

J E P S E N

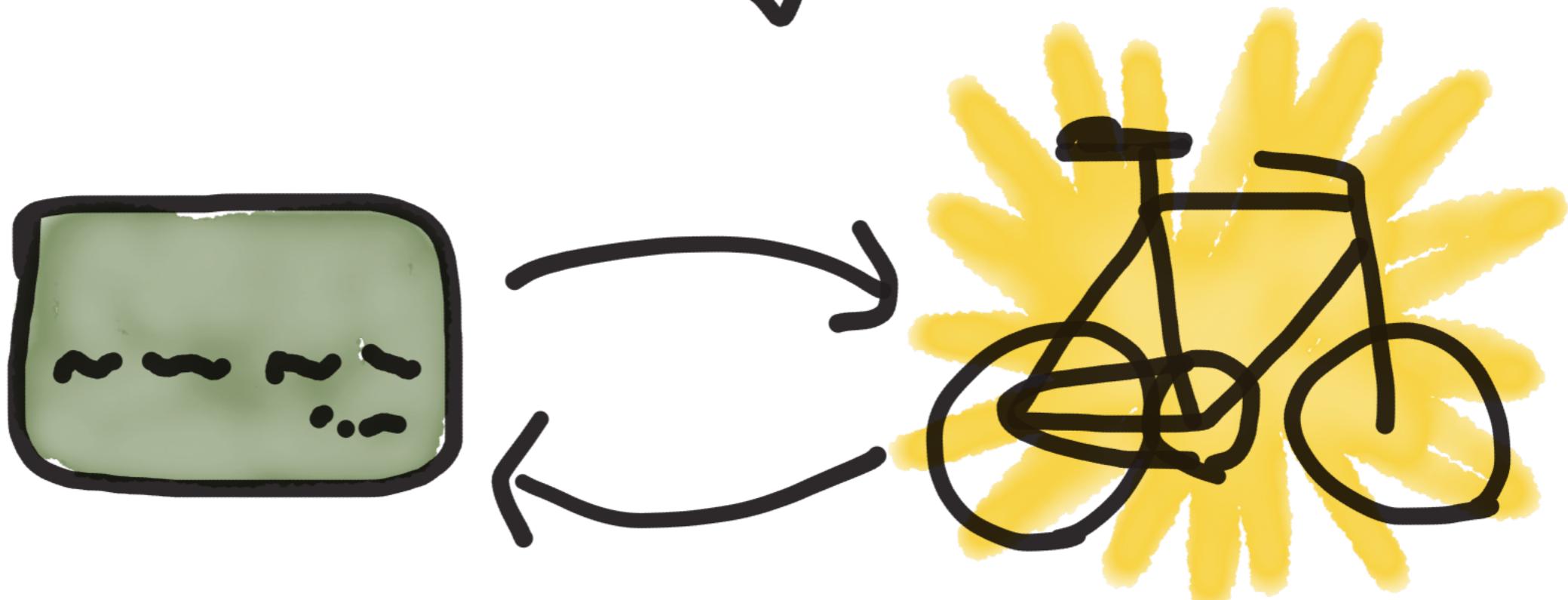
IV

Hope Springs
Eternal



@aphyr
(Kyle Kingsbury)

straße



Public
API {



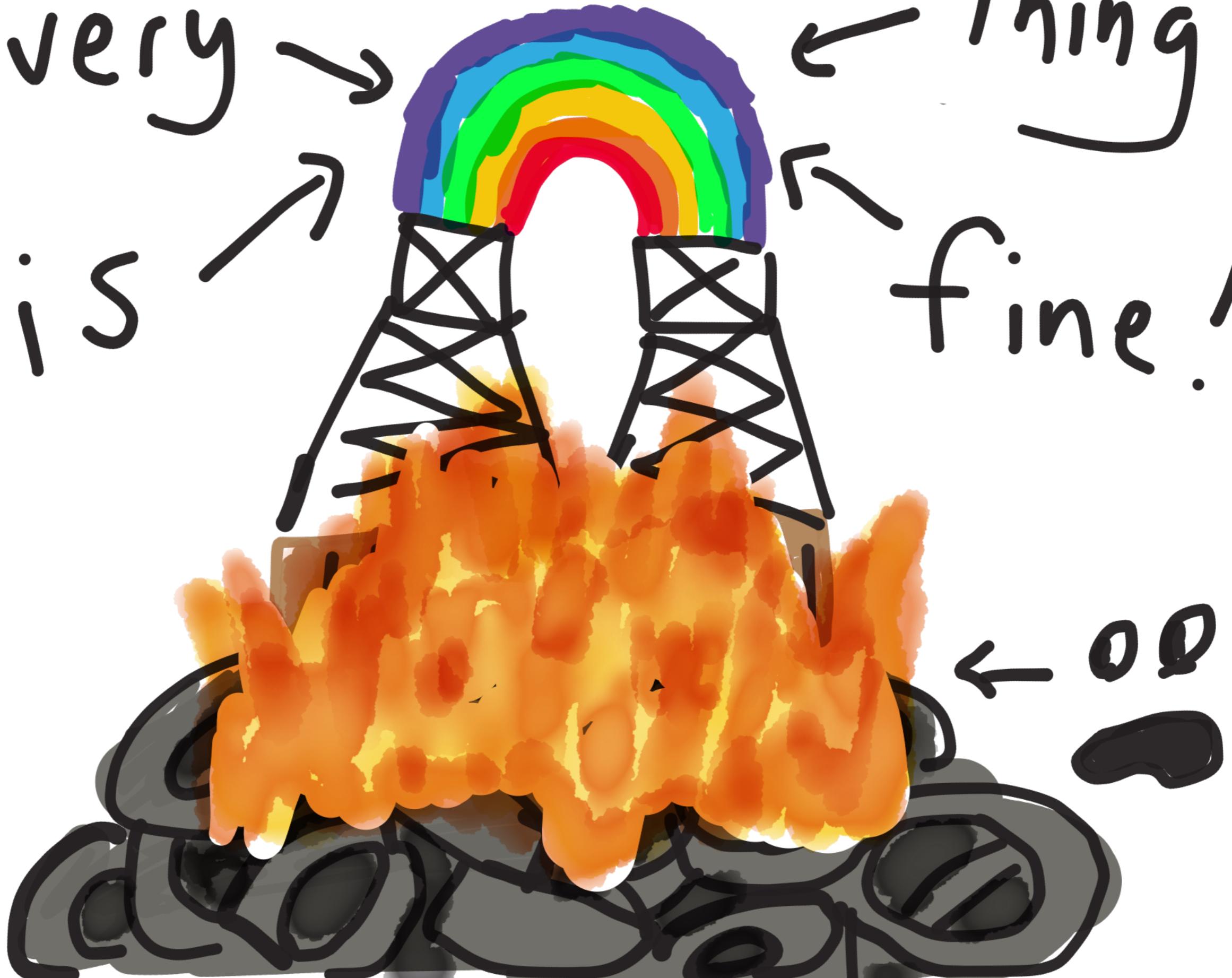
} API code

Ruby {

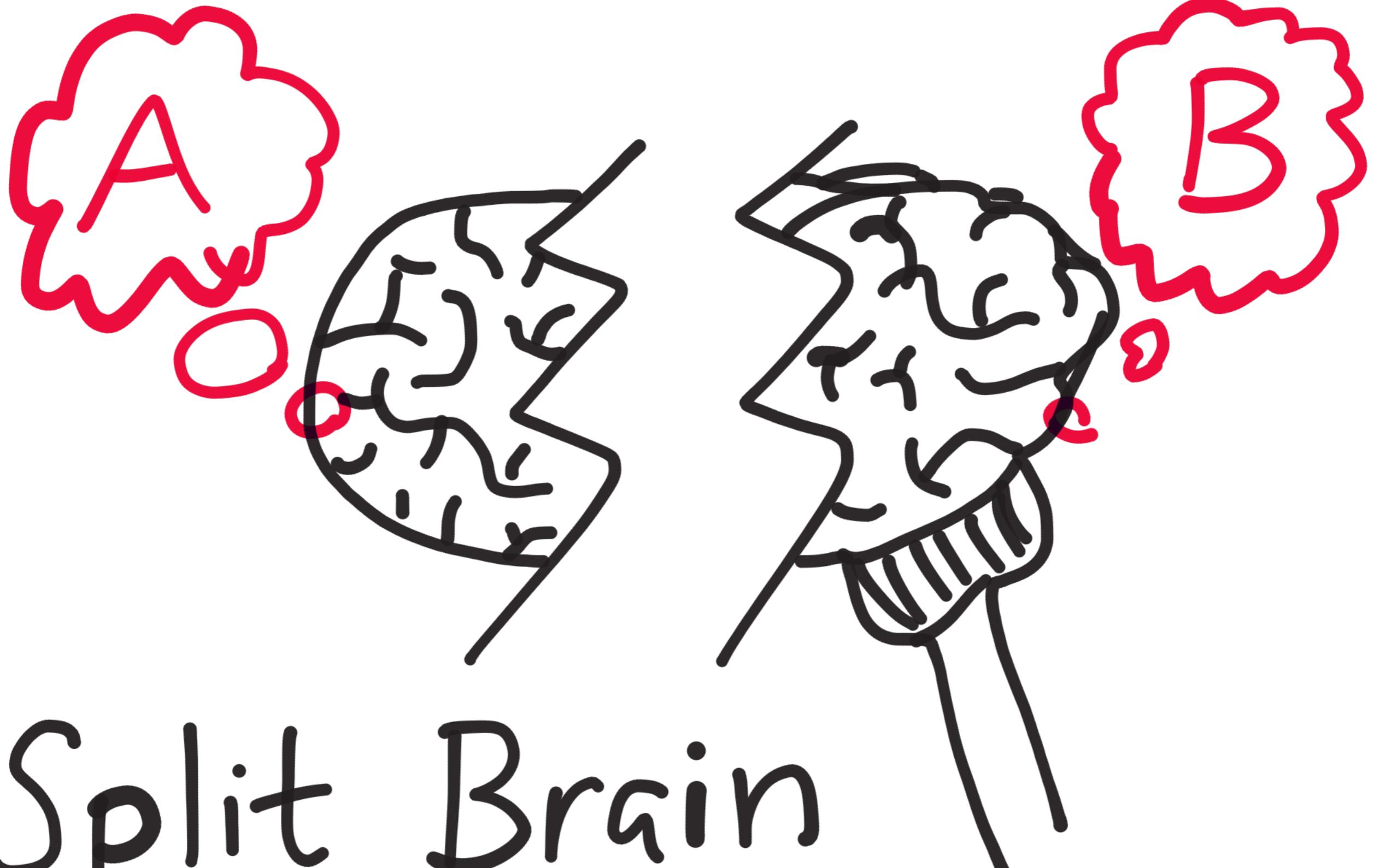


Every →
is →

← Thing
← fine !

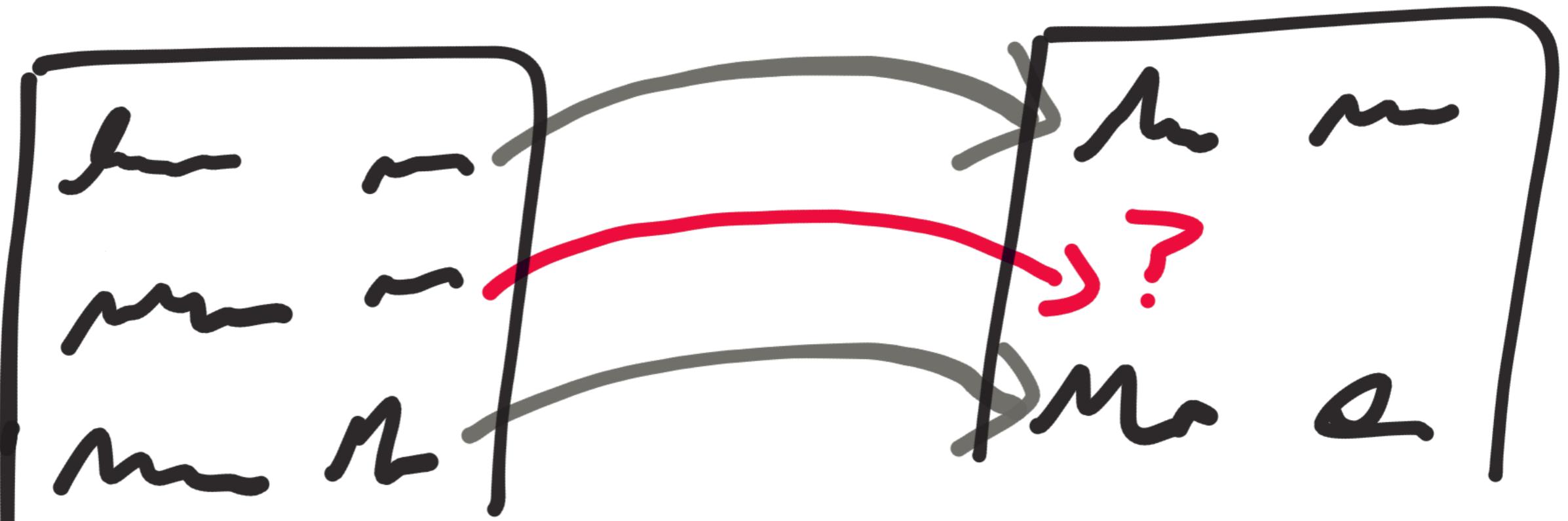


Databases! /
Queues! /
Discovery! /
THE HORROR

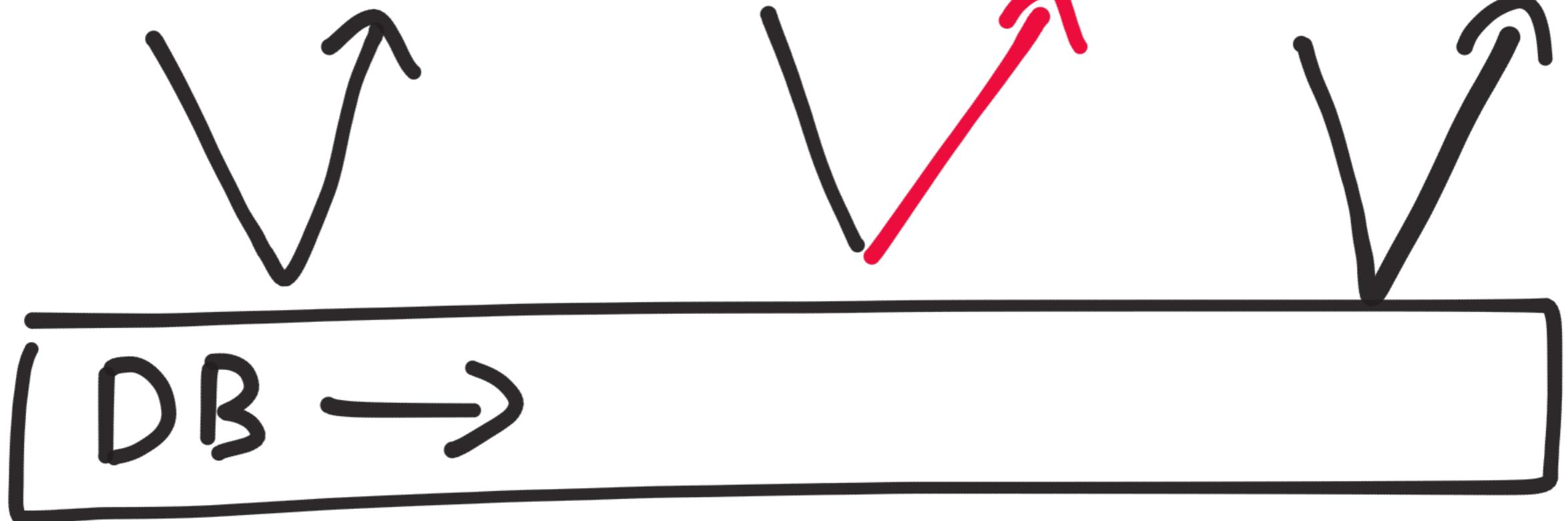


Split Brain

Broken Foreign Keys



write ok read \emptyset read ok



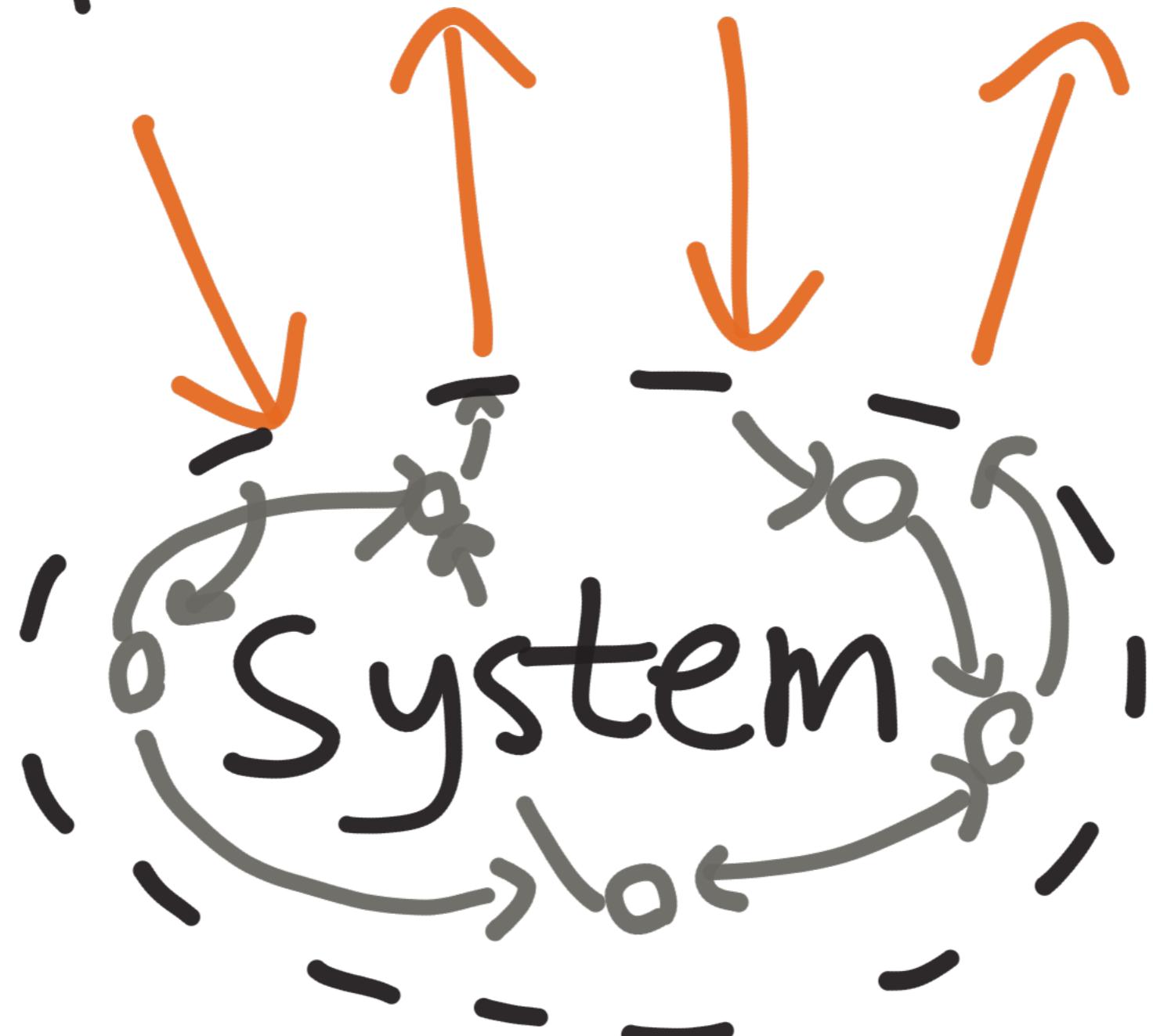
Anomalies

How do you
know if a
system is safe?

