

J E P S E N

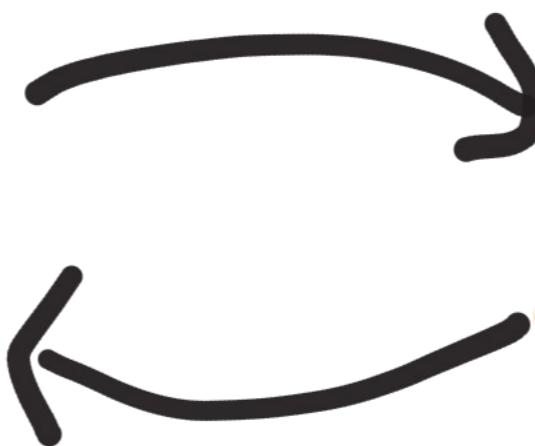
IV

Hope Springs
Eternal

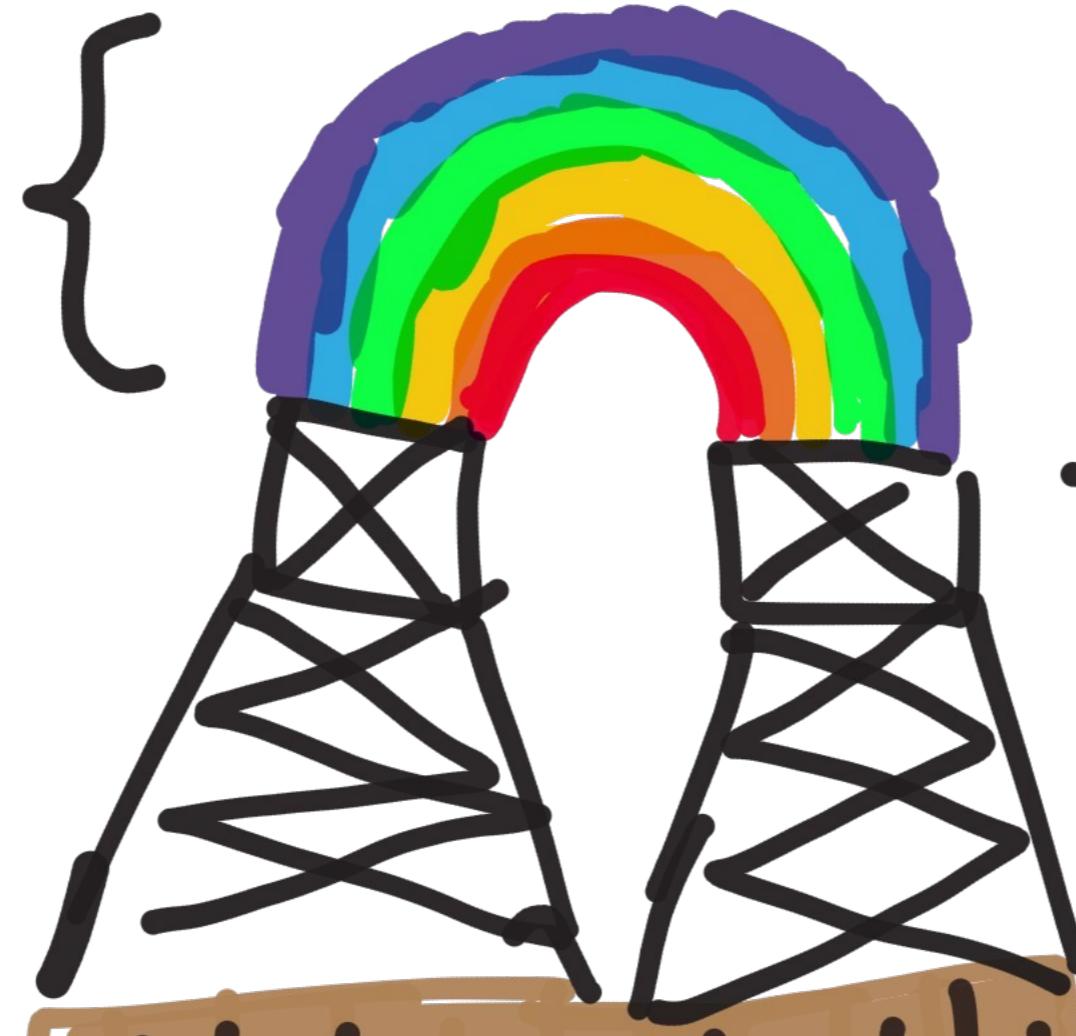


@aphyr
(Kyle Kingsbury)

