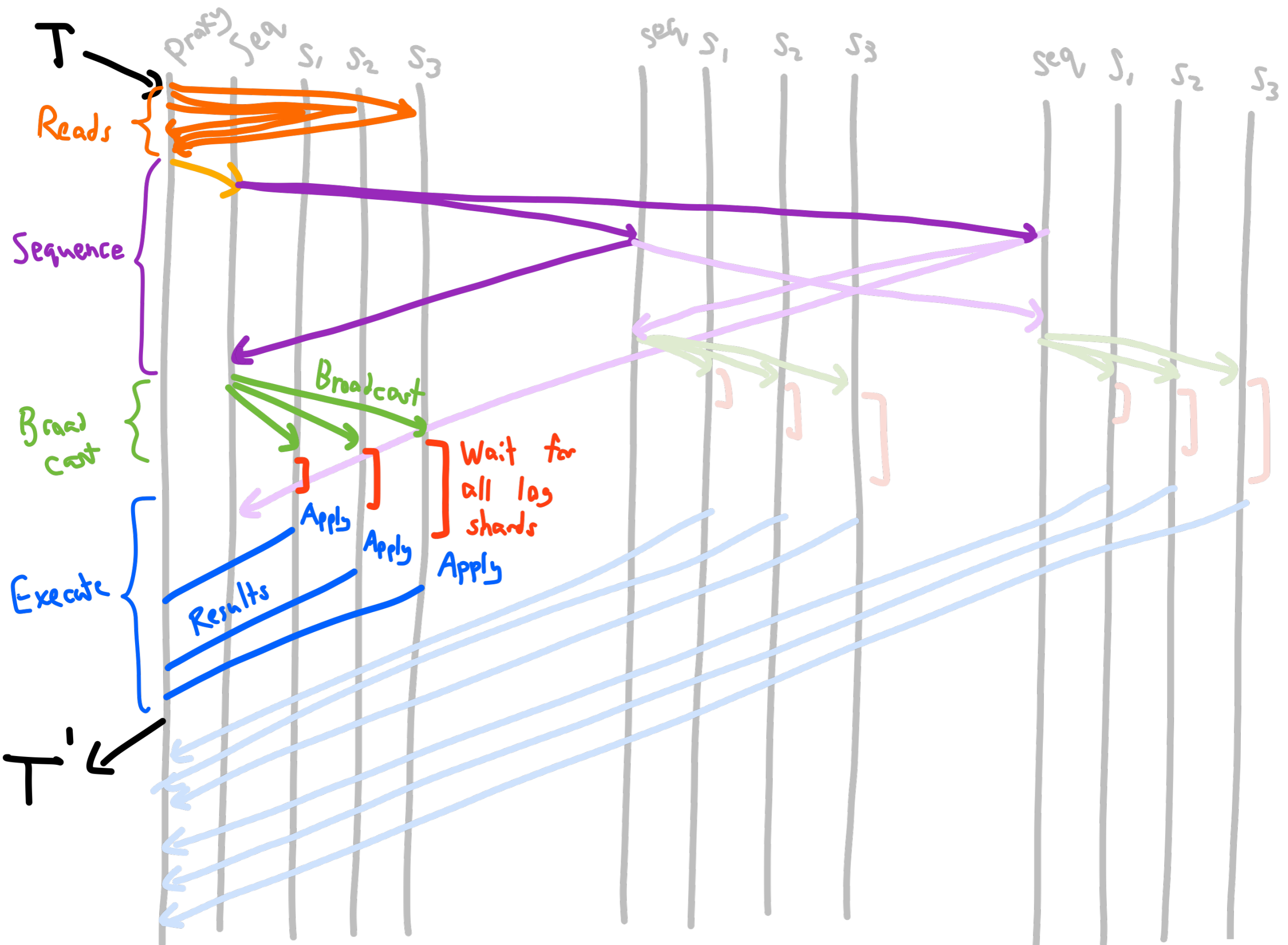
Need to block for epoch
to complete

No Interactive Transactions

- Some queries inexpressible
- Back off to $r + CAS$
- Can extend Calvin w/extra round trips to support interactive transactions?

Fauna DB

- Commercial impl of Calvin
- Raft for consensus
- SI, serializable, Strict-1SR
- Dynamic R/W sets



2 cross-DC hops

7 local hops

to
Review

Latency

Cross-DC

Local

Percolator	8-14?	10-0?
T:DB	14	0
Spanner	8	1
YugabyteDB	2+8	2+2
Cockroach DB	6	4
Calvin	2	5
Fauna DB	2	7

Linearizable	Read	Cast (xDC)
--------------	------	------------

Percolator

2?

TiDB

2?

Spanner

0

Yugabyte DB

2?

Cockroach DB

N/A?

Calvin

2

Clack Errors

- TiDB: Possibly? Not observed
- Spanner: State reads, ???
- Cockroach: Caused reverse, ???
- YB: State reads, ???
- Calvin: No safety issues

Interactive Transactions

- ✓ Percolator
- ✓ TiDB
- ✓ Spanner
- ✓ Cockroach DB
- ✗ Yugabyte DB
- ✗ Calvin
- ✗ Fauna DB

Isolation

Min

Max

Percolator

SI

SI

TiDB

SI

SI

Spanner

SI (r)

Strict-2SR

Cockroach

SI (r)

Serializable

YugaByte

SI

SI

FawnDB

SI

Strict-2SR

Thanks

- Fauna

\$

- Cockroach

\$

- MongoDB

\$

- Yugabyte

\$

- Volt DB

\$

- PingCAP

\$

Thanks

- Camille Fournier
- Kena (@kenz)
- Tang Liu
- Nathan Van Benschoten
- Karthik Ranganathan

Thanks

- @sergei
- Matt Freels
- Evan Weaver

Questions???