Jepsen

github.com/aphyr/jepsen

Environment





— INVARIANTS —

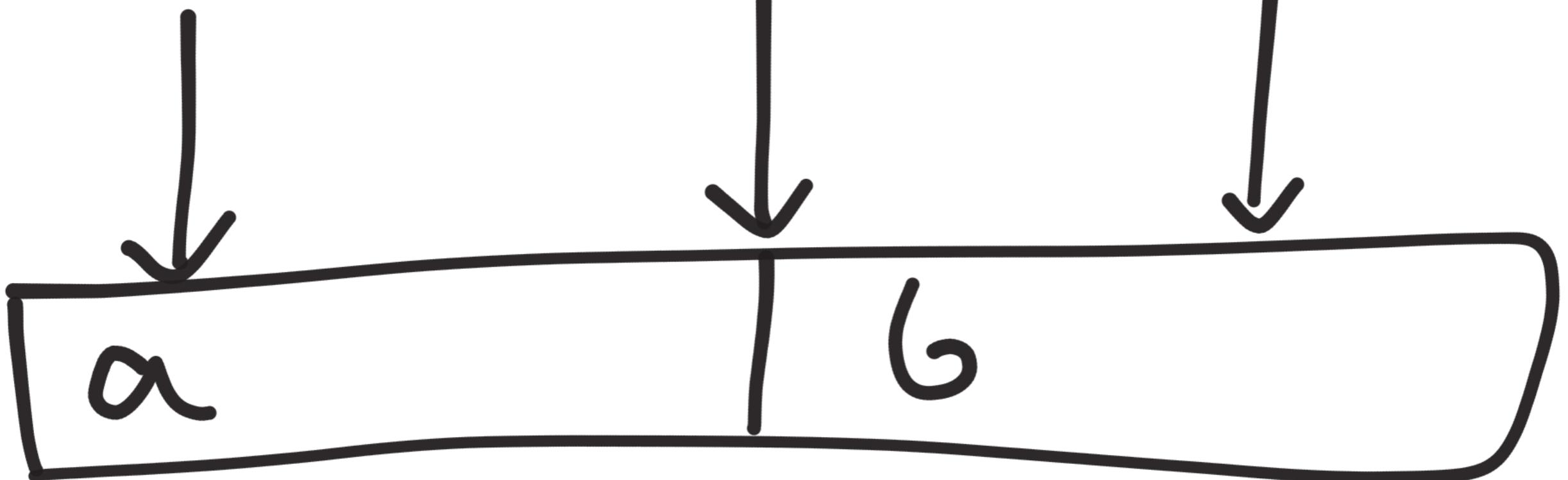


var $X = \alpha$

print X

$X = b$

print X

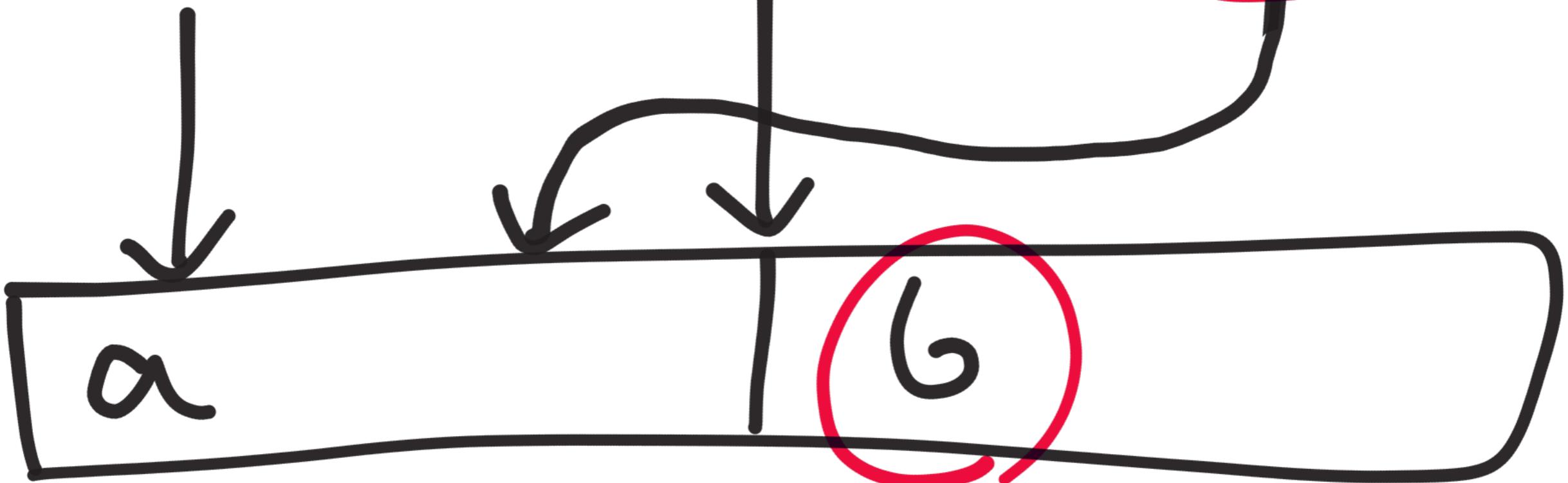
$r(a)$ $w(b)$ $r(b)$ 

time →

$r(a)$

$w(g)$

$r(a)$



time →

Invariants

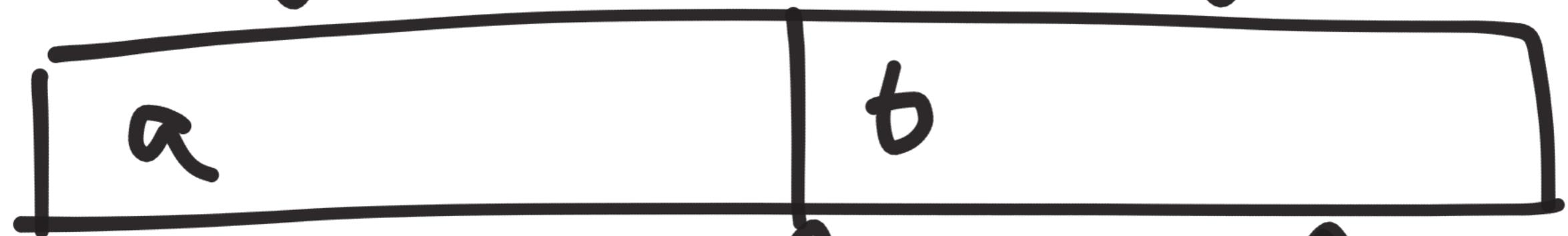
Constraint

Histories

What about

concurrency?

$r(a)$



$r(b)$



$w(b)$



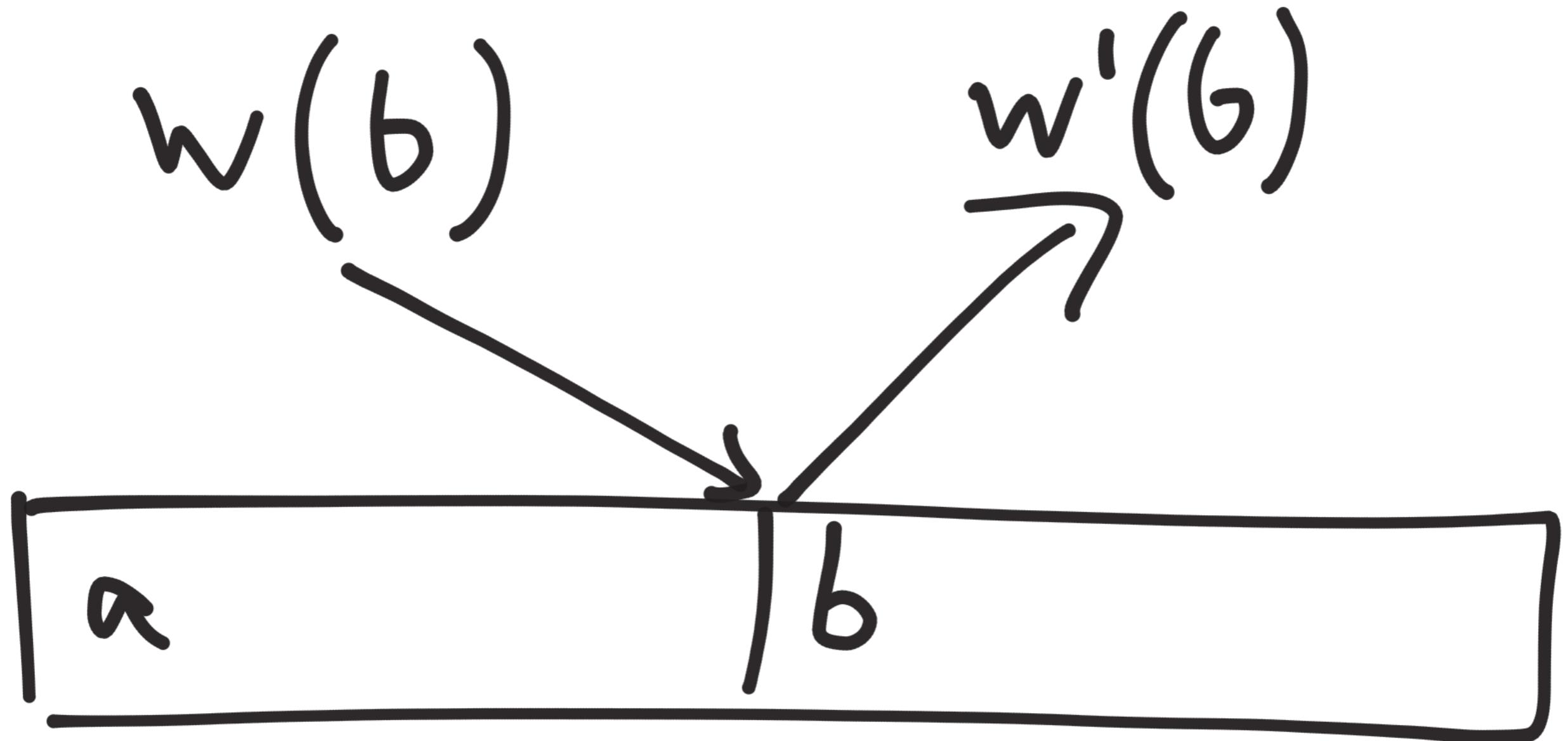
$r(b)$



Distributed

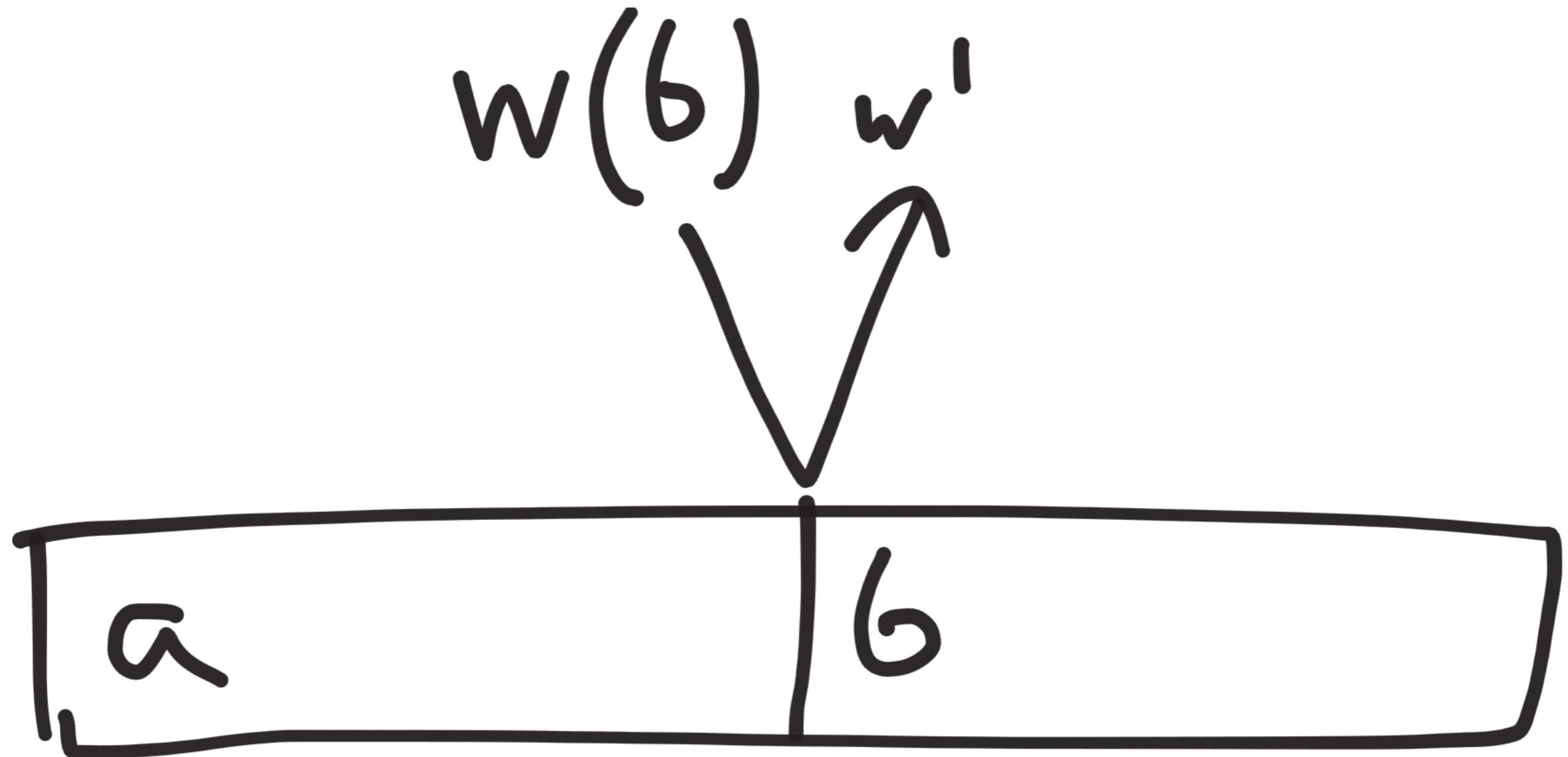


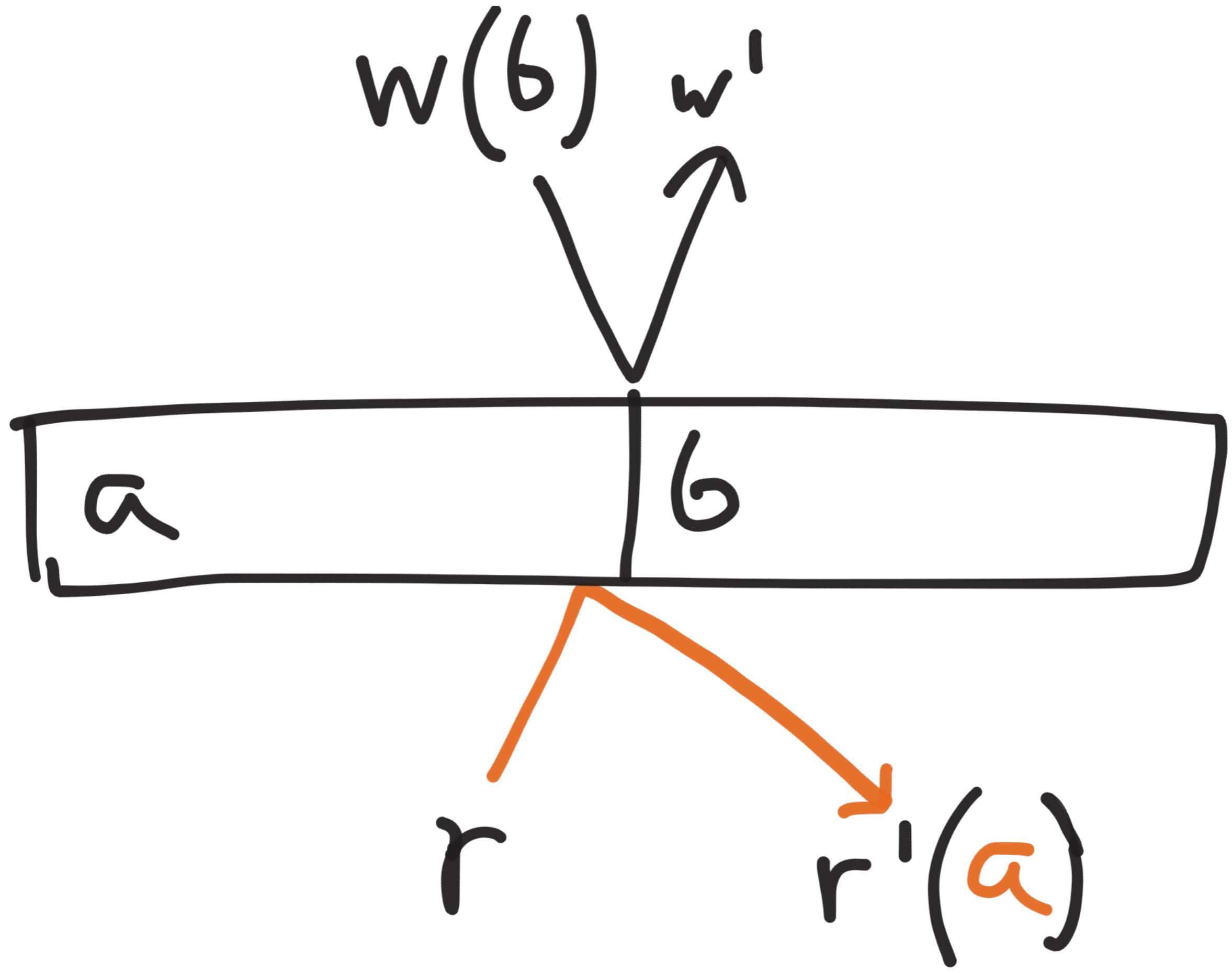
Distance



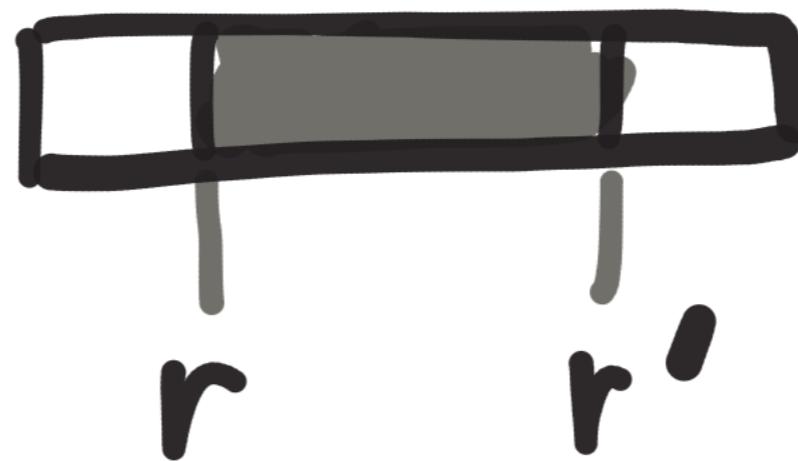
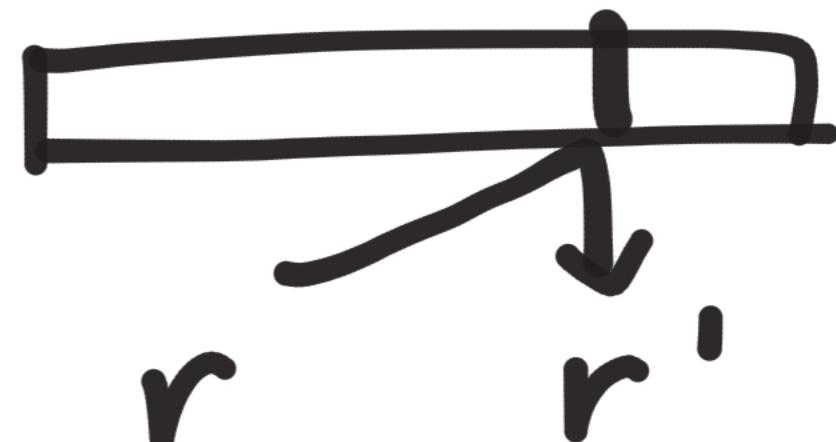
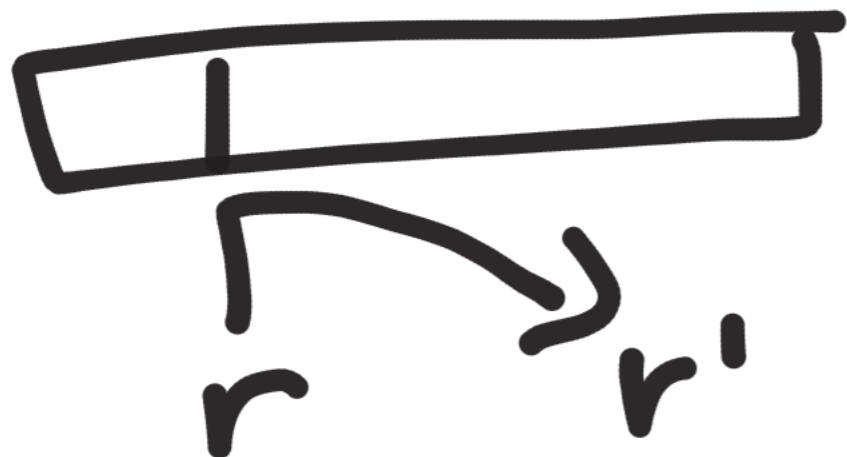
Mess@ger take time

Delays imply
ambiguous orders





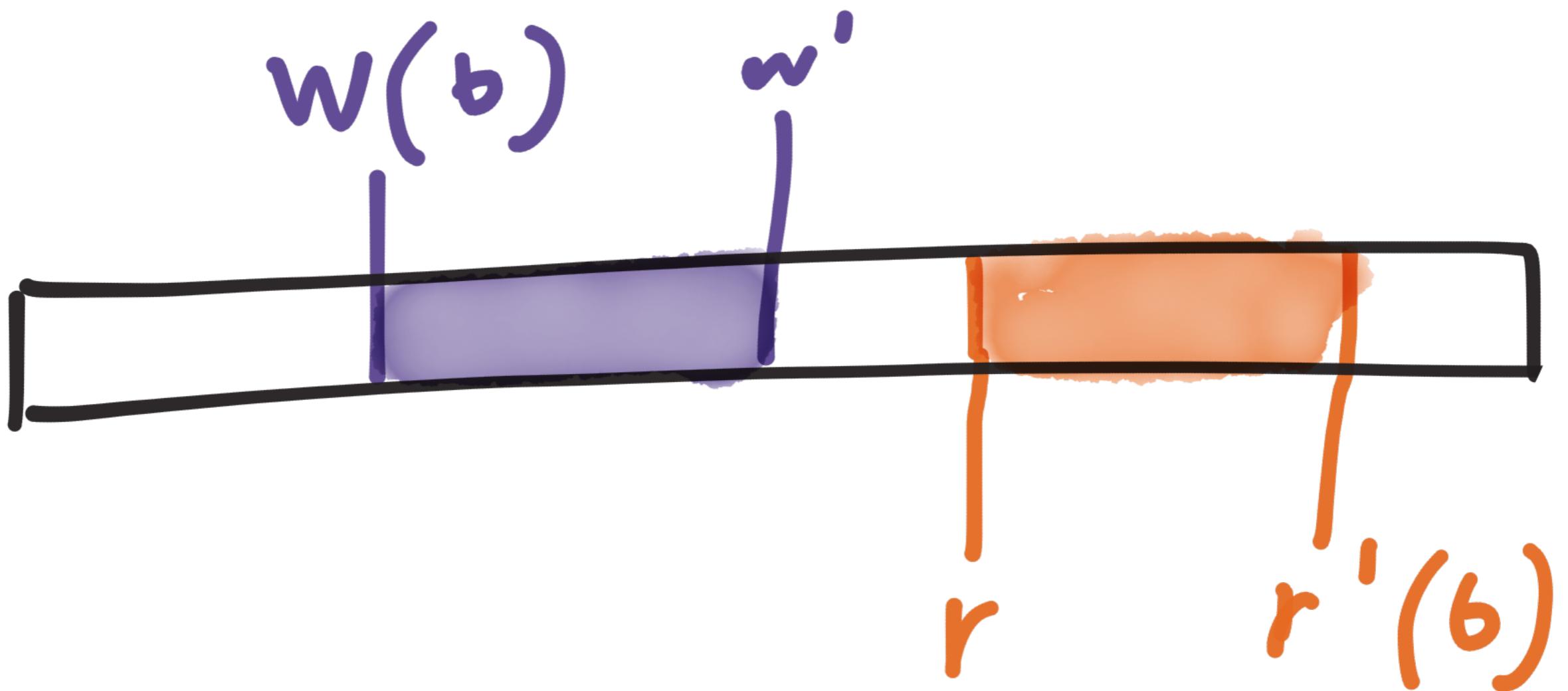
But there are finite
bounds on ambiguity!