style



Public
API {



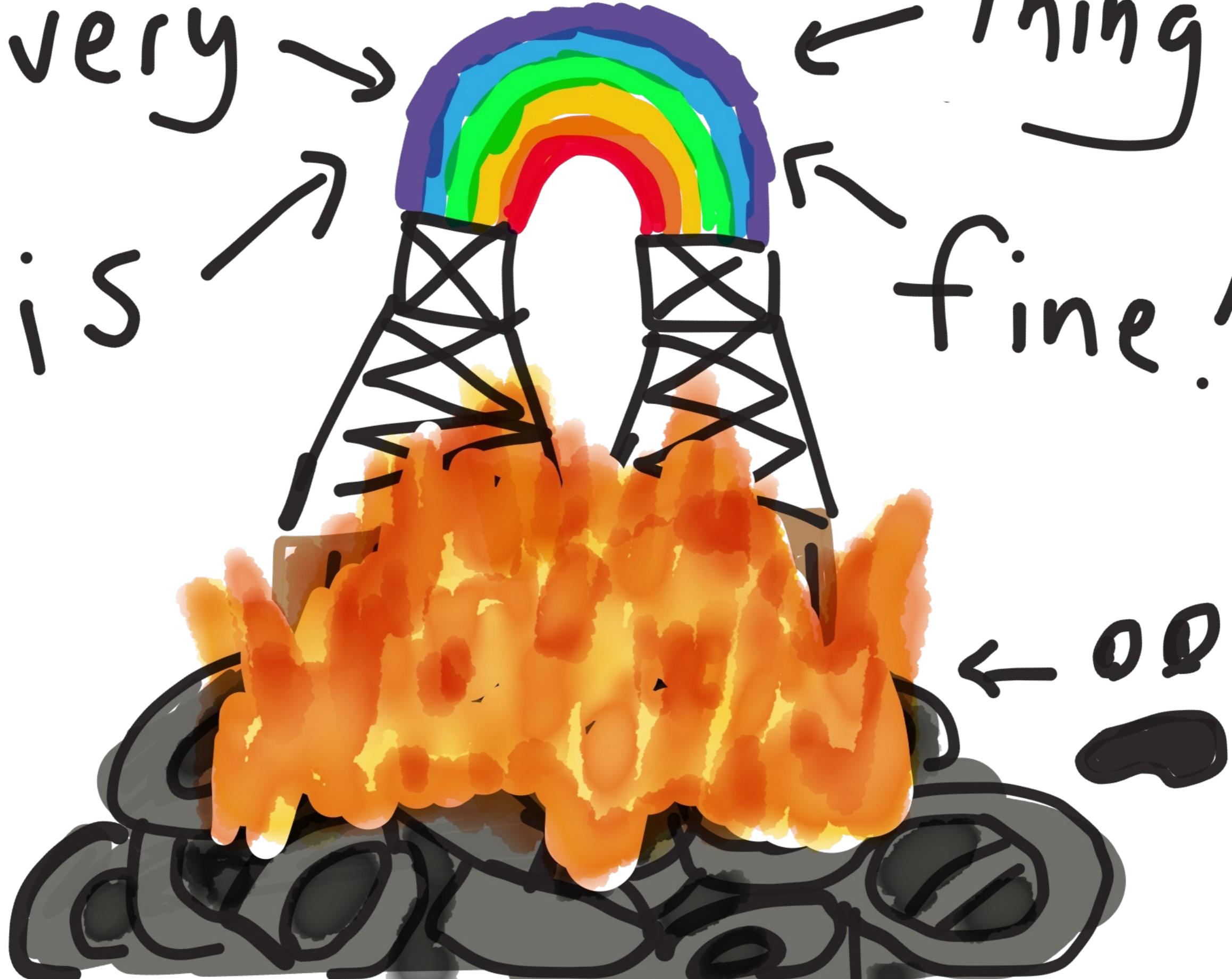
} API code

Ruby {



DBs

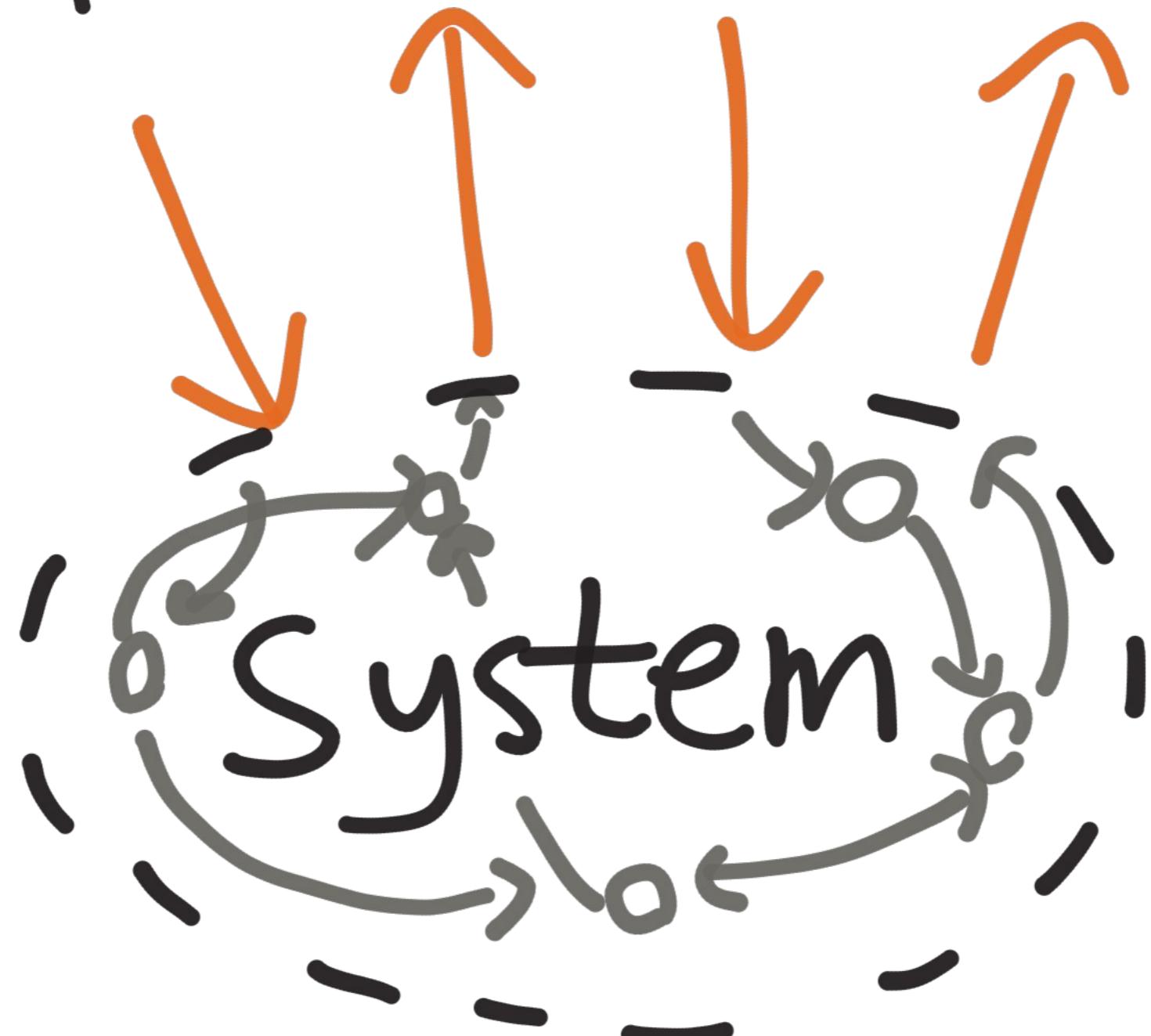
Every →
is →
thing ←
fine !