Linearizability



After op is complete,
guaranteed visibility!

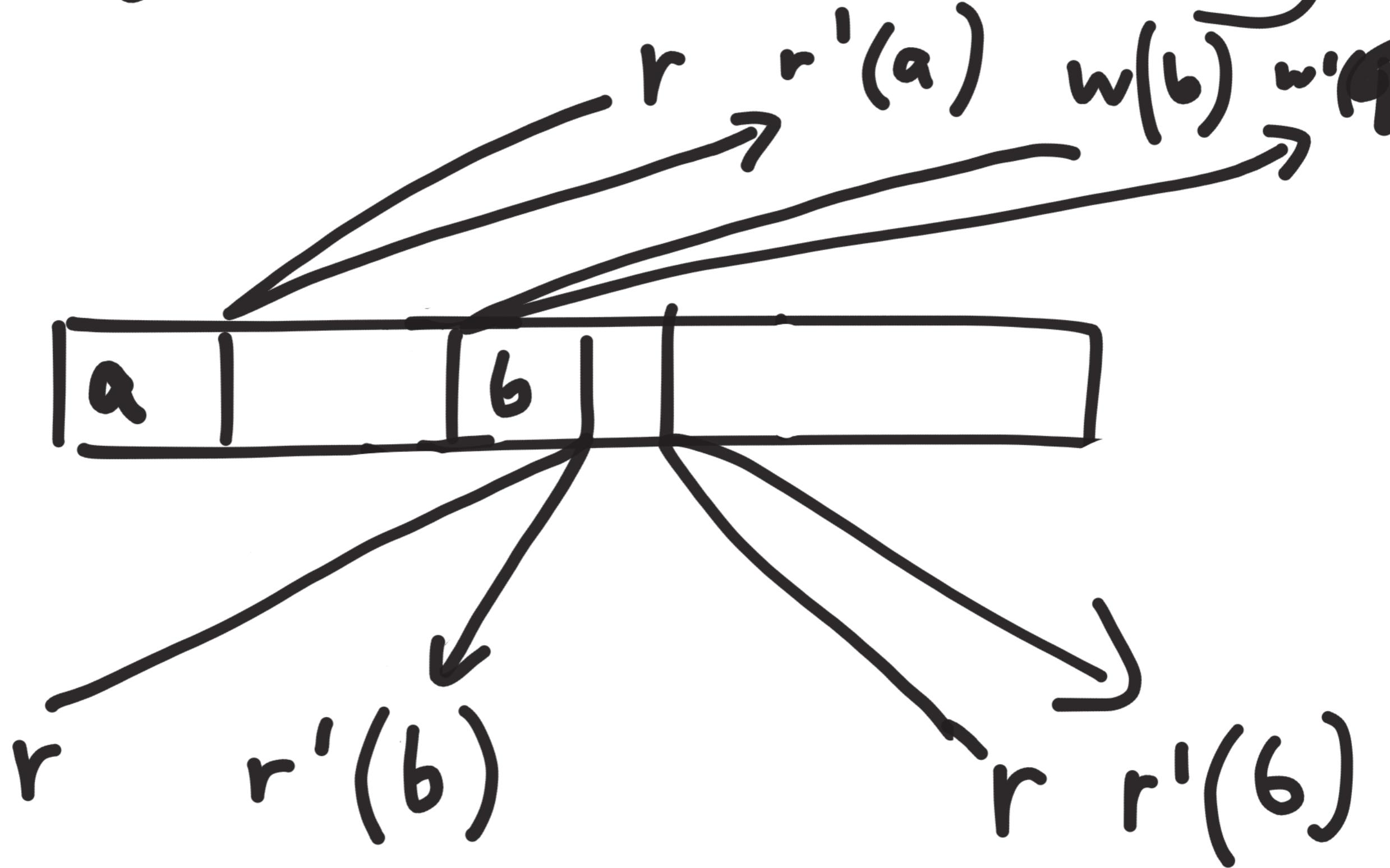


- Everyone agrees on order
- Ops are guaranteed to be visible to ALL clients once complete
- No stale reads

What if...
we relaxed the

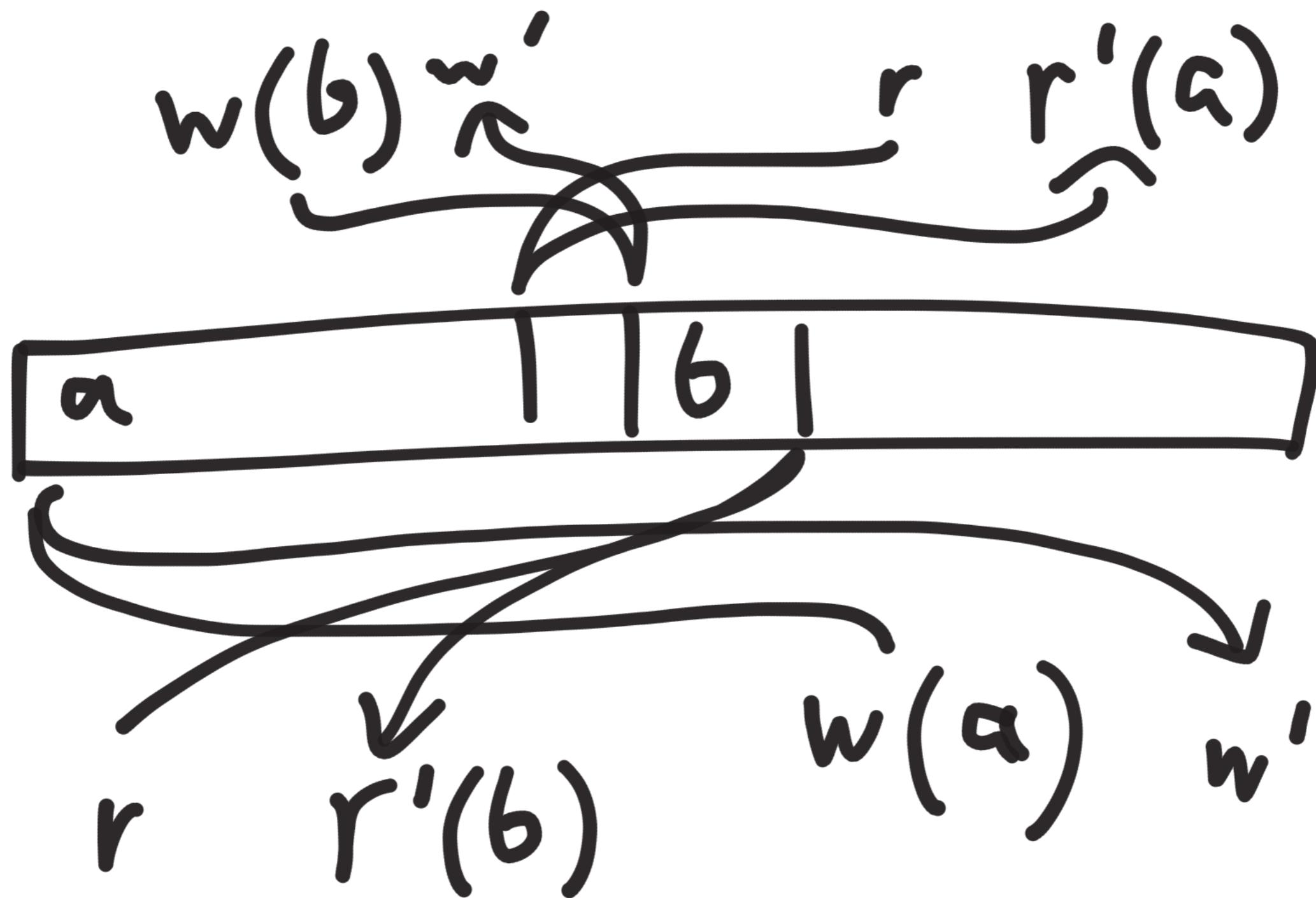
invariants?

Sequential Consistency



All processes agree
on Order of ops, but
may be delayed in
time.

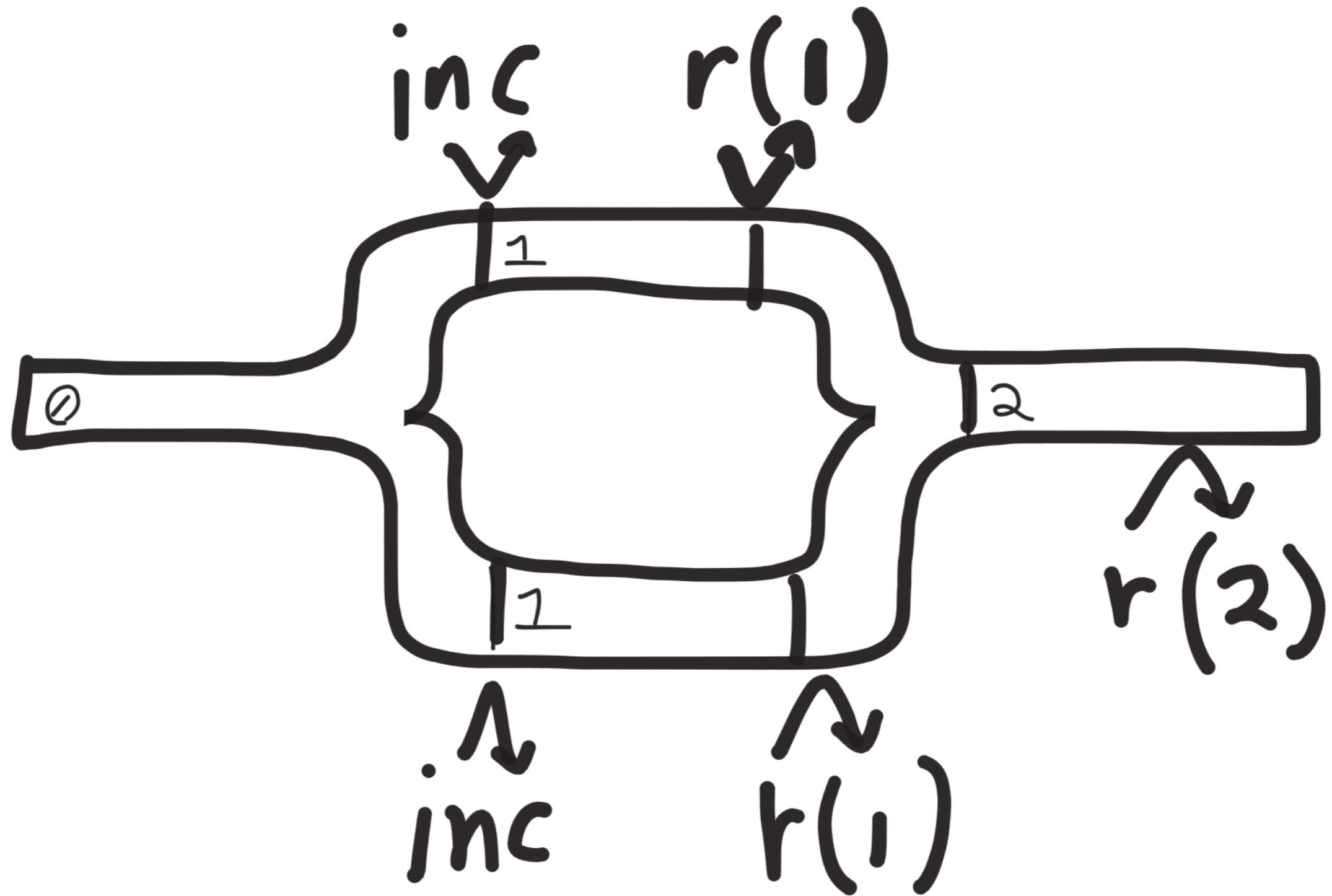
Serializable Consistency



Ops happen in SOME

Order, but nobody
has to agree on what
that was.

Eventual Consistency

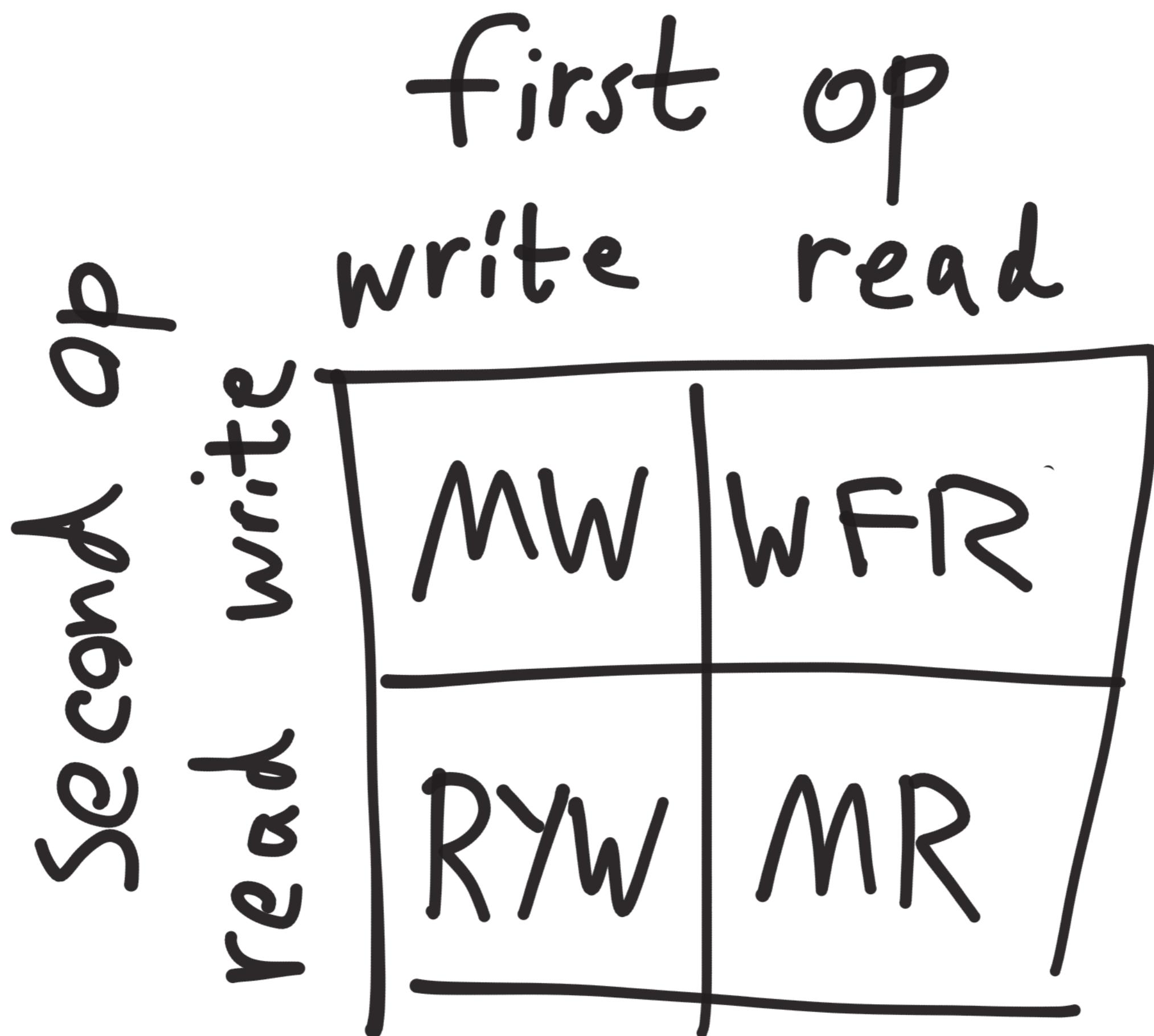


There doesn't have

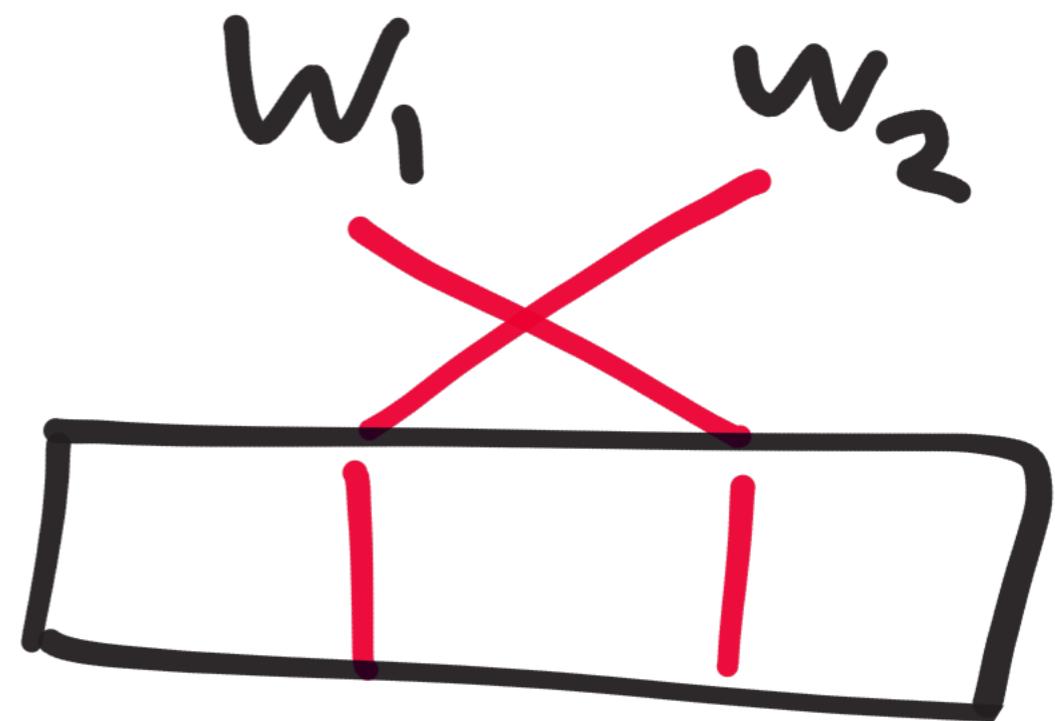
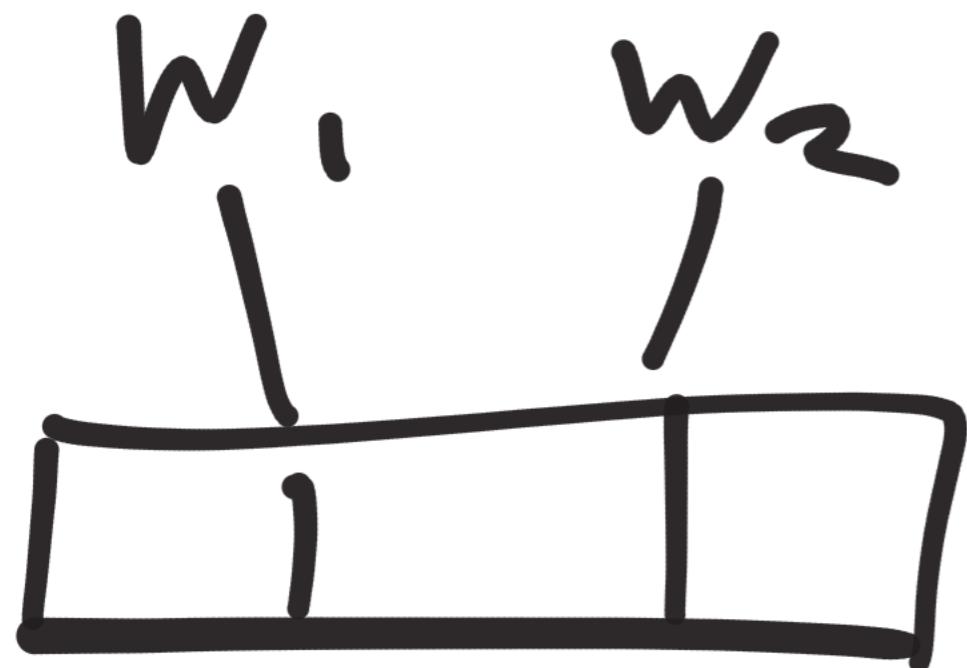
to be any single

order so long as

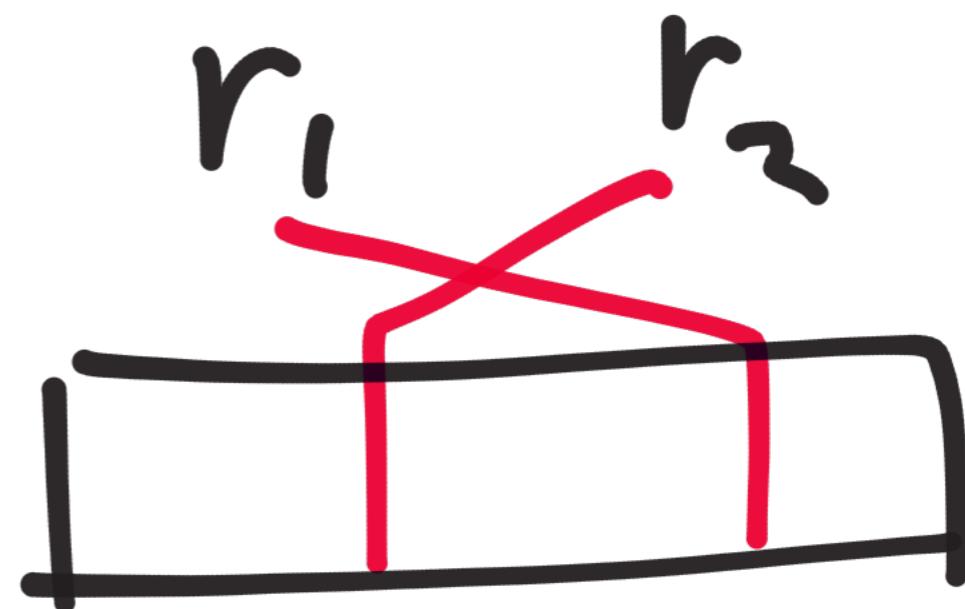
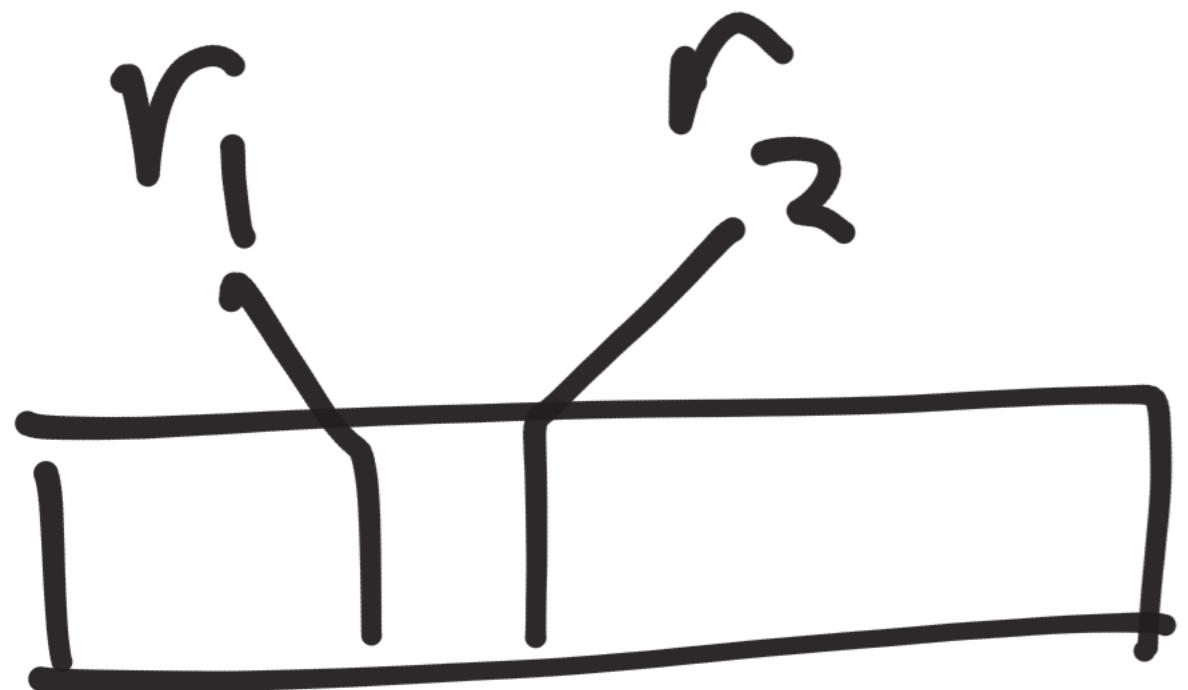
histories converge



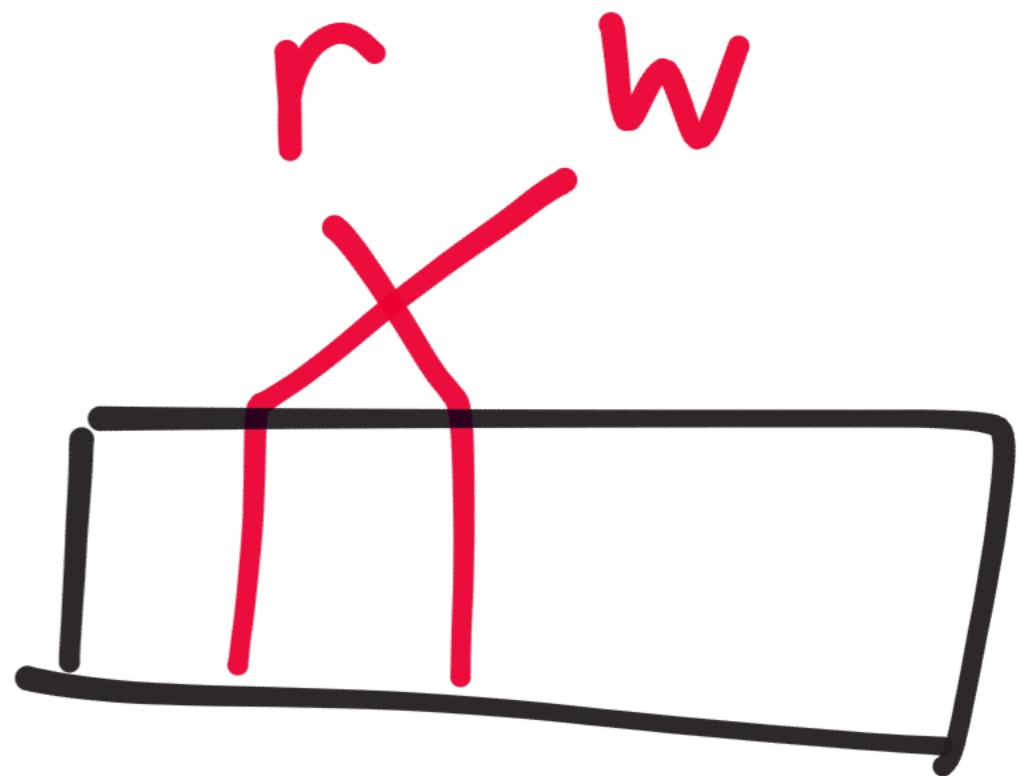
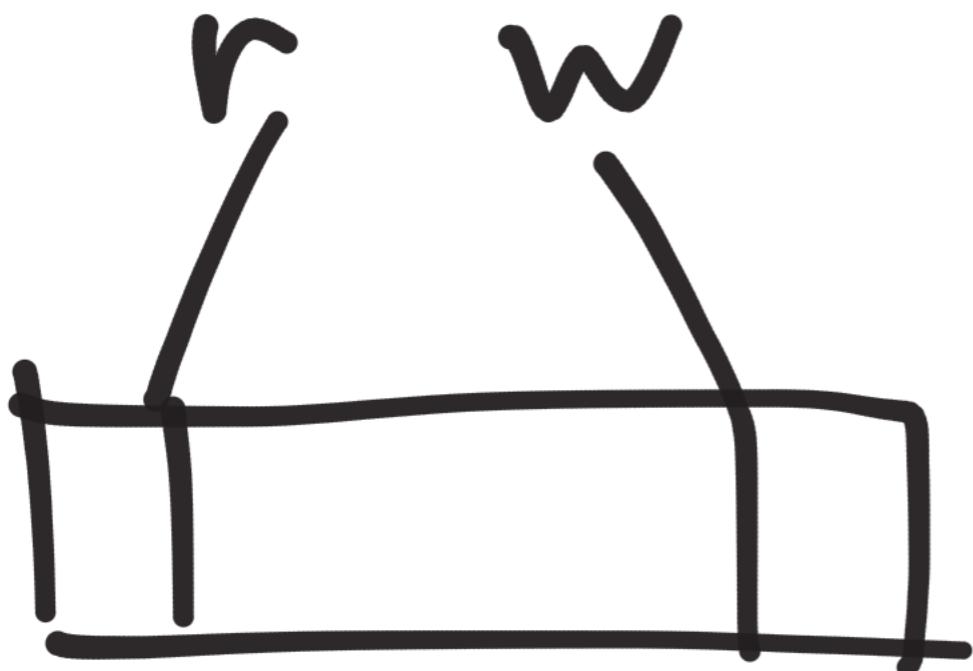
MW: Monotonic Write



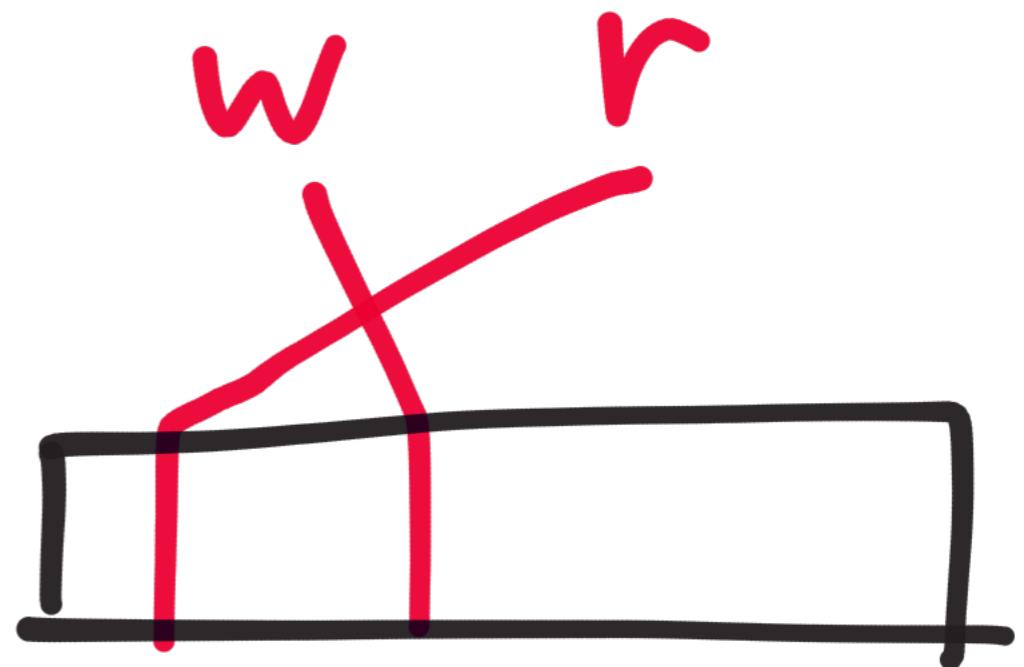
MR : Monotonic Read



WFR: Writes Follow Reads

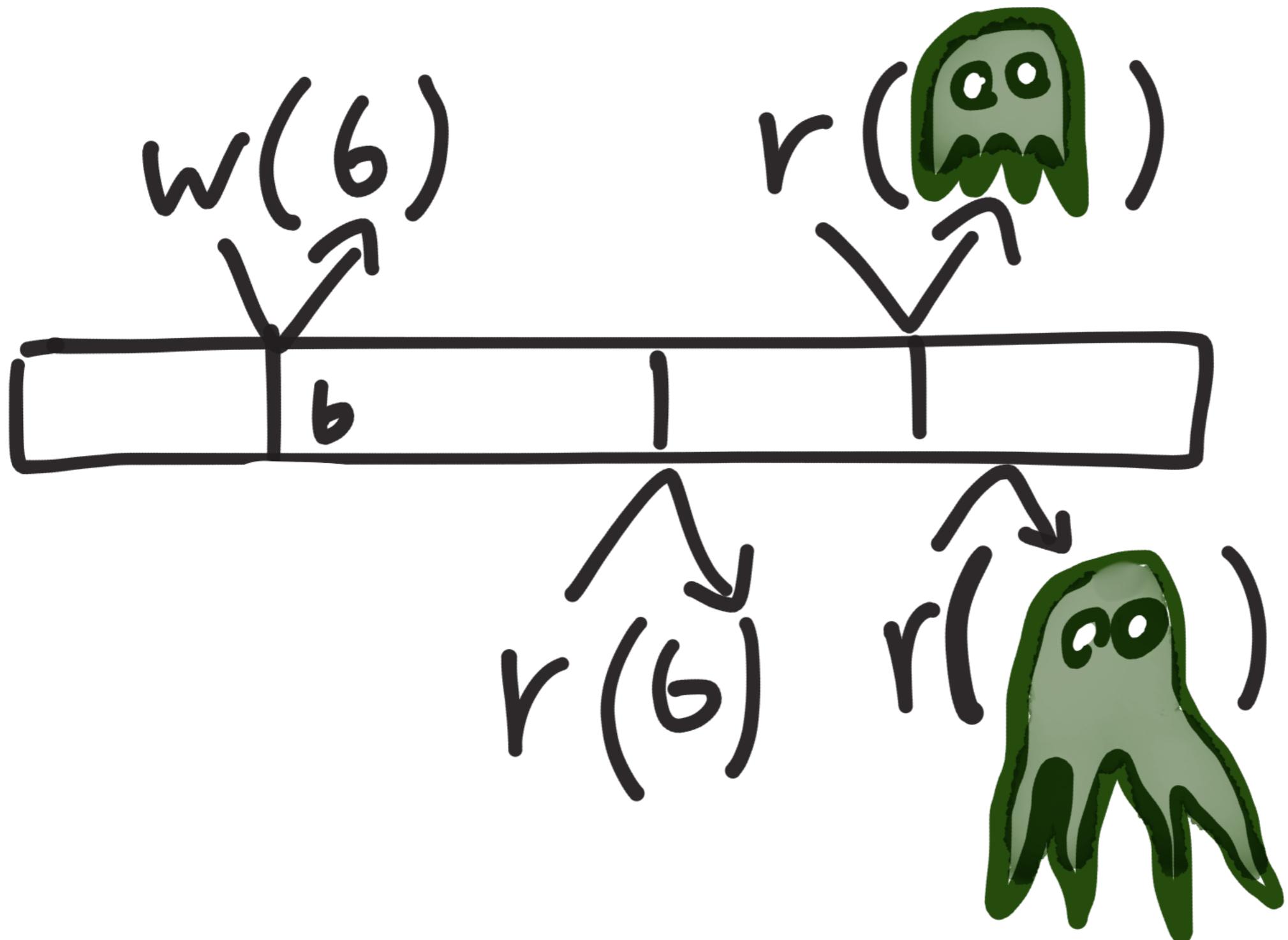


RYW: Read Your Writes



Infinitely many
Consistency classes

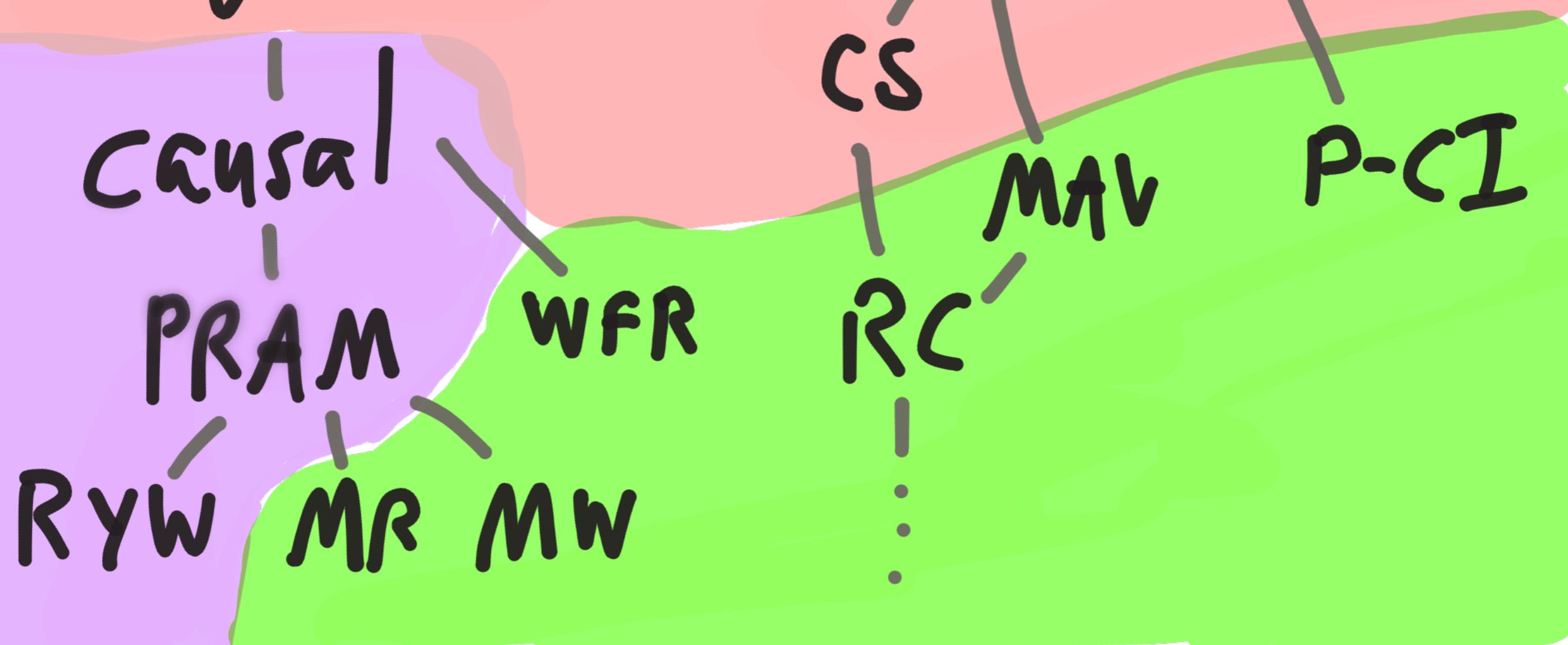
Cthulhu Consistency



strong serializable

linearizable

sequential



"weaker" consistency
models are more
available in failure

"Weaker" models

are less

intuitive

"weaker" models

need less coordination,

so they're faster!

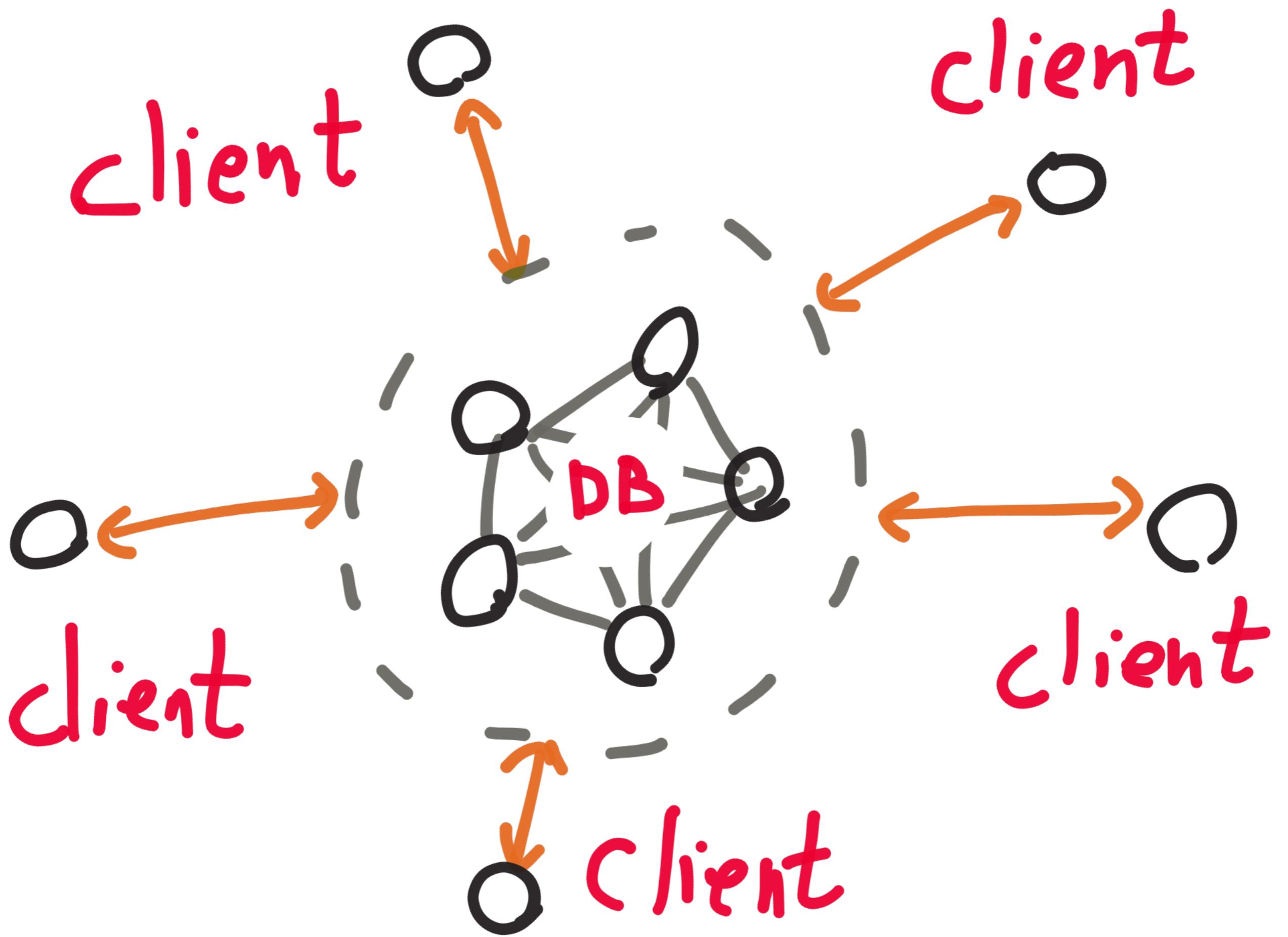
→ CPU memory model

"Weak" \neq unsafe

- Diff algos require
diff. levels of
consistency

SO: Consistency

models constrain
the system-environment
interactions.



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

The diagram illustrates three client requests to a central database system. Each request is shown as a horizontal line segment with arrows at both ends, indicating a transaction or query. The segments are labeled $w-w'$, $r-r'$, and $w-w'$ from left to right. All three segments converge towards a central, blurred, cloud-like shape labeled "DB", representing the database.

DB :

client : w — w' w — ...

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



invoke

ok

invoke

fail

invoke ?

? info

?

?