Jepsen

github.com/aphyr/jepsen

Environment





— INVARIANTS —

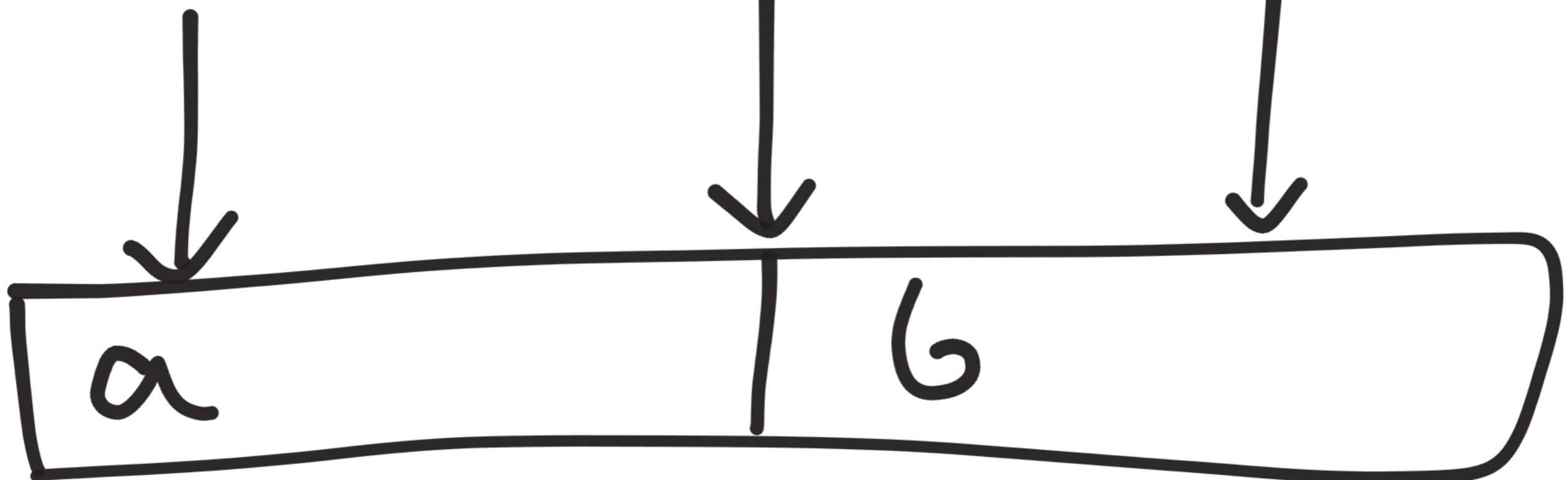


var $X = \alpha$

print X

$X = b$

print X

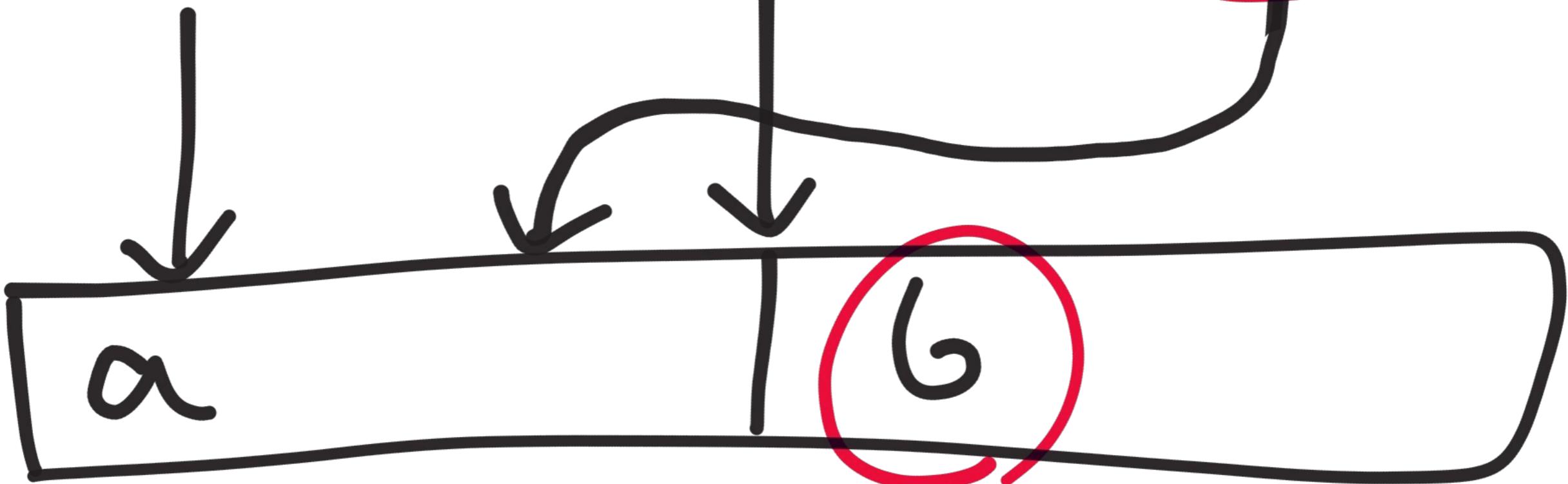
$r(a)$ $w(b)$ $r(b)$ 

time →

$r(a)$

$w(g)$

$r(a)$



time →

Invariants

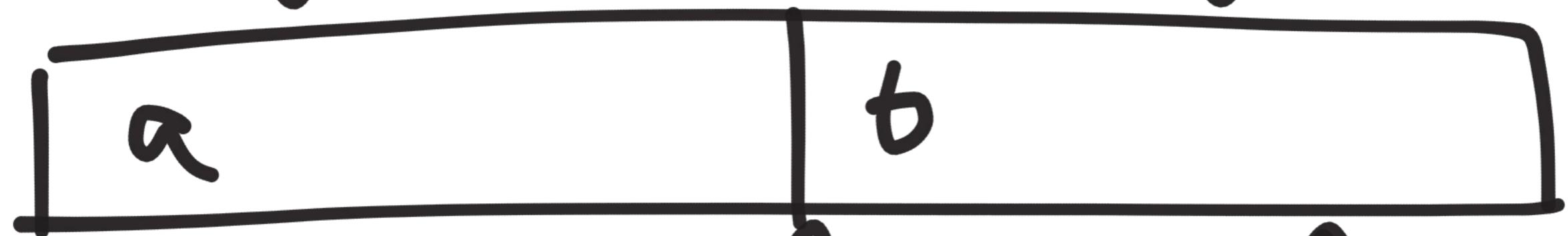
Constraint

Histories

What about

concurrency?

$r(a)$



$r(b)$



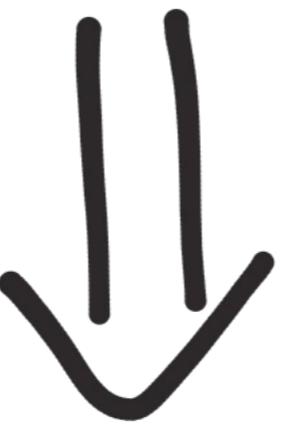
$w(b)$



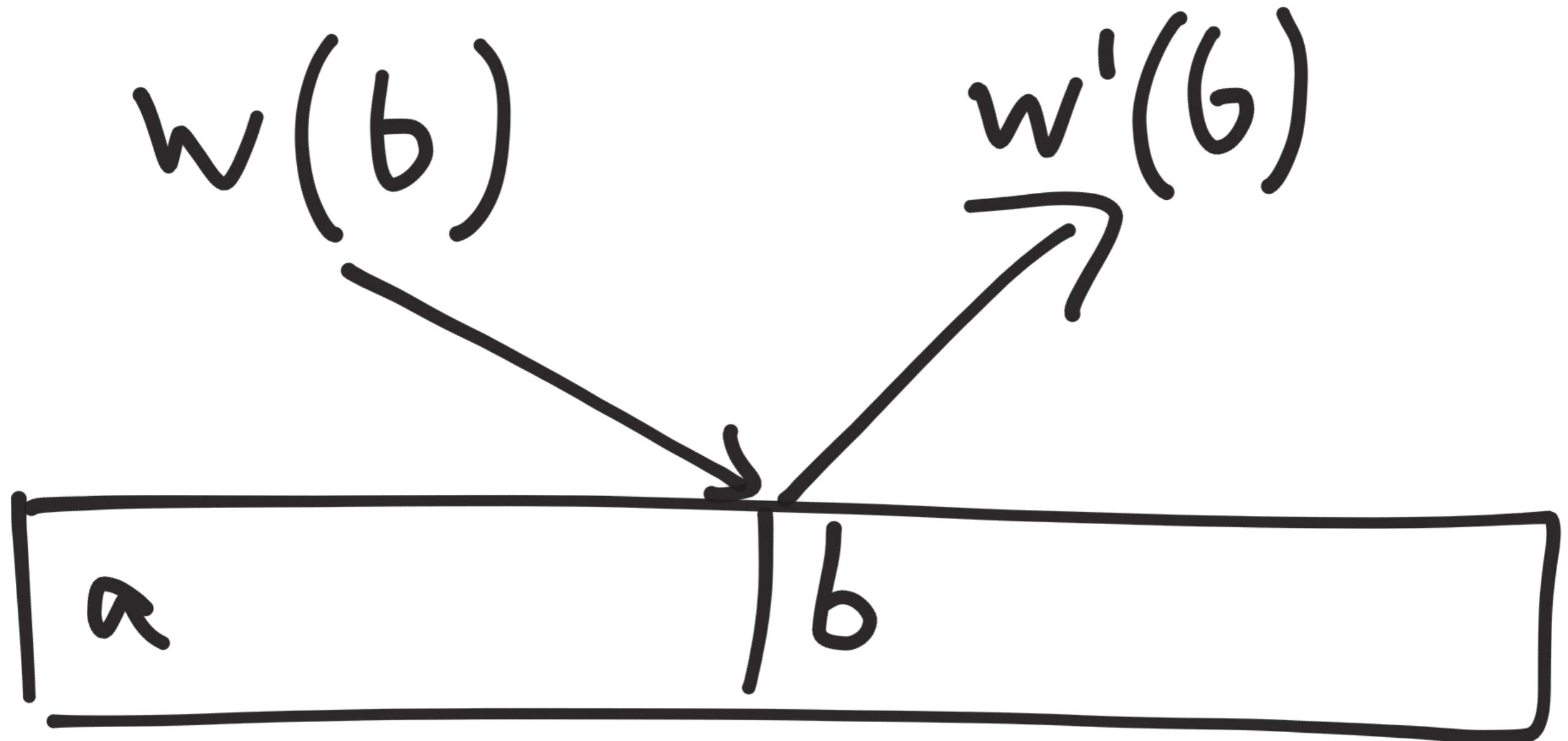
$r(b)$



Distributed

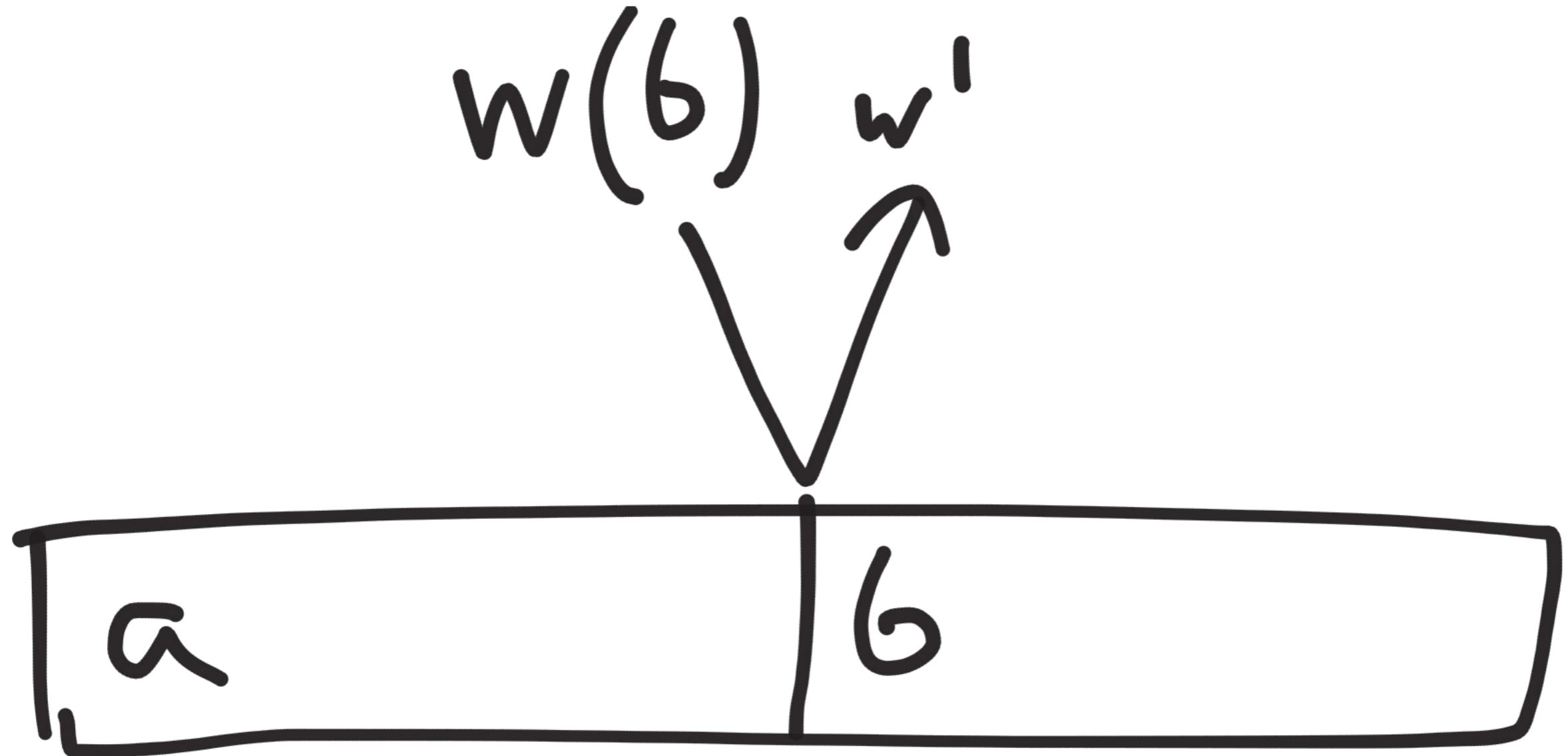


Distance



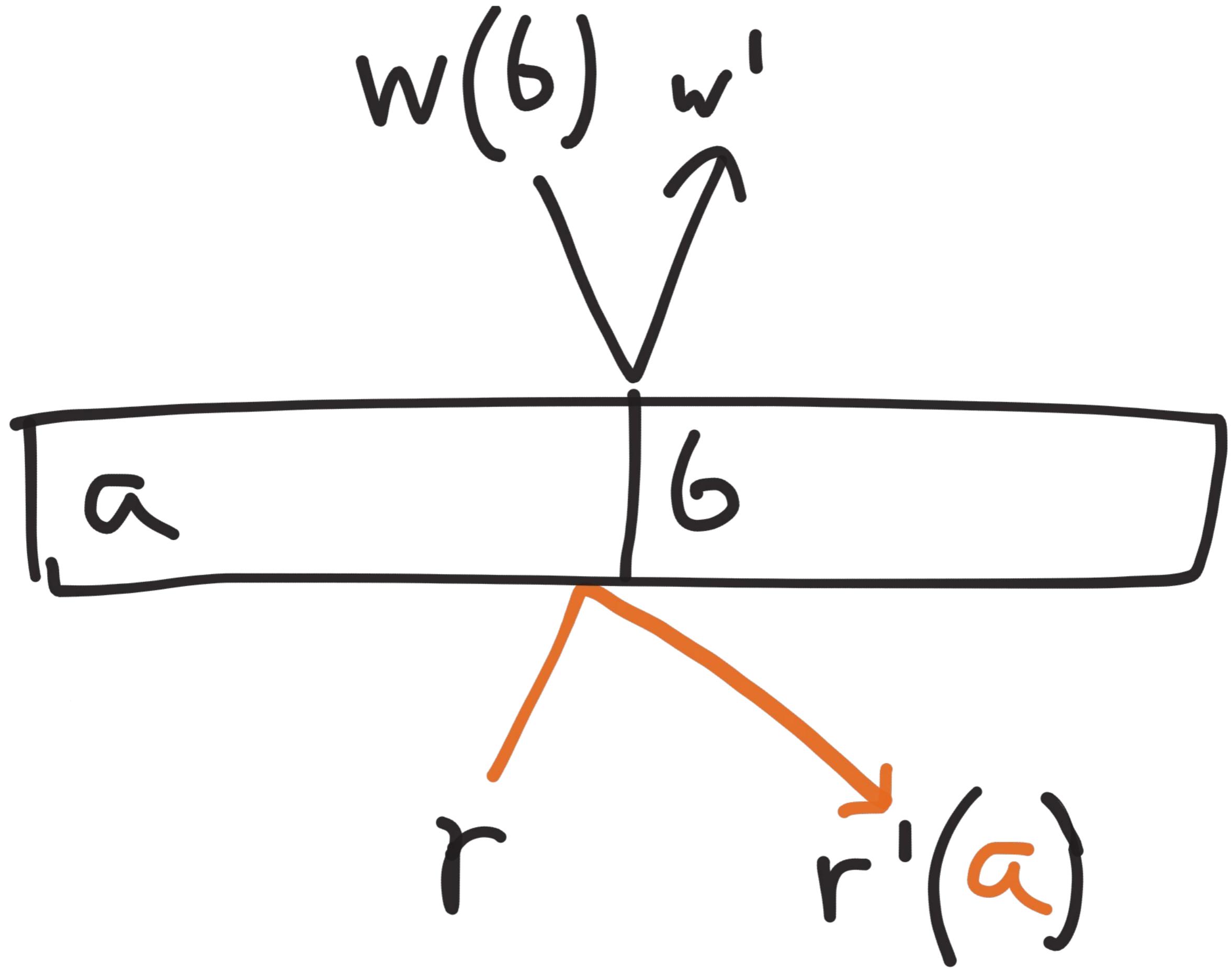
Mess@ger take time

Delays imply
ambiguous orders

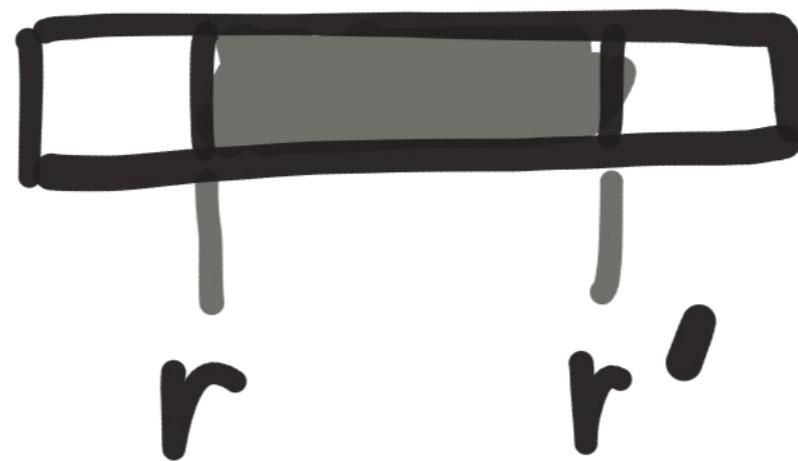
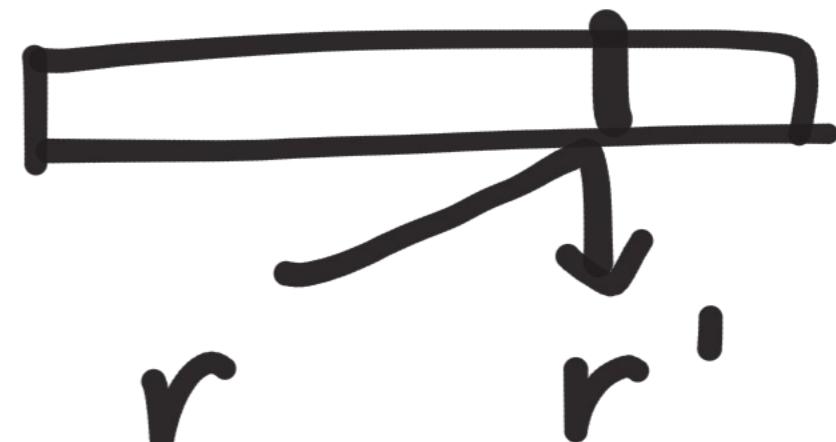
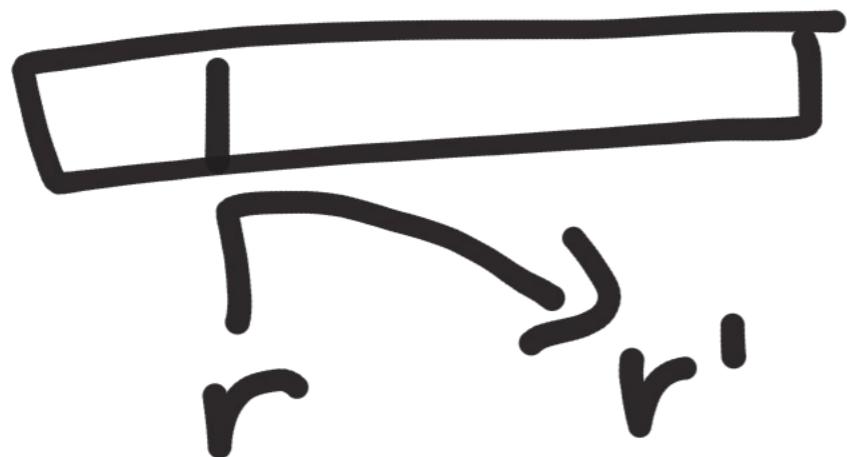