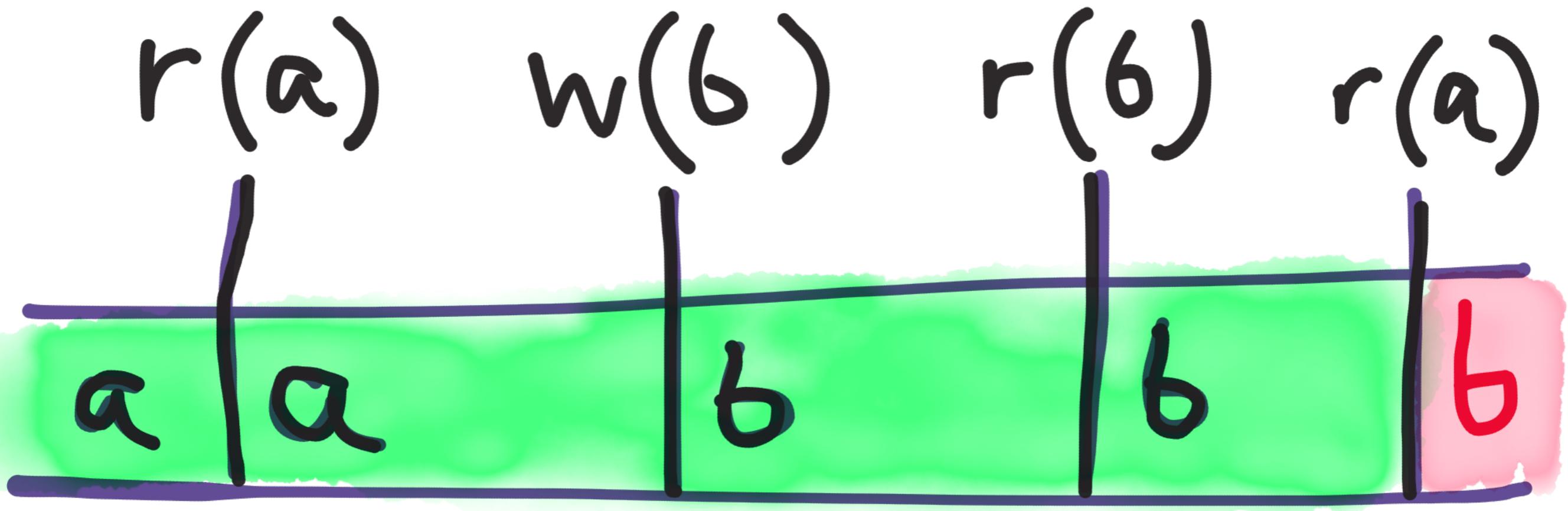
?

invoke ok

invoke fail

invoke ok

invoke info



Use a singlethreaded
model to verify path

invoke ok

invoke fail

invoke ok

invoke info

invoke ok

fail

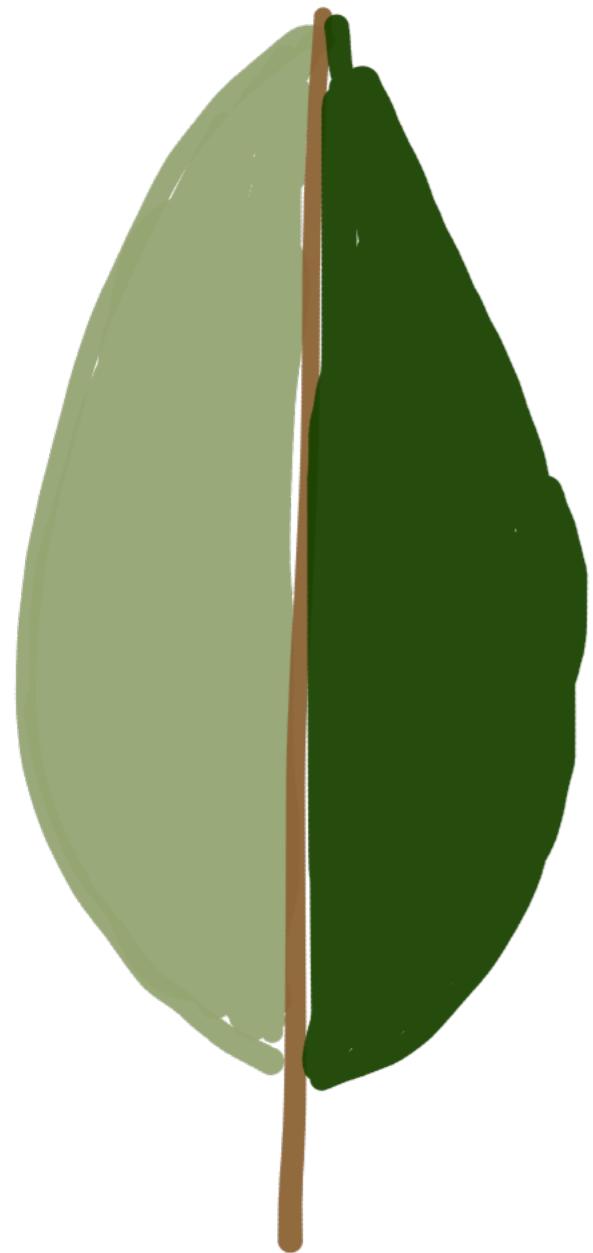
invoke ok

invoke info

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

So, what
have you
found?





mongoDB

First Jepsen test in
May 2013

~~#~~

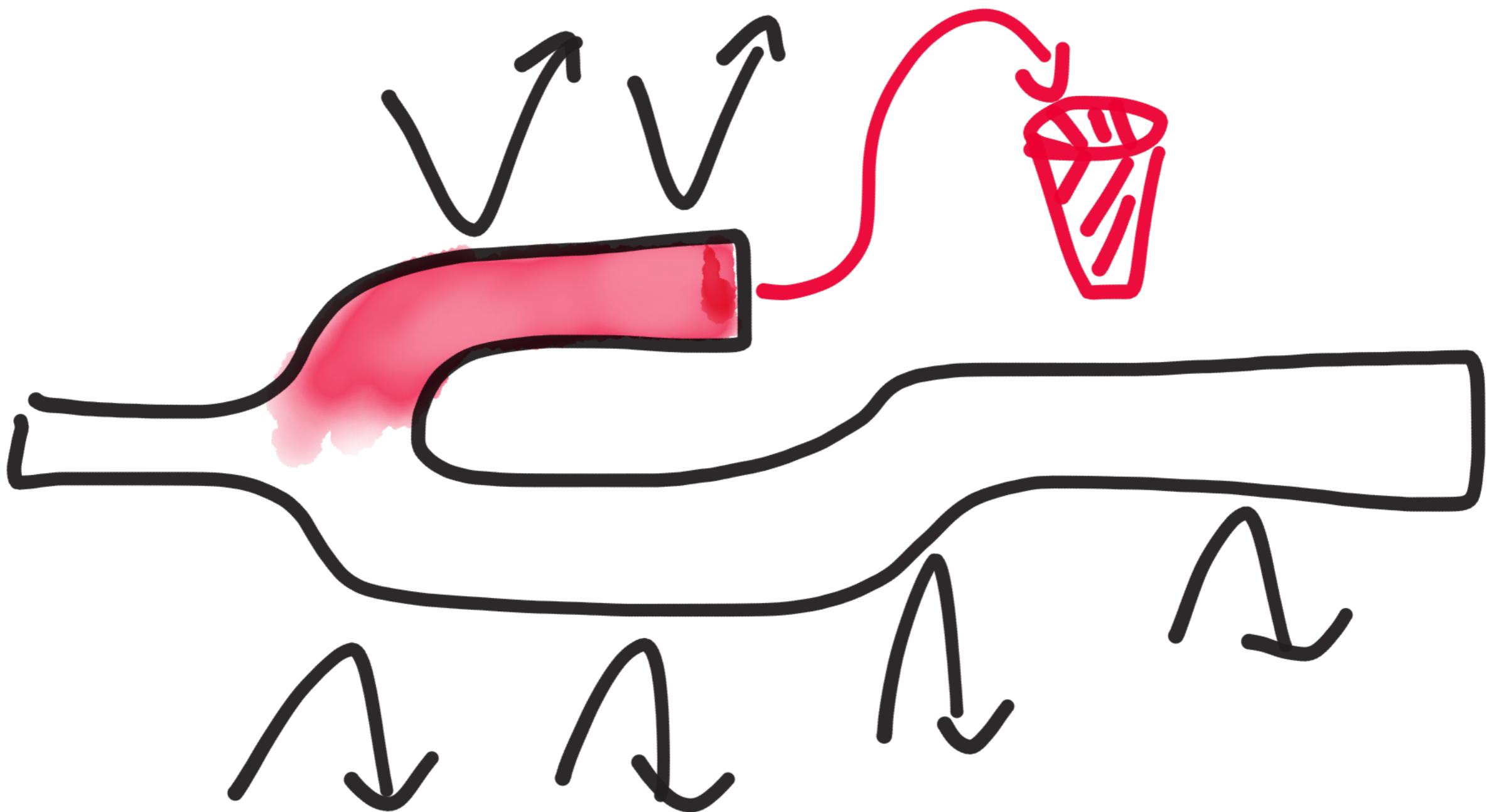
Data loss at all
WriteConcern levels

Unsafe defaults

- Some clients didn't check for errors at all

Unsafe defaults

- Better clients only wait for 1 or nodes, instead of a majority!



Rollbacks

Even Majority unsafe!

Network failures were
considered successful
responses



But These bugs
are fixed now!



Ok! Let's

try 2.6.7

Linearizable

CaS Registers

~~0~~ → write(1) → 1

$\emptyset \rightarrow \text{read}(\emptyset) \rightarrow \emptyset$

\emptyset $\rightarrow \text{read}(1) \rightarrow \boxed{\text{err}}$

0 → cas(0, 1) → 1

1 → cas(1, 2) → err

An Anomaly
Emerges!