r

$r'(b)$



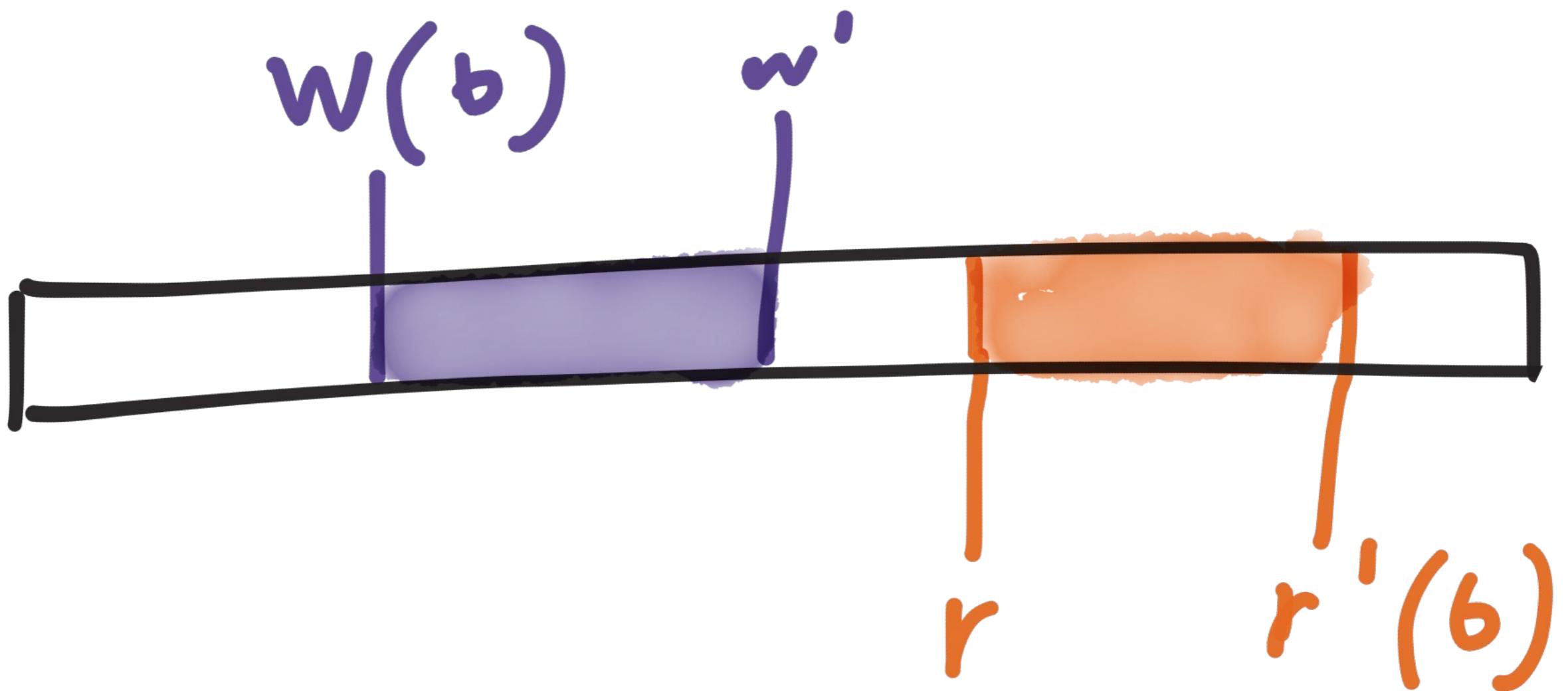
But there are finite
bounds on ambiguity!



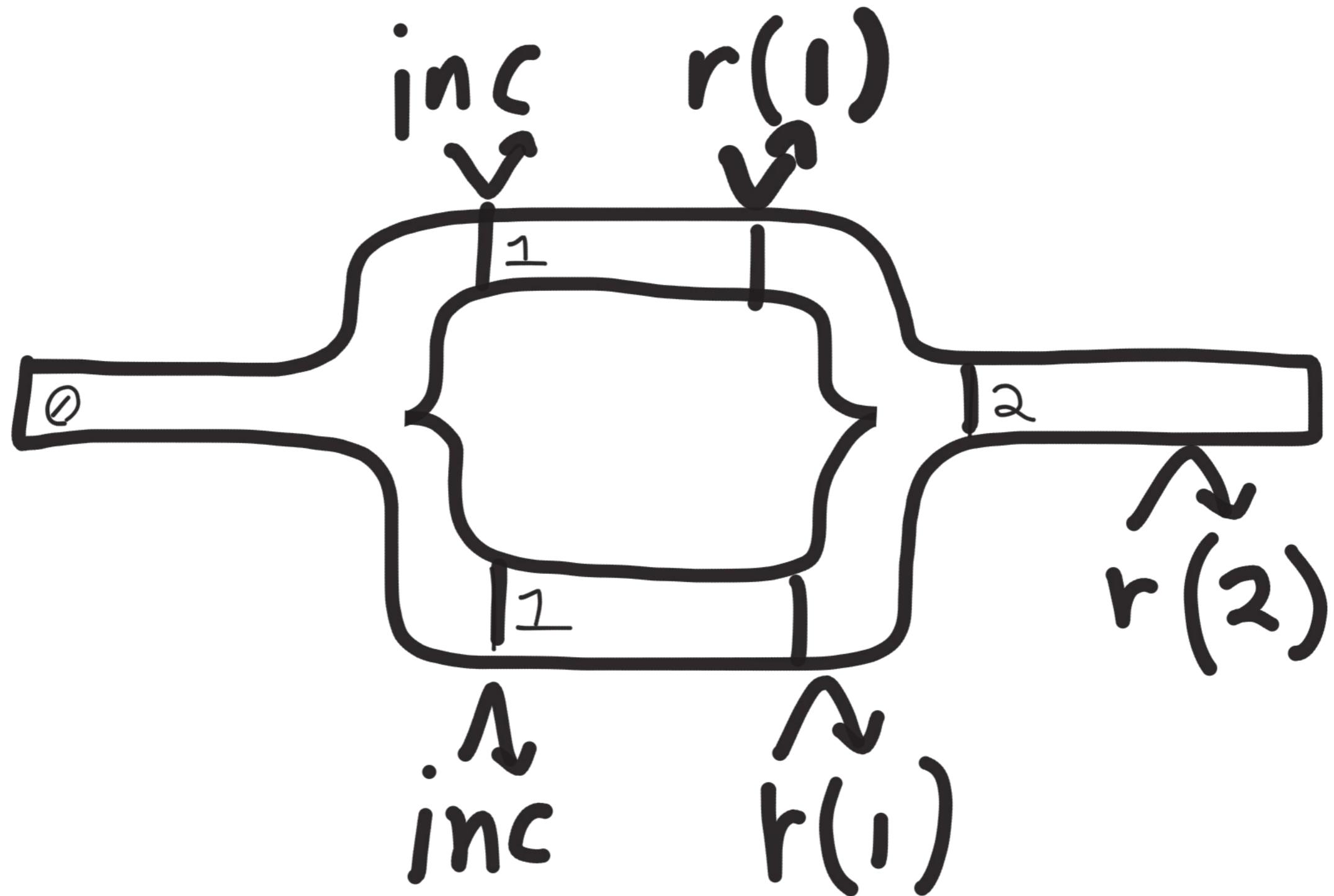
Linearizability



After op is complete,
guaranteed visibility!



Eventual Consistency



strong
linearizable
sequential

serializable

serializable

RR

SI

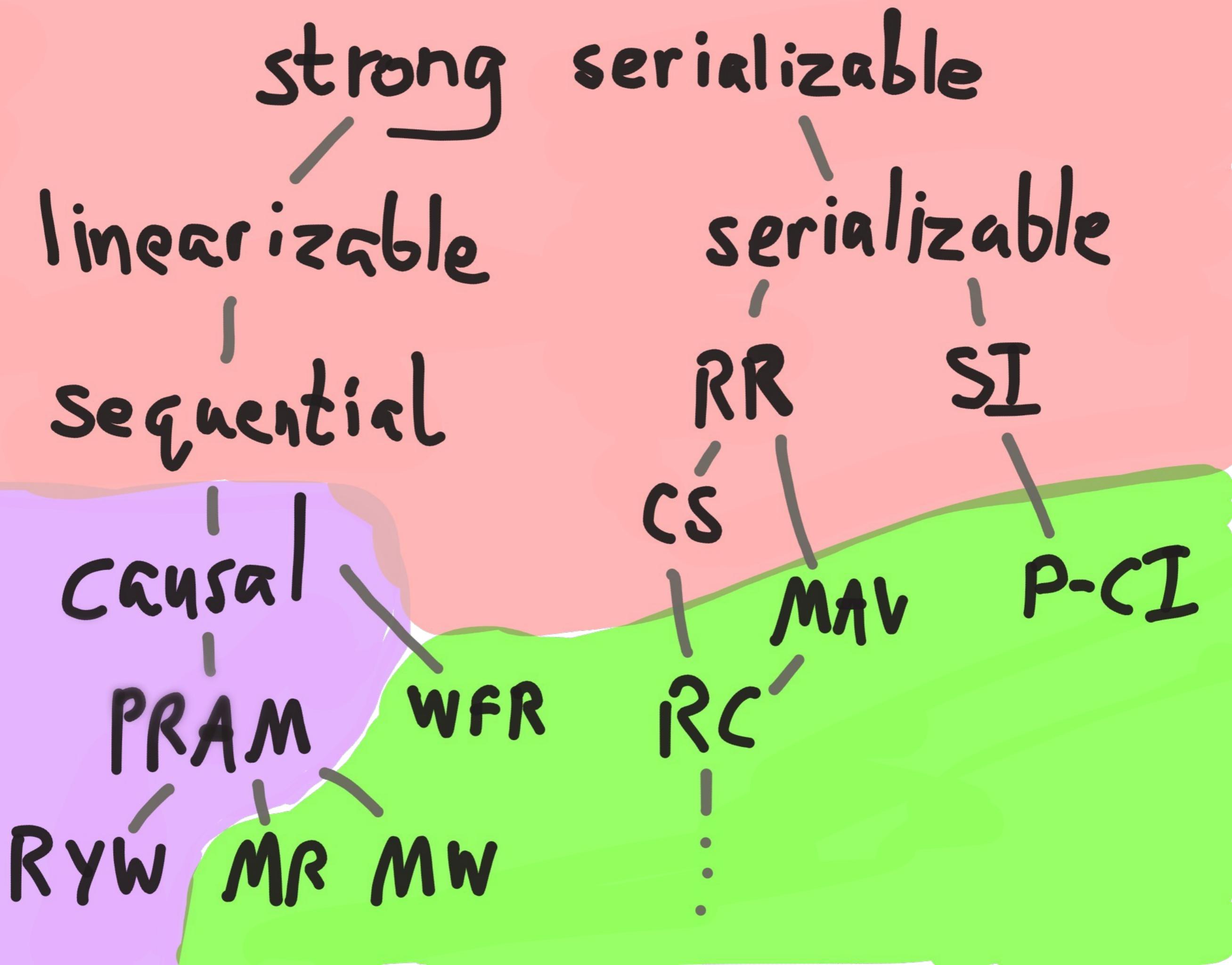
CS

P-CI

RC'

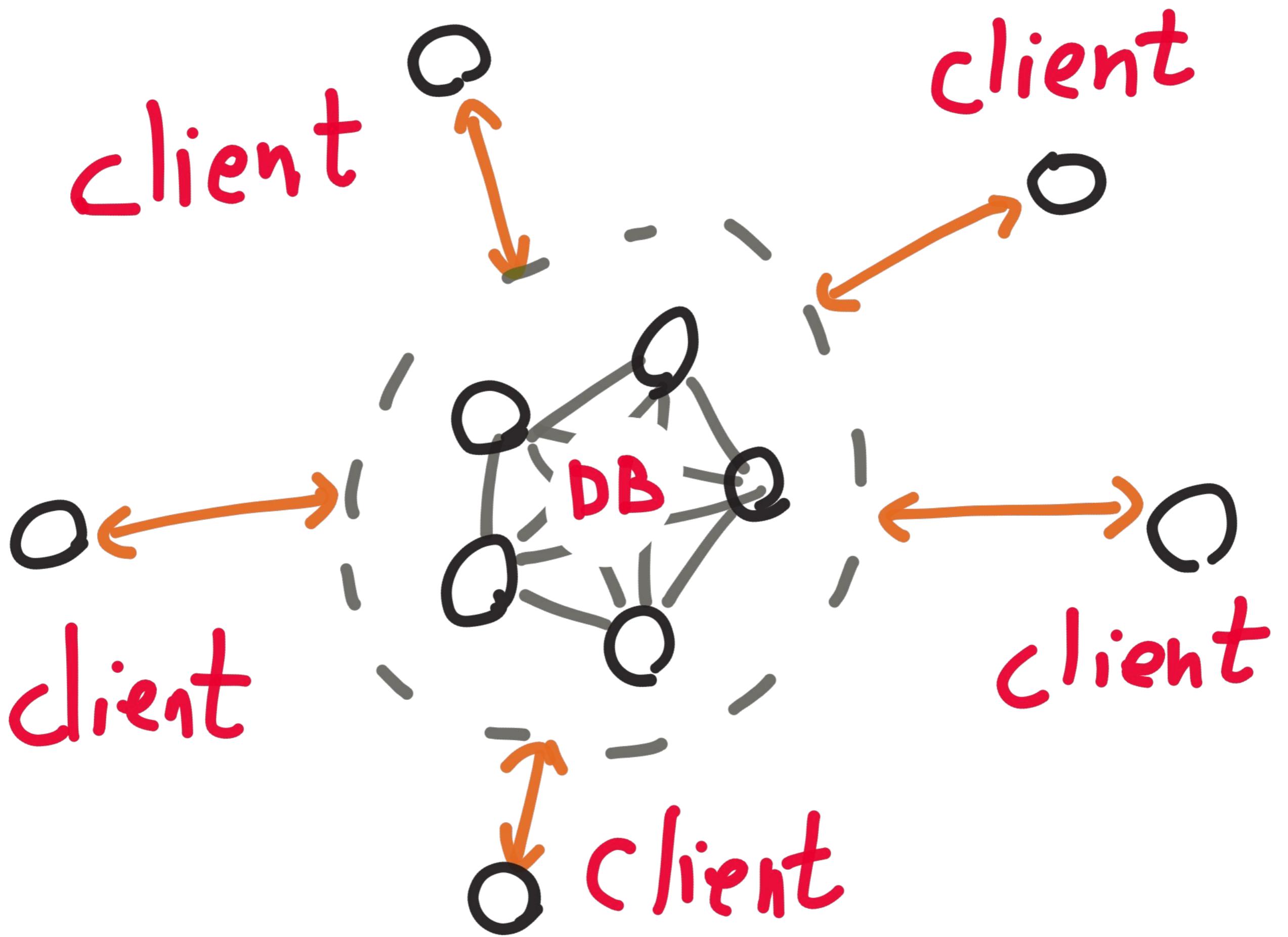
⋮

causal
PRAM
RYW MR MW
WFR



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 representing the database.

DB :

client : w — w' w — ...

The diagram shows the state of the database after the client requests have been processed. The database is represented by a central, blurred, cloud-like shape. A horizontal line segment connects the first two nodes, labeled w and w' , with a break symbol \dots indicating further nodes. Arrows point from the w node to the cloud and from the cloud to the w' node, illustrating the update or propagation of changes from the client to the database.

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

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

So, what
have you
found?