47

CAS 0→4

0

1

r4

r3

CAS 3→1

2

r0

Write 4

3

CAS 4→2

4

Write 3

5

read 0

0

1

2

3

4

5

CAS 0→4

r4

r3

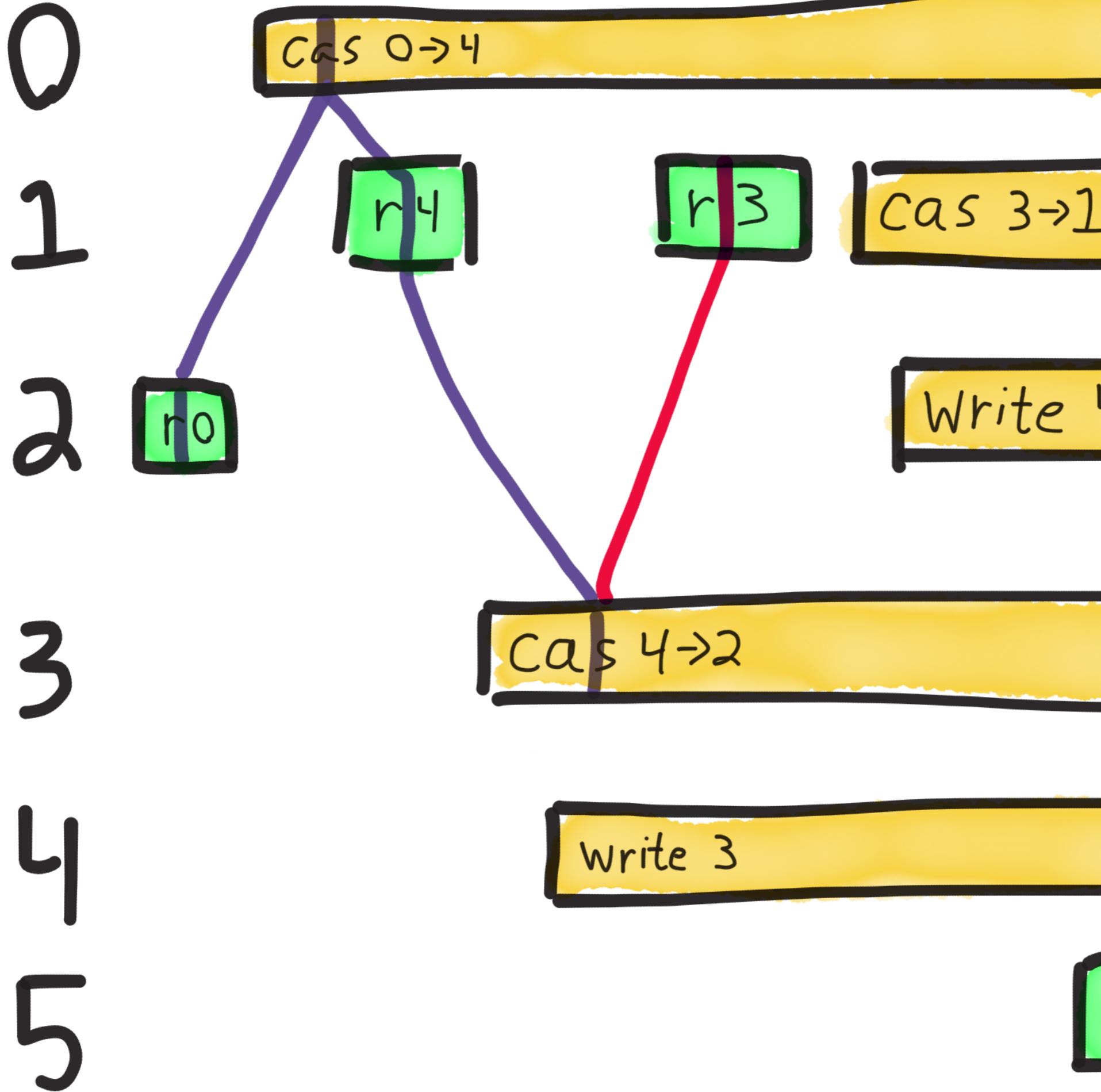
CAS 3→1

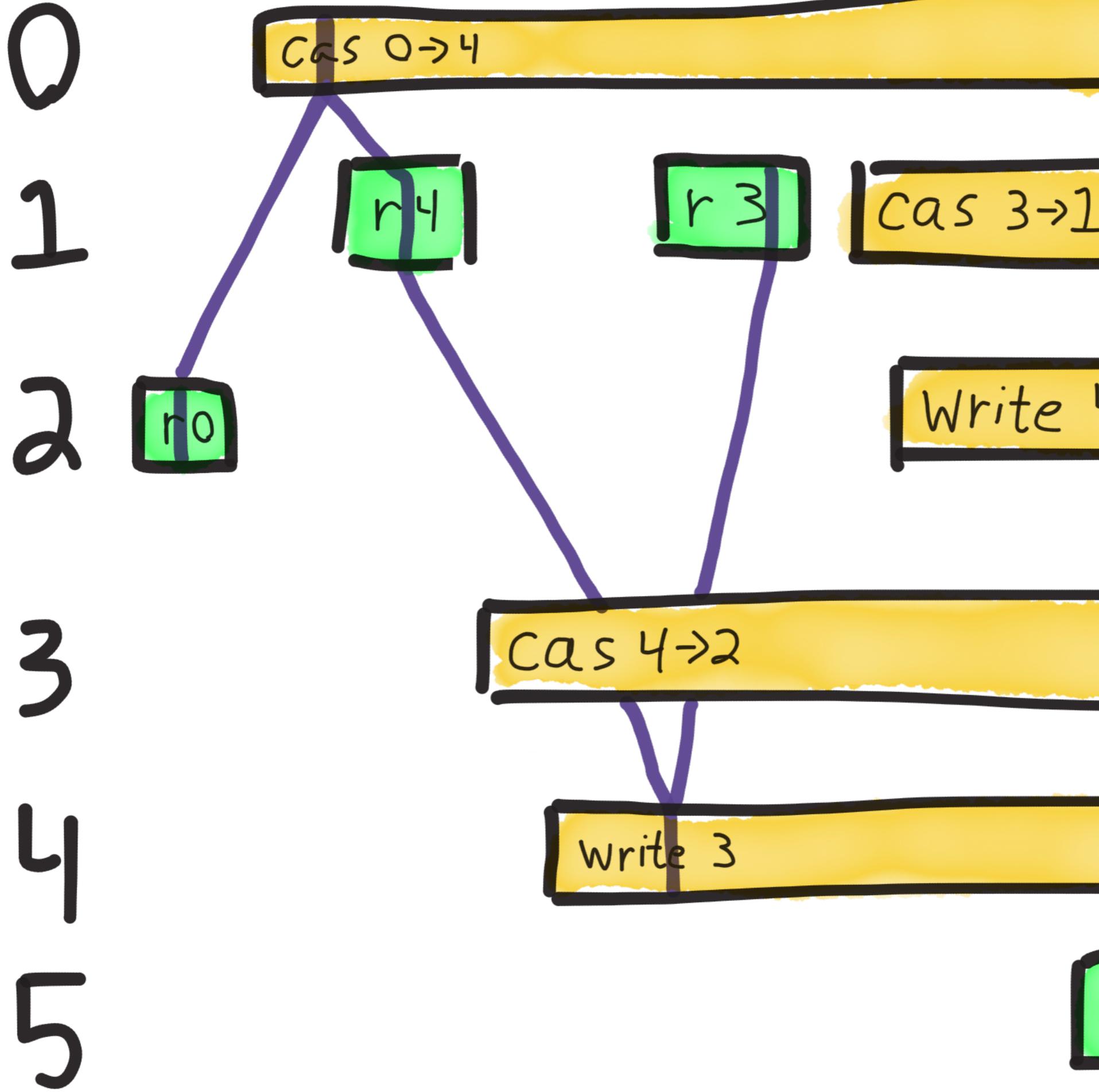
Write 4

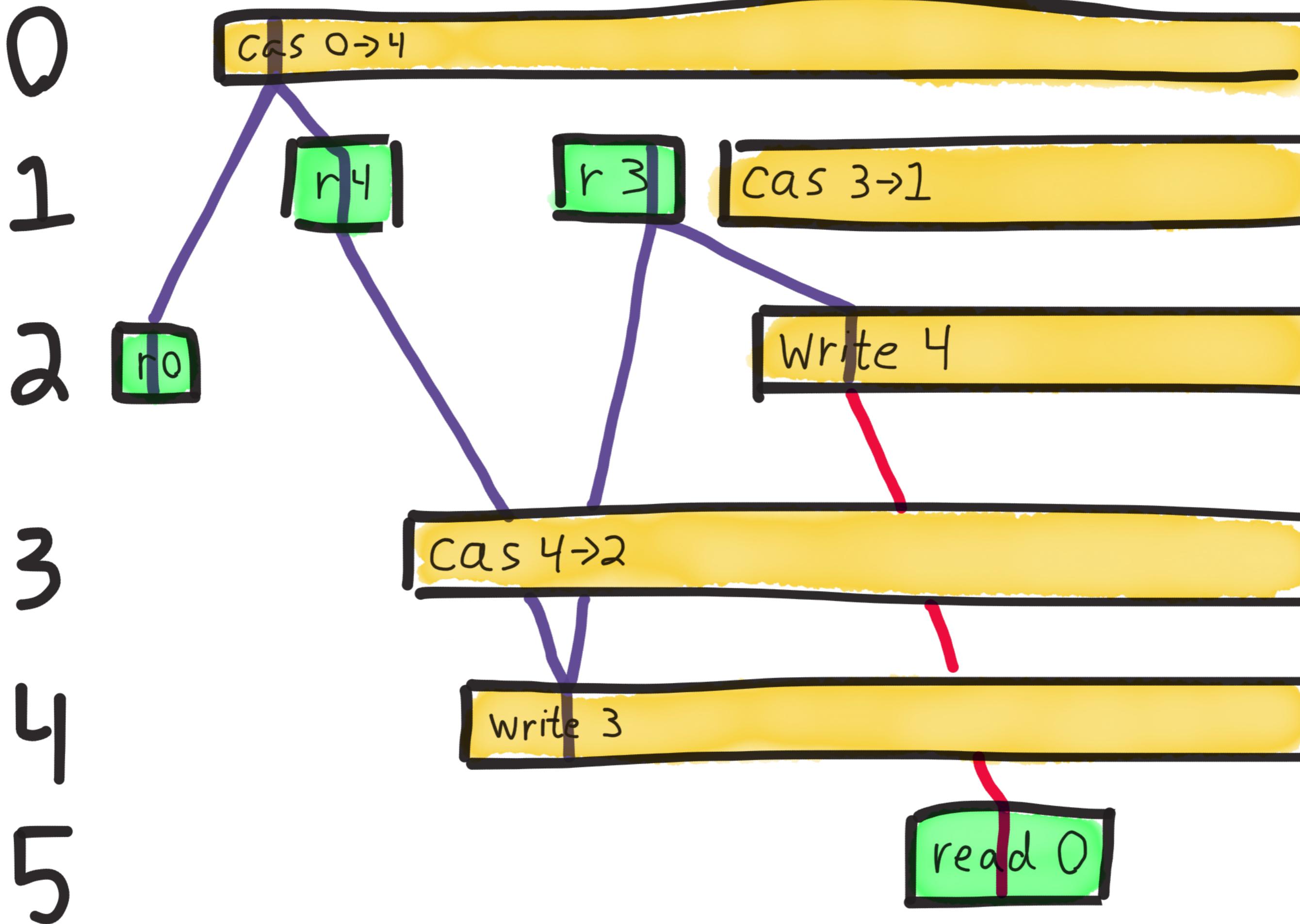
CAS 4→2

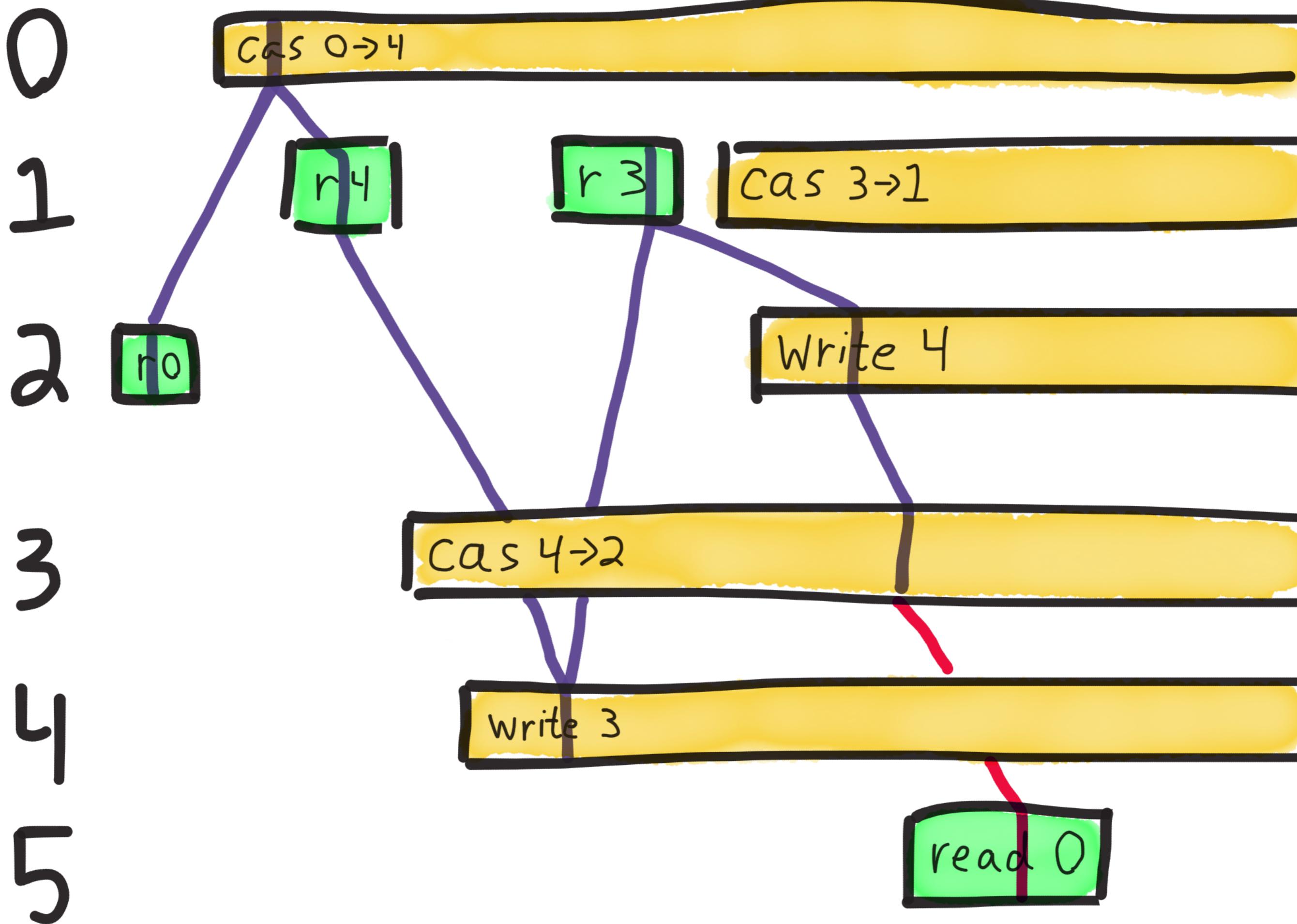
Write 3

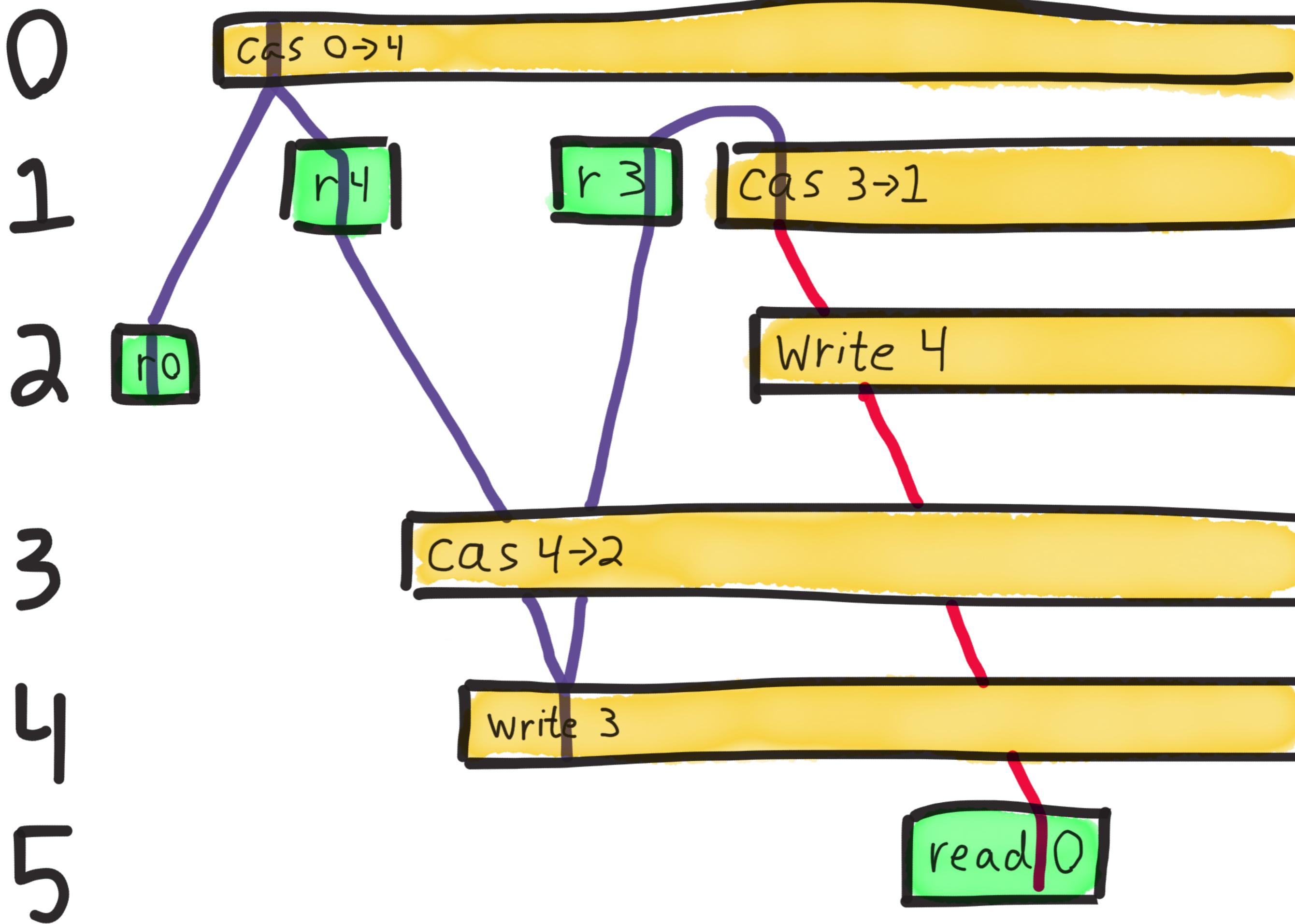
read 0

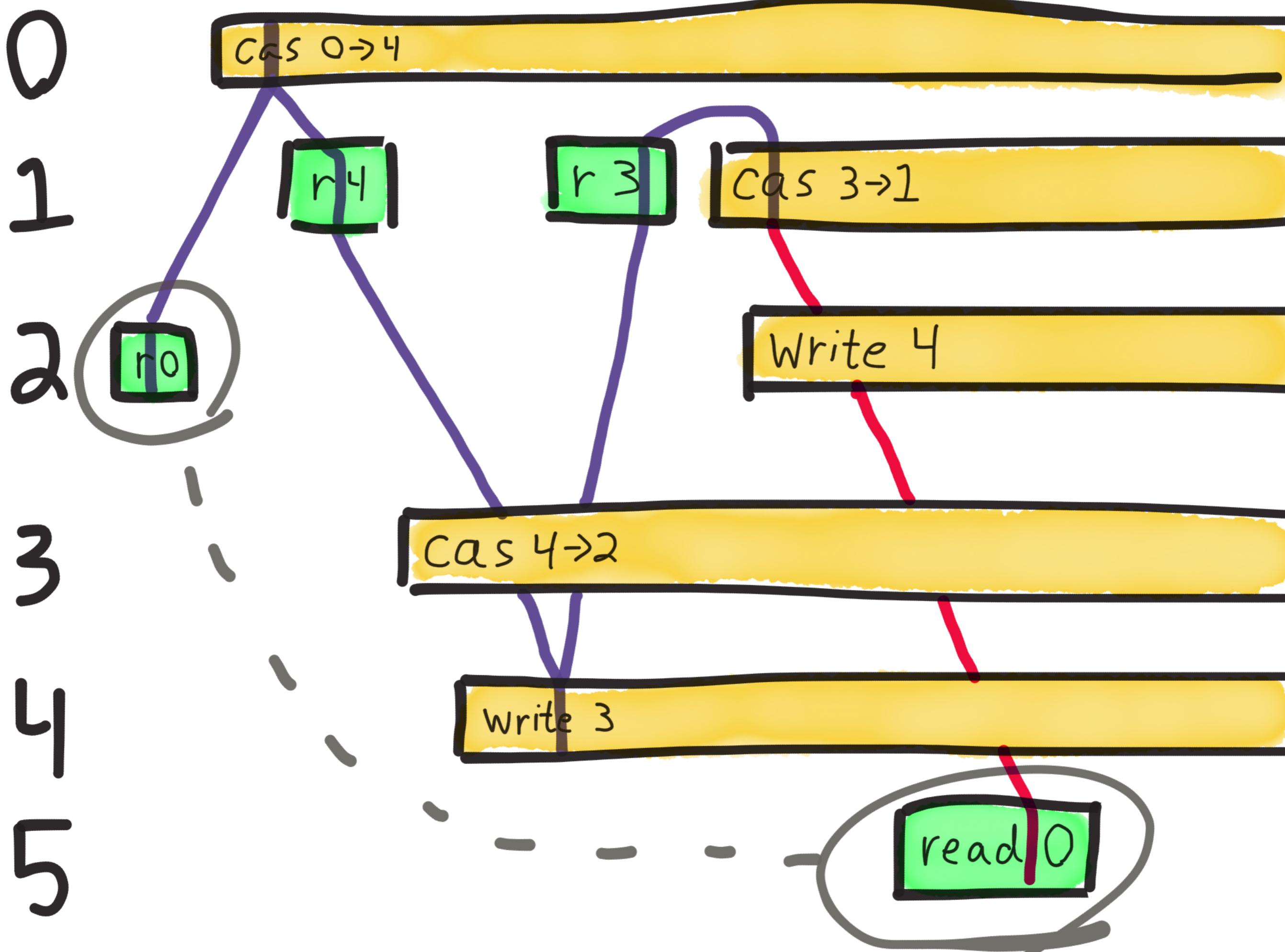




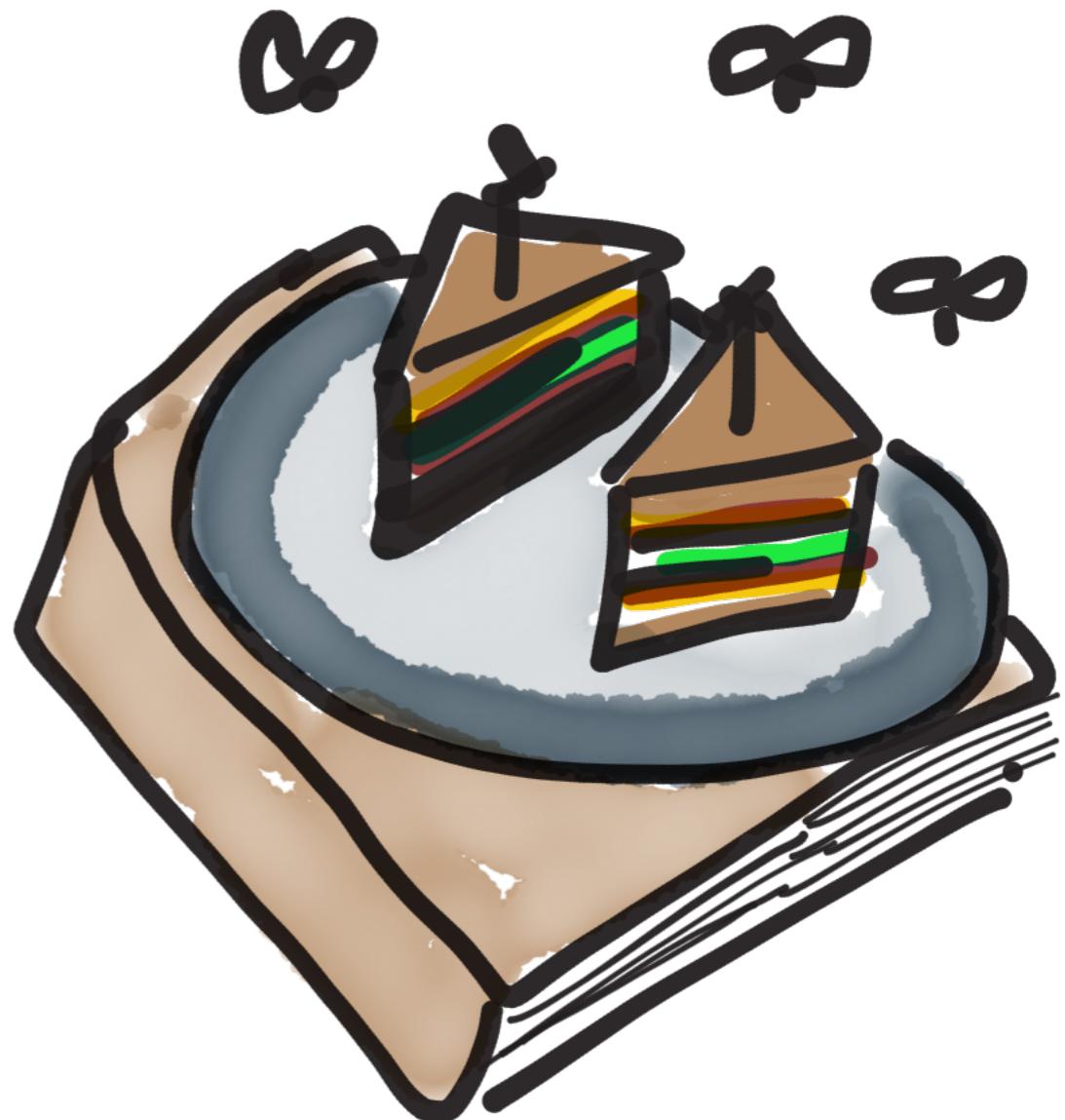








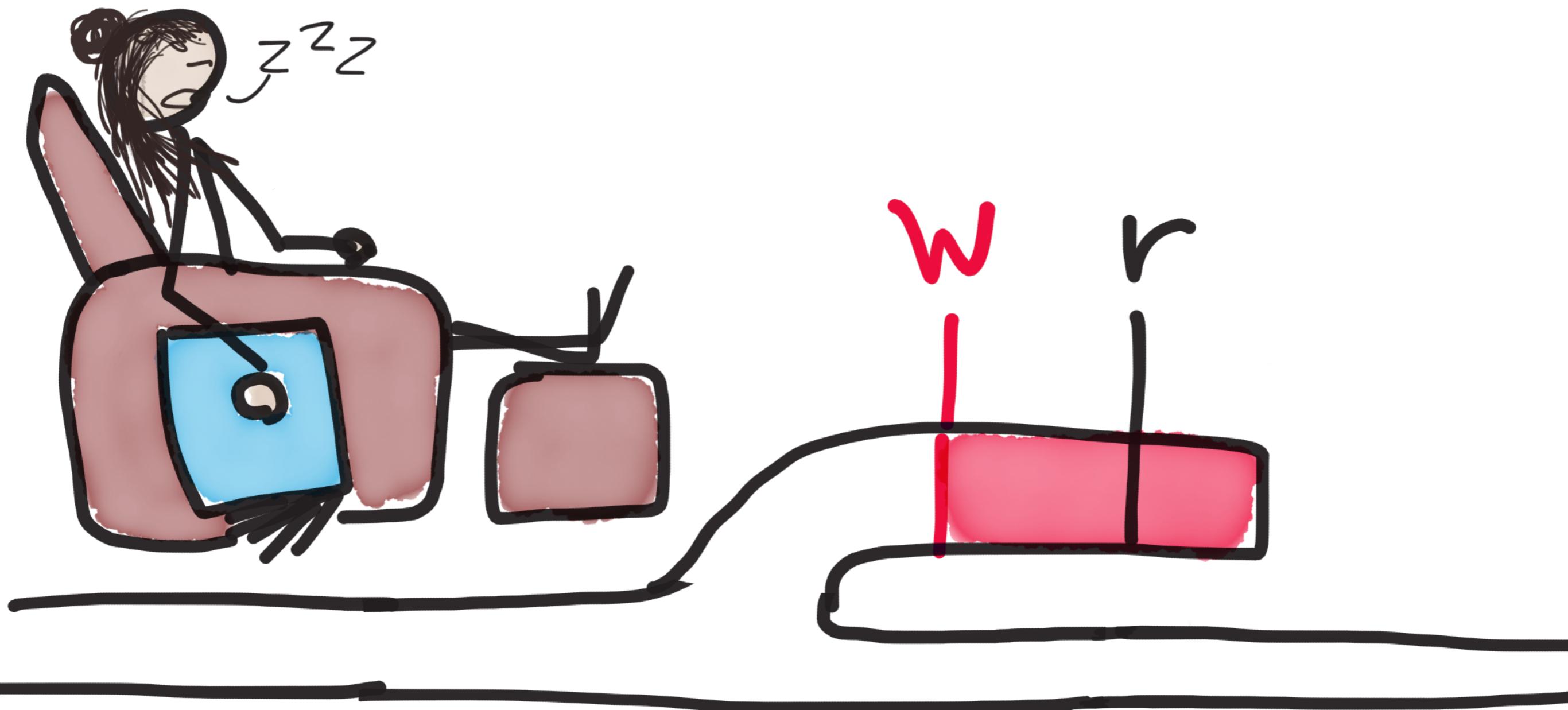
Stale
Reads



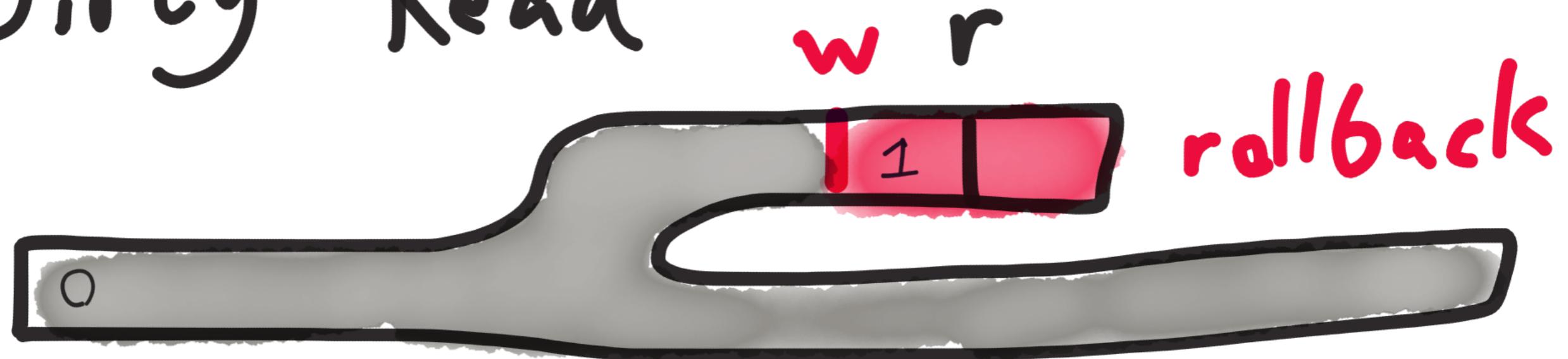
w

r

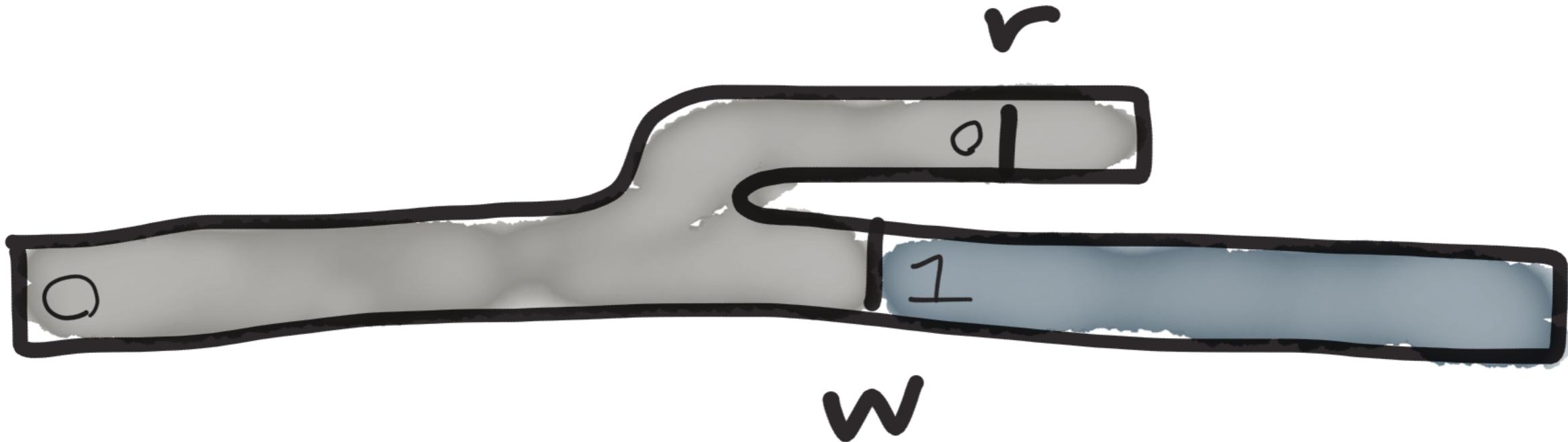
Read Uncommitted



Dirty Read



Stale Read

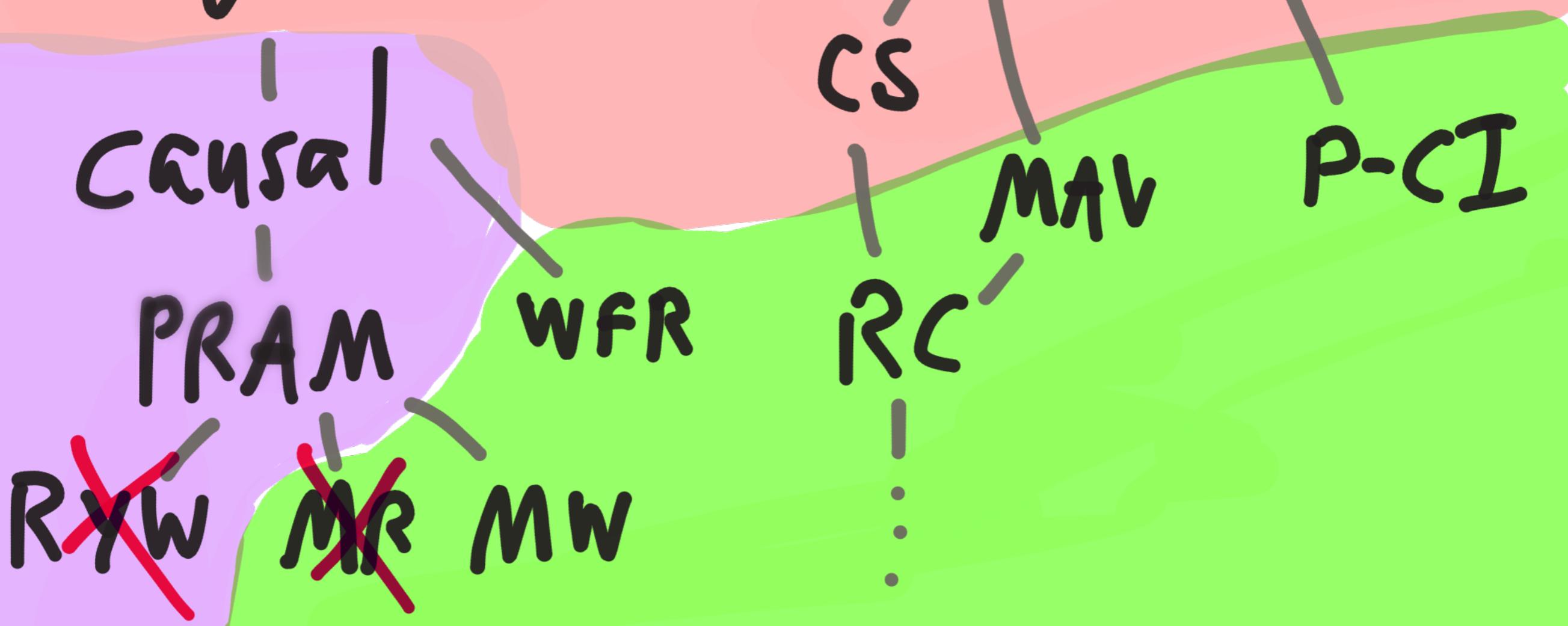


Server-17975

strong serializable

linearizable

sequential



~~strong~~ ~~serializable~~

~~linearizable~~

~~sequential~~

RR

CS

MAV

P-CI

RC'