2013

- Redis Sentinel
- MongoDB
- Riak

2013

- Kafka
- NuoDB
- Cassandra

2014

- RabbitMQ
- etcd + Consul
- Elasticsearch

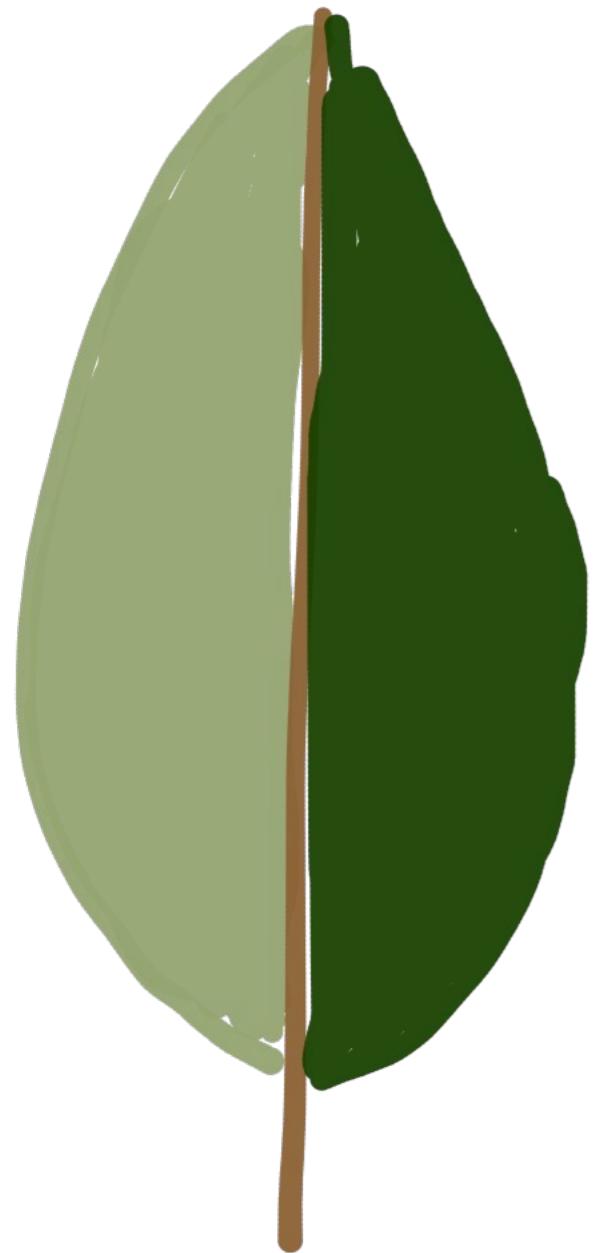
And now, for

something

COMPLETELY

~~DIFFERENT~~

THE SAME



mongoDB

First Jepsen test in
May 2013

~~#~~

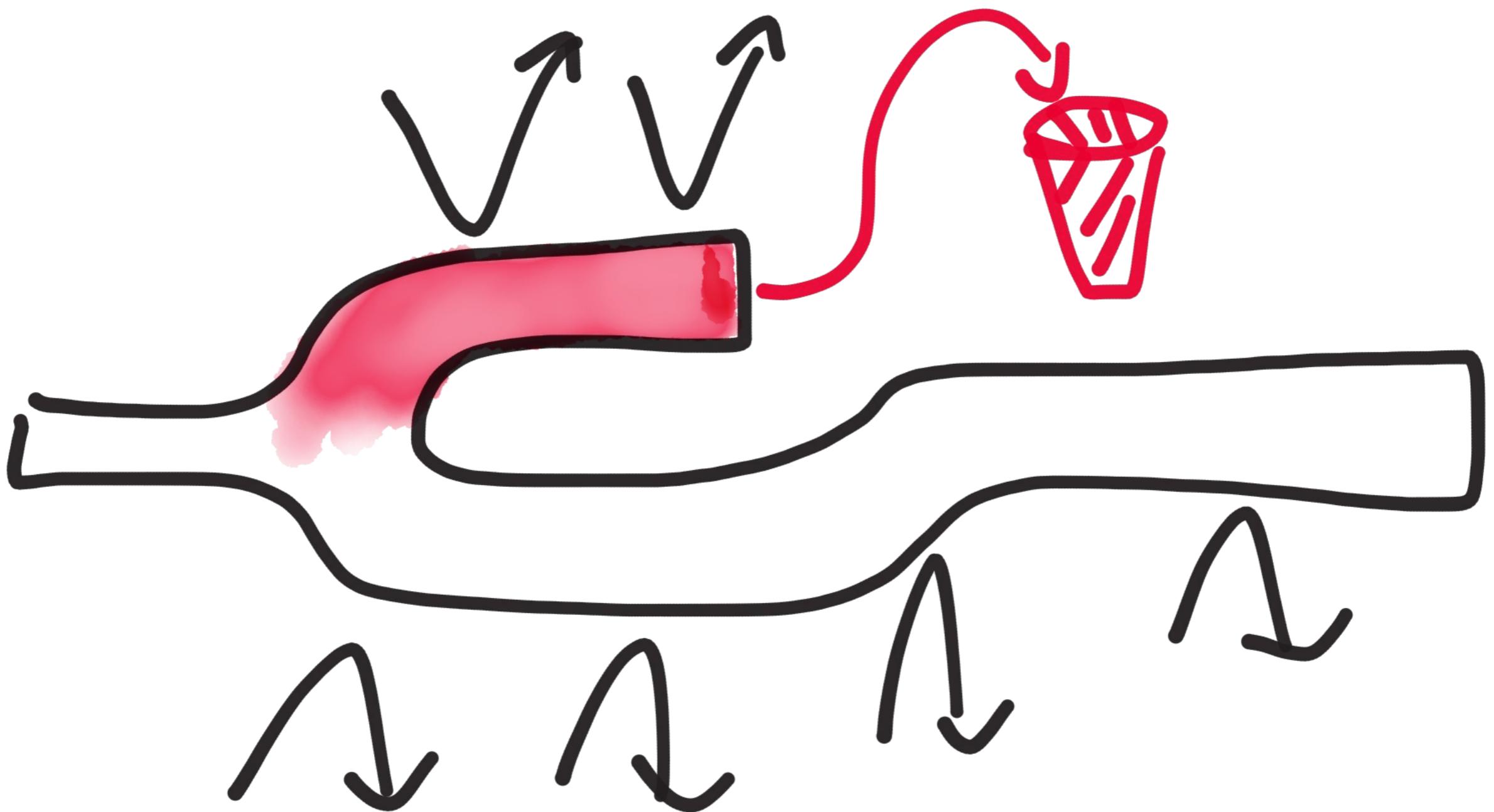
Data loss at all
WriteConcern levels

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

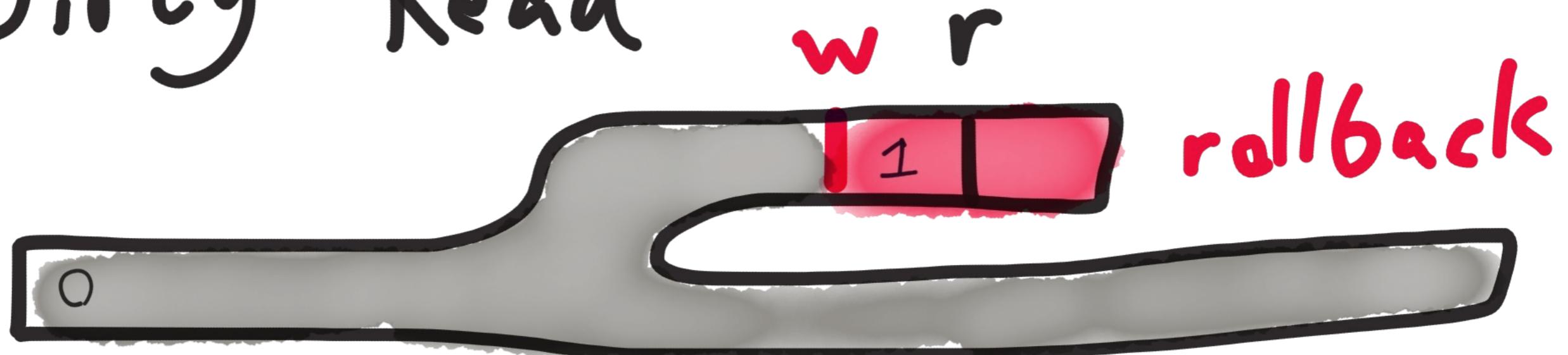
Linearizable

CaS Registers

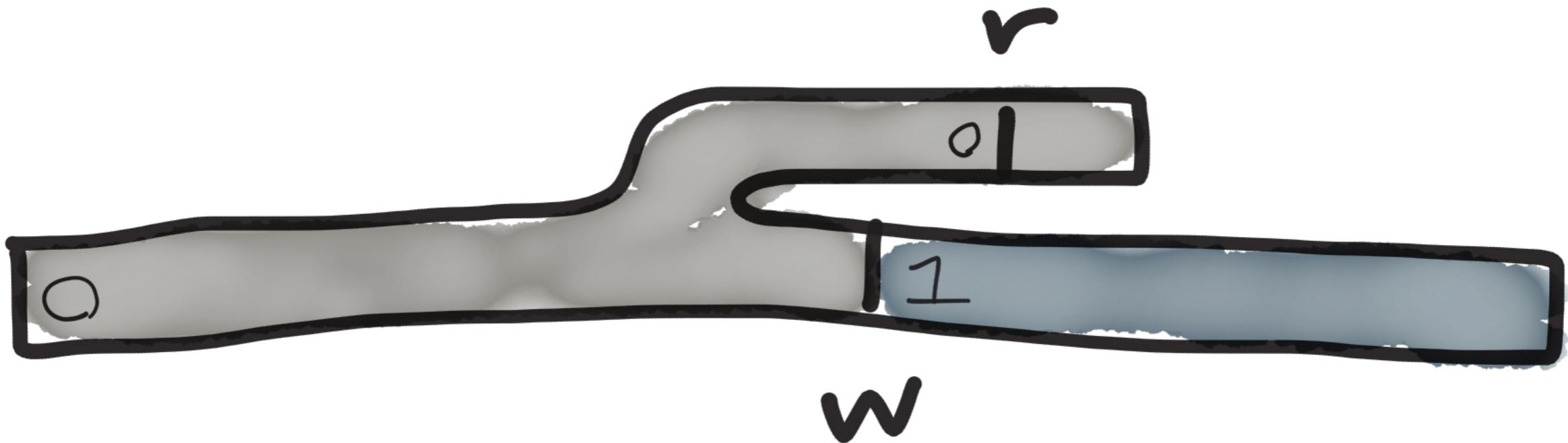
An Anomaly
Emerges!