⋮

~~Causal~~

~~PRAM~~

~~RYW~~

~~MR~~

~~MW~~

WFR

~~strong~~ ~~serializable~~

~~linearizable~~

~~sequential~~

~~Causal~~

~~PRAM~~

~~RW~~

~~MR~~

~~MW~~

WFR

RU

~~CS~~

~~RC~~

I

~~MAV~~

P-CI

~~RR~~

~~SI~~

~~strong~~ ~~serializable~~
~~linearizable~~
~~sequential~~



~~strong~~ ~~serializable~~

~~linearizable~~

~~sequential~~

~~Causal~~

~~PRAM~~

~~RXYW~~

~~MR~~ ~~MW~~

WFR

~~serializable~~

~~RR~~

~~CS~~

~~MAV~~

P-CI

~~RC~~

RU

Mongo
vs. docs

Mongo Recommendations

- Plan for stale reads
- Use CAS to verify
read values are
legitimate

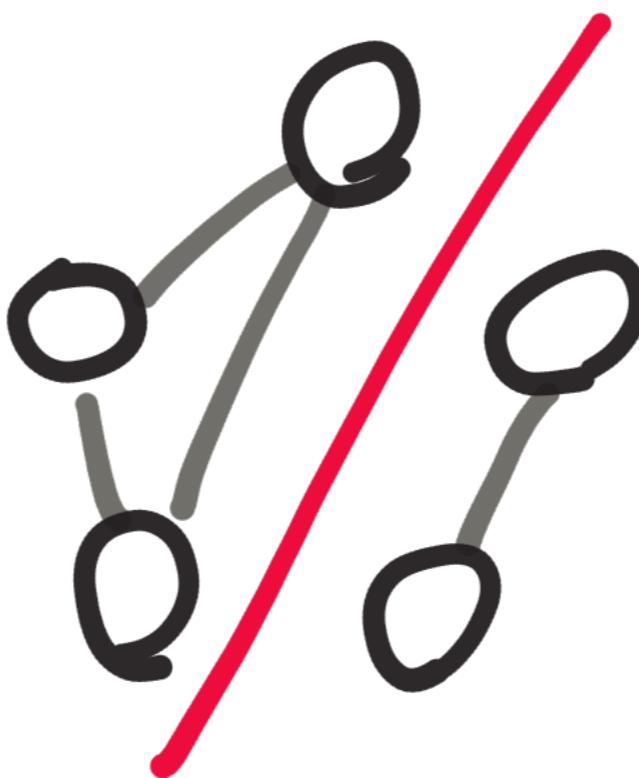


elastic
[search]

First Elasticsearch
tests: June 2014

v 1.1.0

Split Brain from



Massive data loss: 90%.

ES added a
great page on
safety issues

Screenshot of
softw) docs

Screenshot of
closed issue

Folks are
still referring
to the last talk

Screenshot of

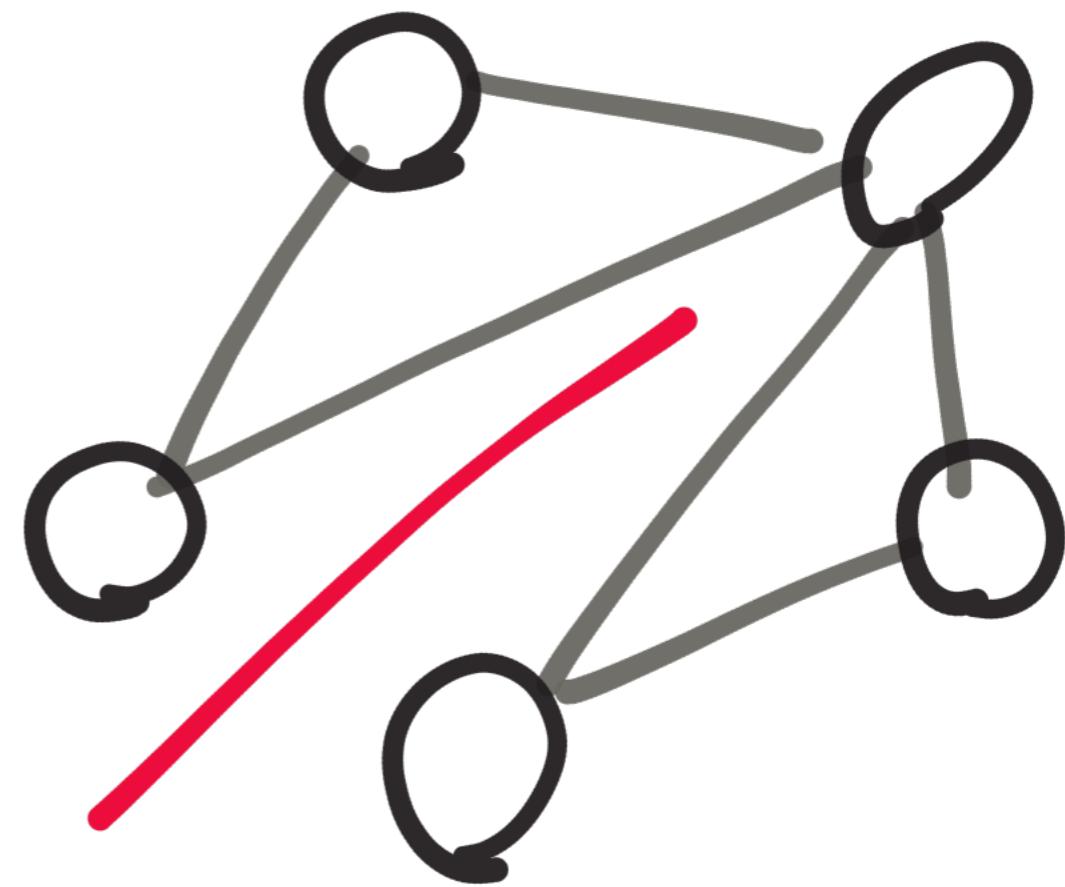
"not a big deal"

Retesting with

ES 1.50

Intersecting
partitions still cause

data loss



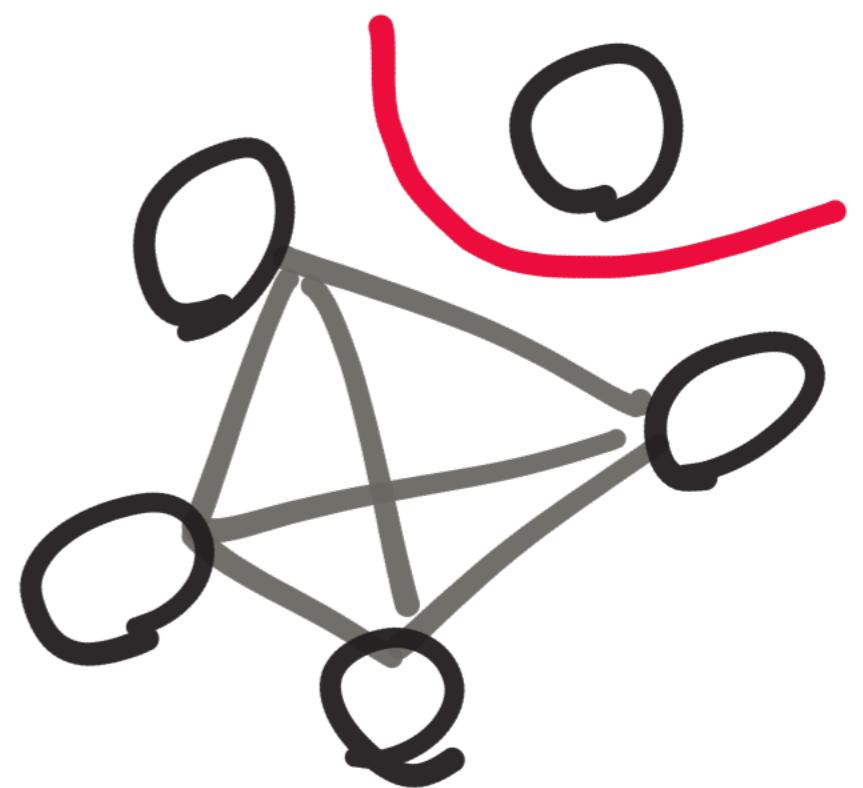
Long-lived split-

brain appears

fixed, gut docs

are still last

Single node
partitions still cause
data loss!



Electing a new
primary still takes

~90 seconds

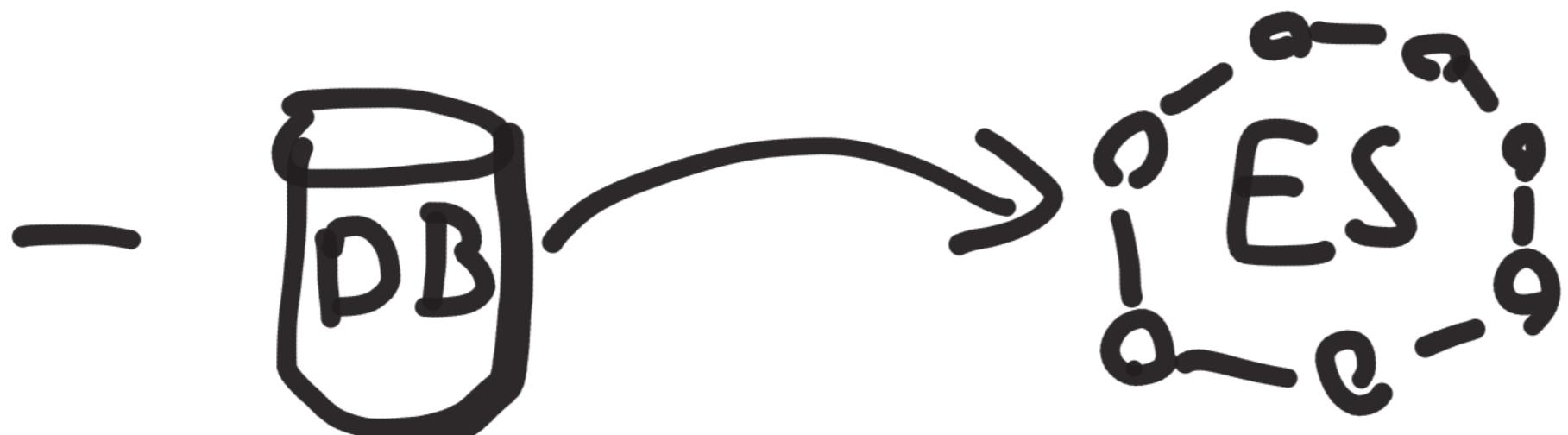
- Writers will stall

Elasticsearch is
doing better, esp.

with documentation !

ElasticSearch Recommendations

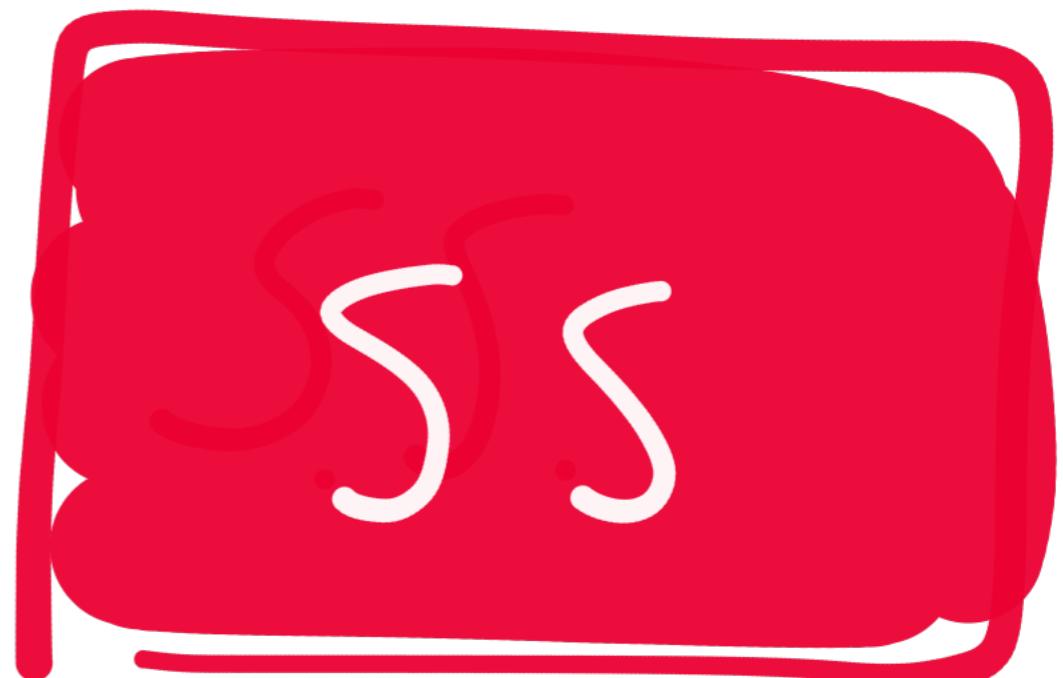
- No longer catastrophic
- But, still loses data in all failure modes



The logo consists of the word "AEROSPIKE" in a white, outlined, sans-serif font. The letter "A" is stylized with a diagonal line through its top-left section. The entire logo is set against a solid red rectangular background.

High-performance
5-D KV store
for online analytics
(esp. in adtech)

"ACID reliability
with no downtime
— ever!"



Screams like a ho

data loss

"No data loss. Row-level locking... immediate consistency (ACID), with synchronous replication"

We have to go

Deeper

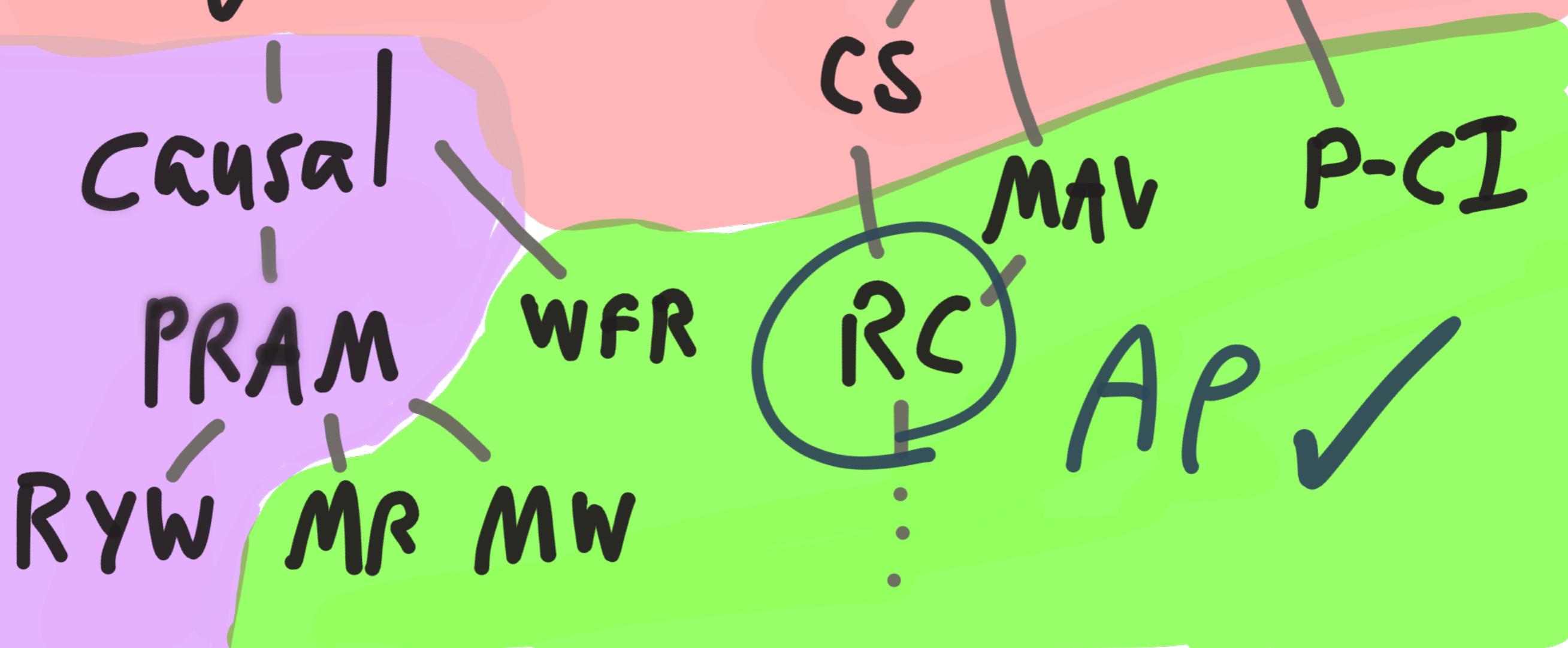
"AS is by and large an
AP system that provides
high consistency"

"Aerospike provides read-committed isolation level using record locks to ensure isolation between [txns]"

strong serializable

linearizable

sequential



"for ops. on single keys
with replication, Aerospike
provides immediate consistency
using synchronous writes to
replicas"

"...read requests are
guaranteed to find
the newly written data..."

strong

serializable

linearizable

sequential

Causal

PRAM

RYW

MR MW

WFR

RR

CS

RC'

⋮

SI

MAV

P-CI

"Virtually eliminate
partition formation
proven by years of
deployment in DC and
Cloud environments."

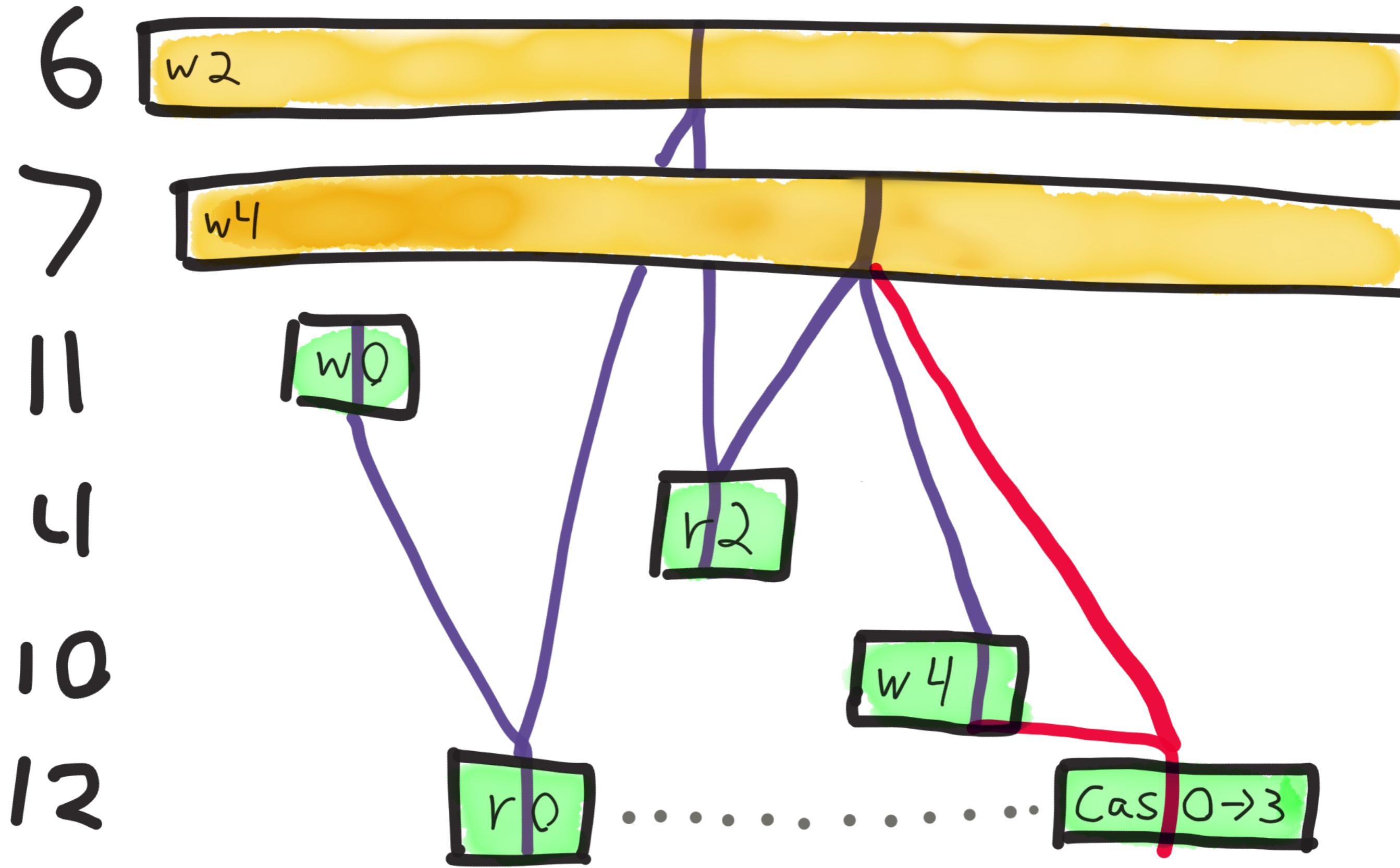
"A key design point of AS
is to setup cluster nodes
that are tightly coupled so
that partitions are virtually
impossible to create"

Screenshot:

Deploy page

So for a CAS

Register...



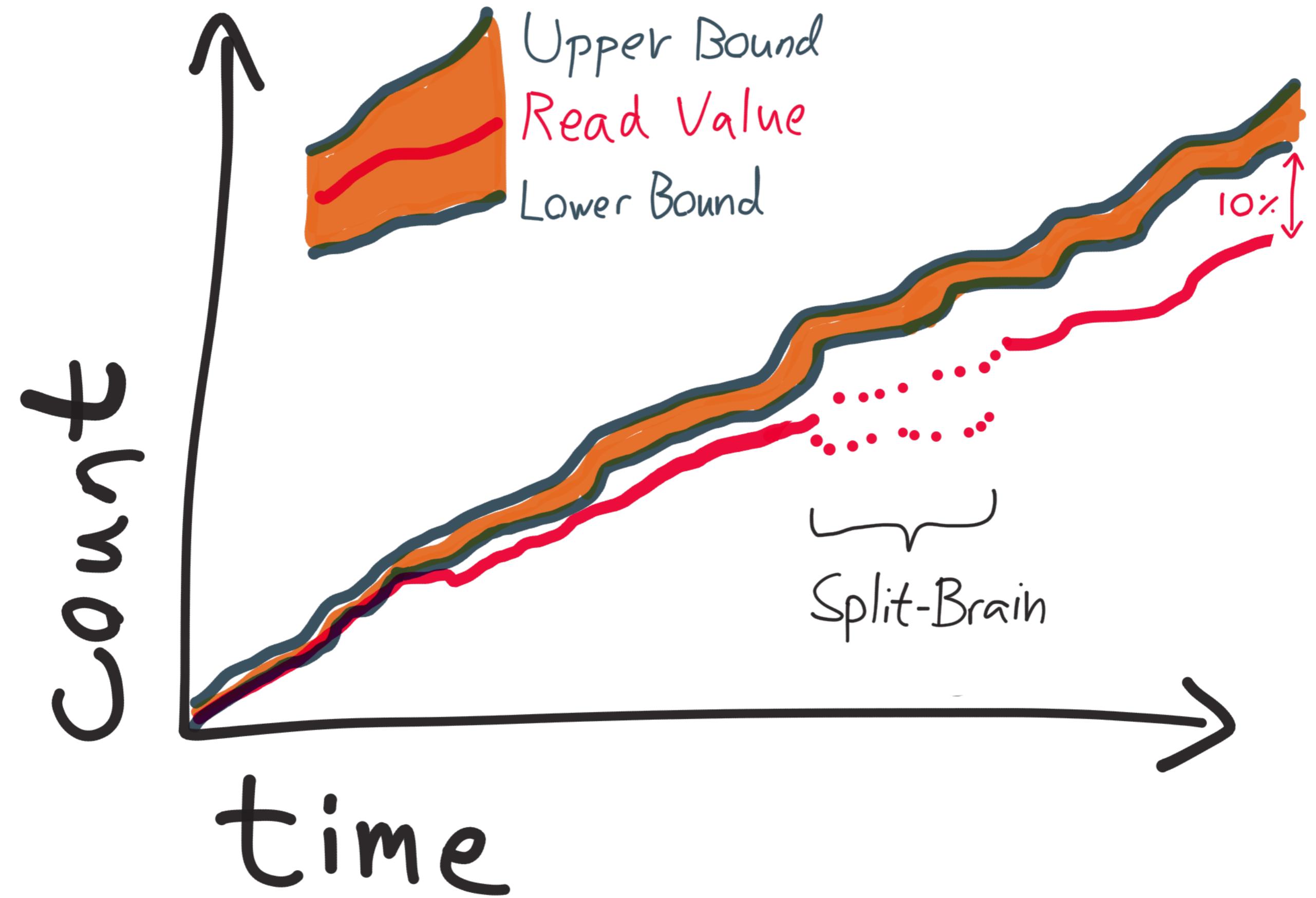
What about
Commutative ops?

$$a+b = b+a$$

Counters will

drop increments

during a partition



"Deployment modes
AP vs. CP"



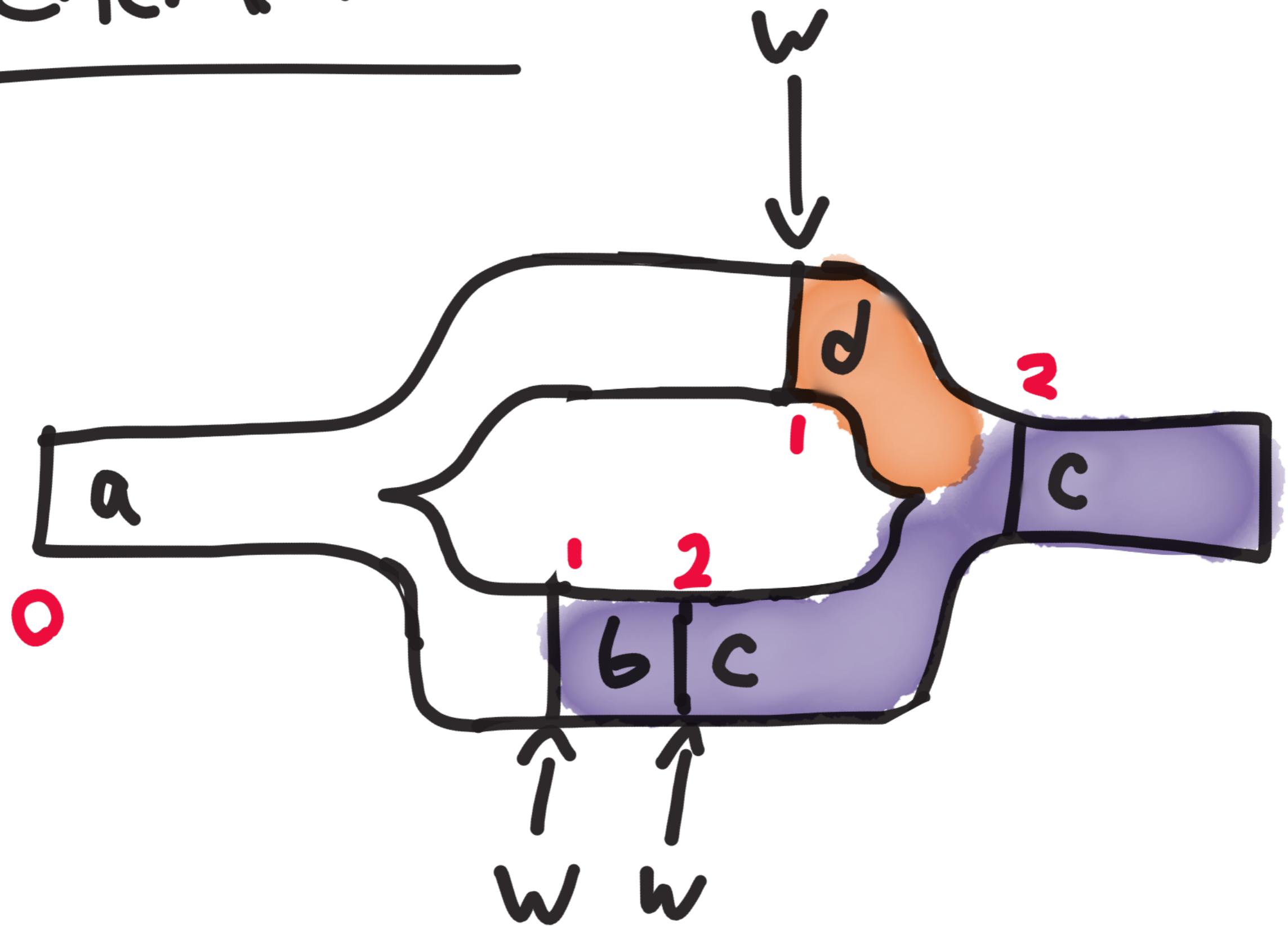
CP is vaporware

OK so how

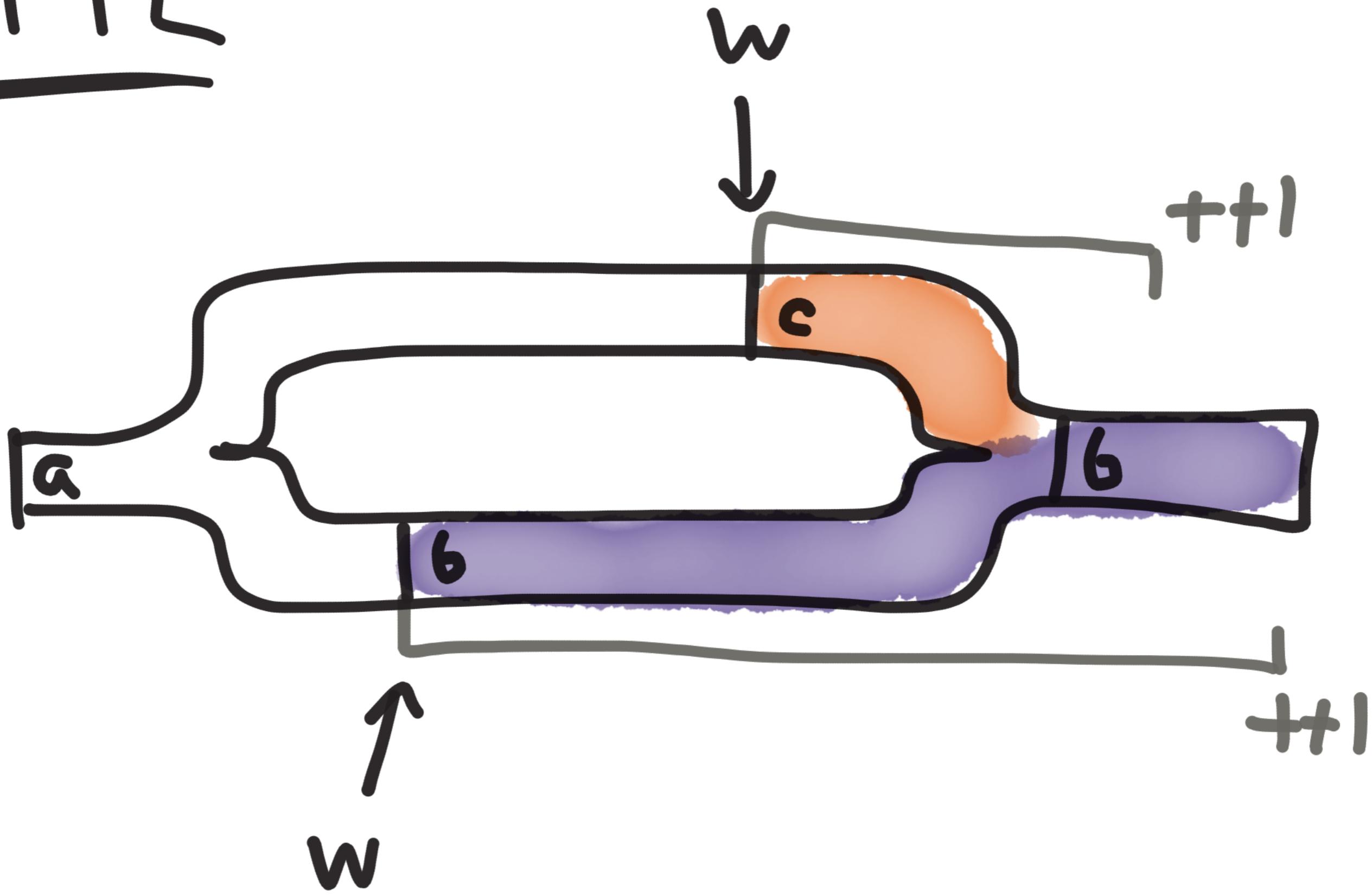
does AP mode

work?

Generations



TTL



UPDATES

~~strong serializable~~

~~Linearizable~~

~~Sequential~~

~~CAUSA~~

PRAM

Ryw

10

MW

RR

2

12

120

M&V

The ST logo consists of the letters 'ST' in a bold, black, sans-serif font, positioned on a diagonal red stripe.

P-Cl

"Why not use your
PAXOS

implementation?"



Probably fine

for stats data,

but expect some loss

Recap



Absence of evidence

— is not —

Evidence of Correctness

Mongo & EJ

have made

improvements!



Still Surprises

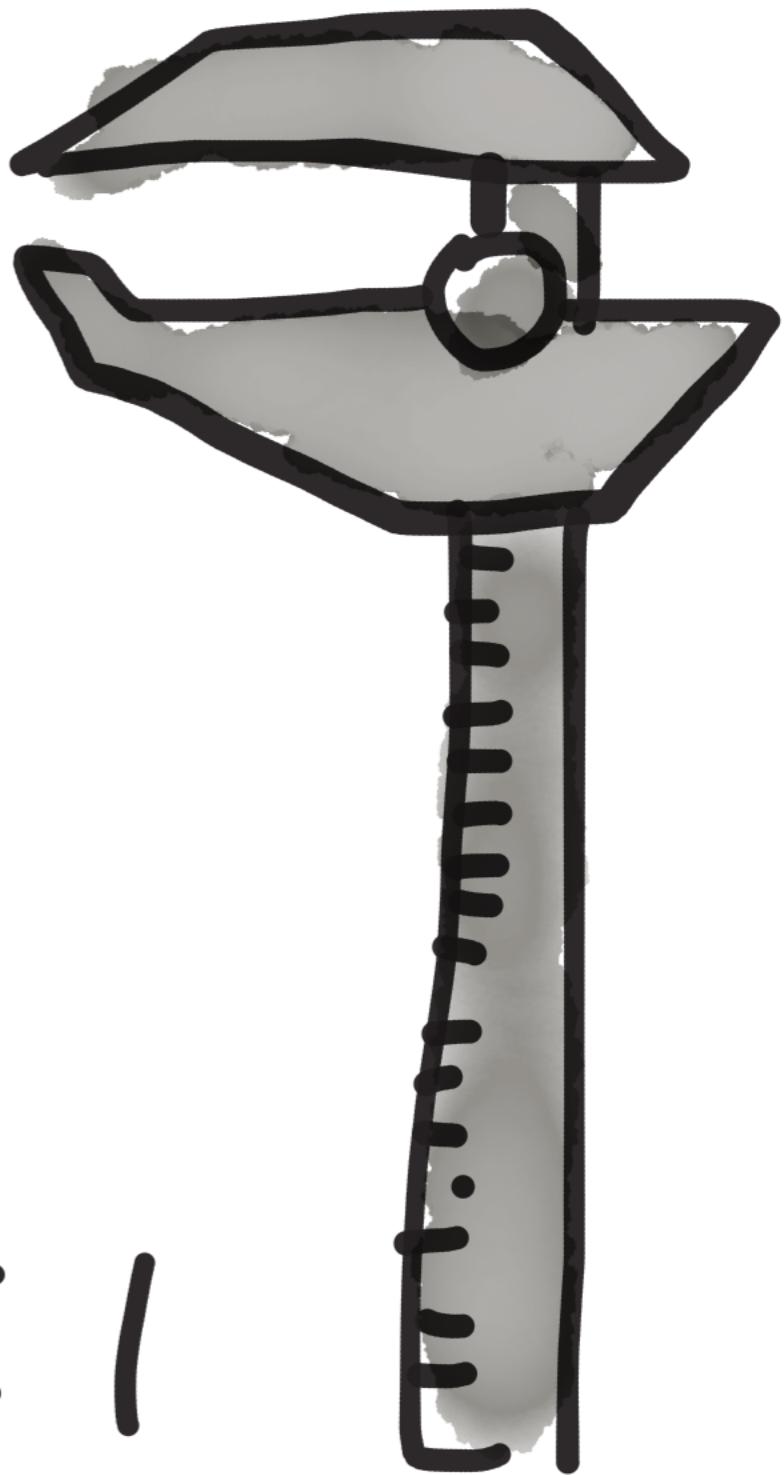
Lurking in the



Read the docs!



Then
MEASURE
— far —
YOURSELF!



*Figure out the
invariants
your system needs*

How much



Can you tolerate?

Consider
your
failure modes

Partitions

I/O/GC pauses

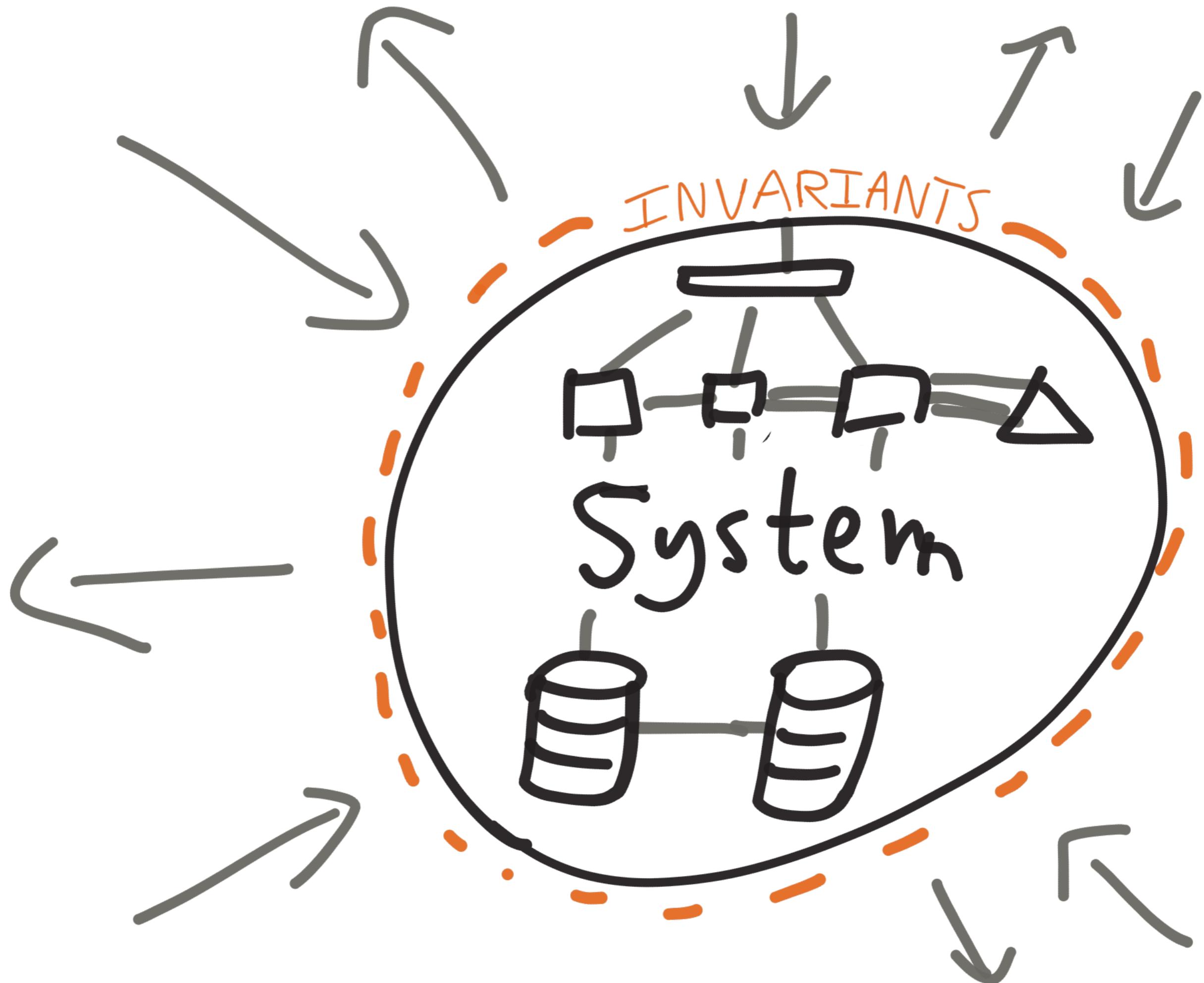
Node failure

Clock skew

Test your

Systems

end ————— to ————— → end



Thanks!

Kevin Porter

Lucien Valmar

Boaz Leskes

Lee Hinman

Matt Kangas

Stripe

Evan Broder

Asya Kamsky

Kelly Stirman

Marc Hedlund