47

Dirty Read



Stale Read

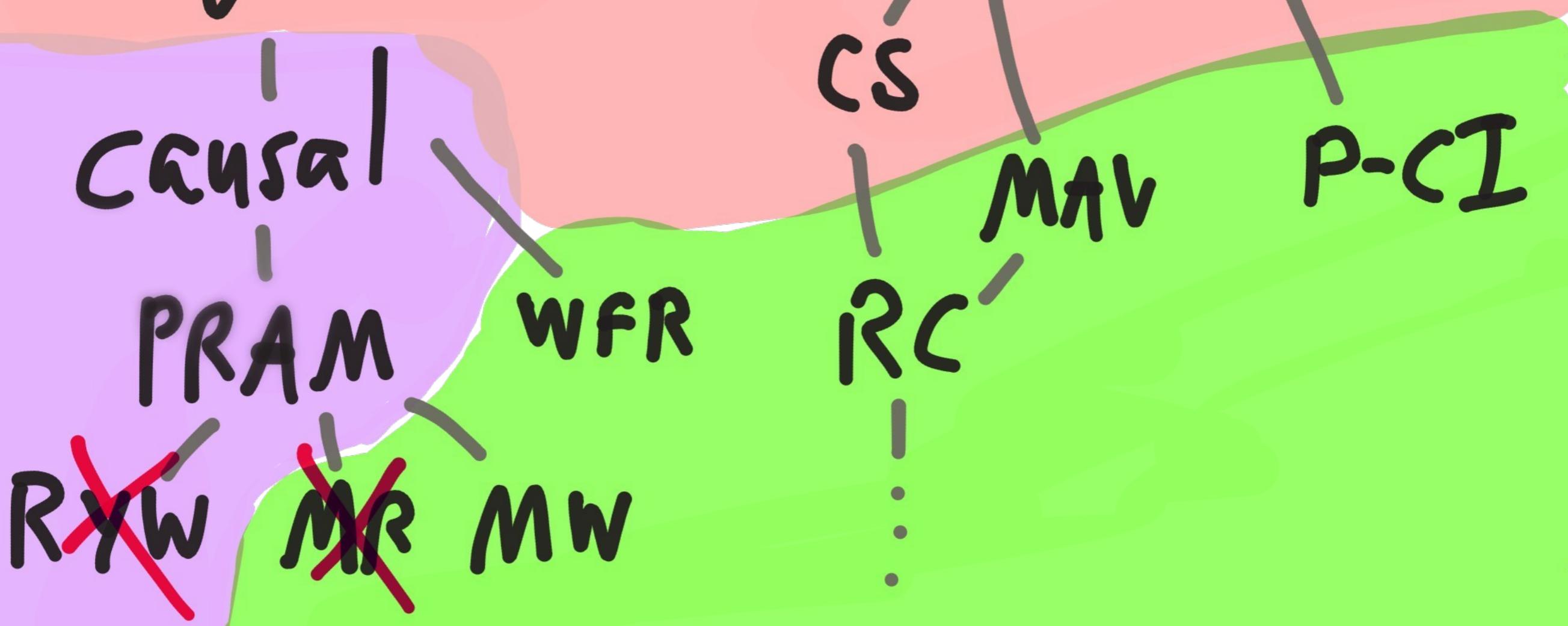


Server-17975

strong serializable

linearizable

sequential



~~strong~~ ~~serializable~~

~~linearizable~~

~~sequential~~

~~Causal~~

~~PRAM~~

~~RW~~ ~~MR~~ ~~MW~~

WFR

RR

CS

RC

MAV

SI

P-CI

⋮

~~strong~~ ~~serializable~~

~~linearizable~~

~~sequential~~

~~Causal~~

~~PRAM~~

~~RW~~ ~~MR~~ ~~MW~~

~~serializable~~

~~RR~~

~~CS~~

~~MLV~~

~~RC~~

~~RU~~

P-CI

WFR

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



~~strong~~ ~~serializable~~

~~linearizable~~

~~sequential~~

~~Causal~~

~~PRAM~~

~~RXYW~~

~~MR~~

~~MW~~

WFR

~~serializable~~

~~RR~~

~~CS~~

~~MAV~~

~~RC~~

RU

P-CI

Mongo
vs. docs

Mongo Recommendations

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

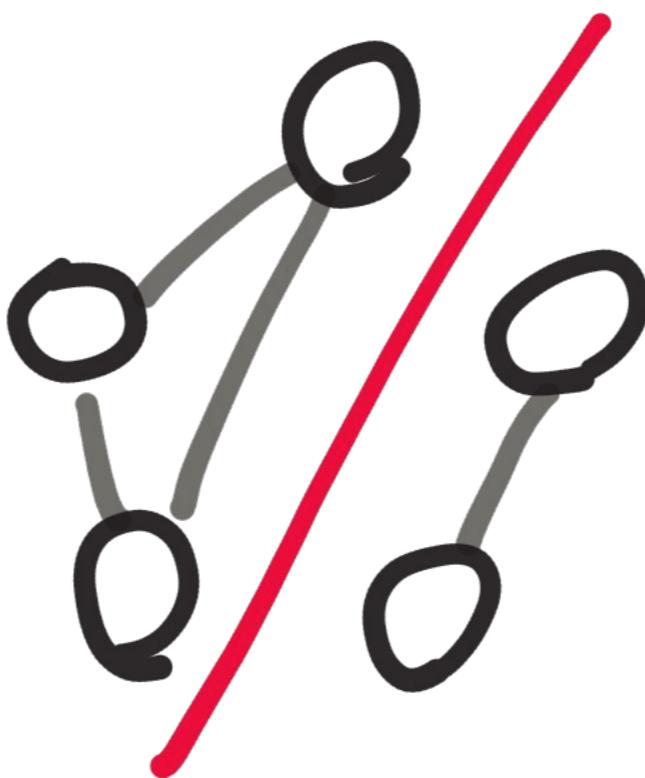
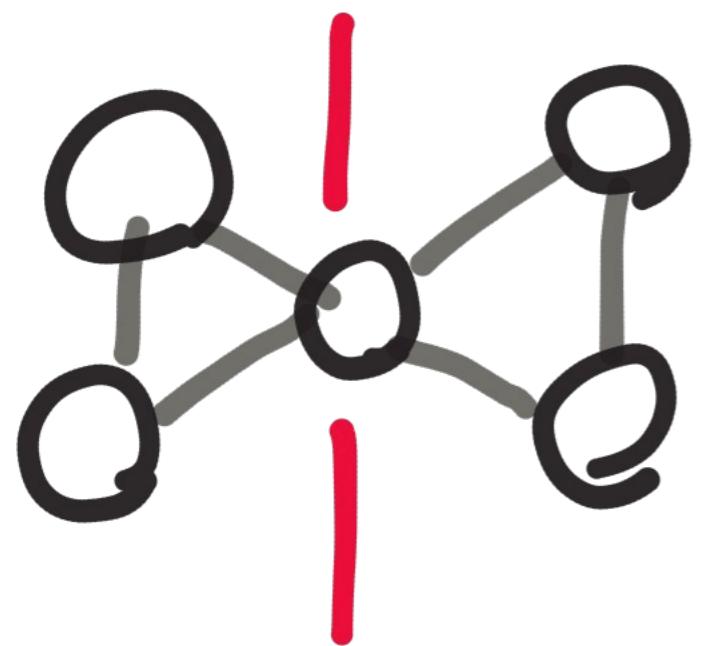


elasticsearch
[search]

First Elasticsearch
tests: June 2014

v 1.1.0

Split Brain from



Massive data loss: 90%.

- Cluster status lies
- 90 seconds for recovery
- Permanently wedged clusters

ES added a
great page on
safety issues



bleskes commented on Sep 1, 2014

Owner

I'm closing this issue, as it is solved, as specified, by the changes made in #7493. Of course, there is more work to be done and the effort continues.

Thx for all the input and discussion.



bleskes closed this on Sep 1, 2014



AeroNotix commented on Sep 1, 2014

@bleskes so it's 100% fixed?



bleskes commented on Sep 1, 2014

Owner

this issue (partial network splits causing split brain) is fixed now, yes.

Folks are
still referring
to the last talk



Mark Walkom

@warkolm



Follow

@falican durability in ES isn't a major problem these days, take a look at elasticsearch.org/guide/en/elast... for info on our view



9:36 PM - 20 Feb 2015

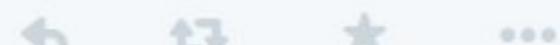


Reply to @warkolm @falican



Rory Hart @falican · Feb 20

@warkolm yeah I've read through that, looks like you're all working very hard. I've got an abundance of caution today. Appreciate the time.



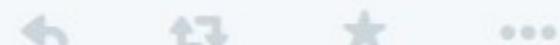
Mark Walkom @warkolm · Feb 20

@falican any time!



Mark Walkom @warkolm · Feb 20

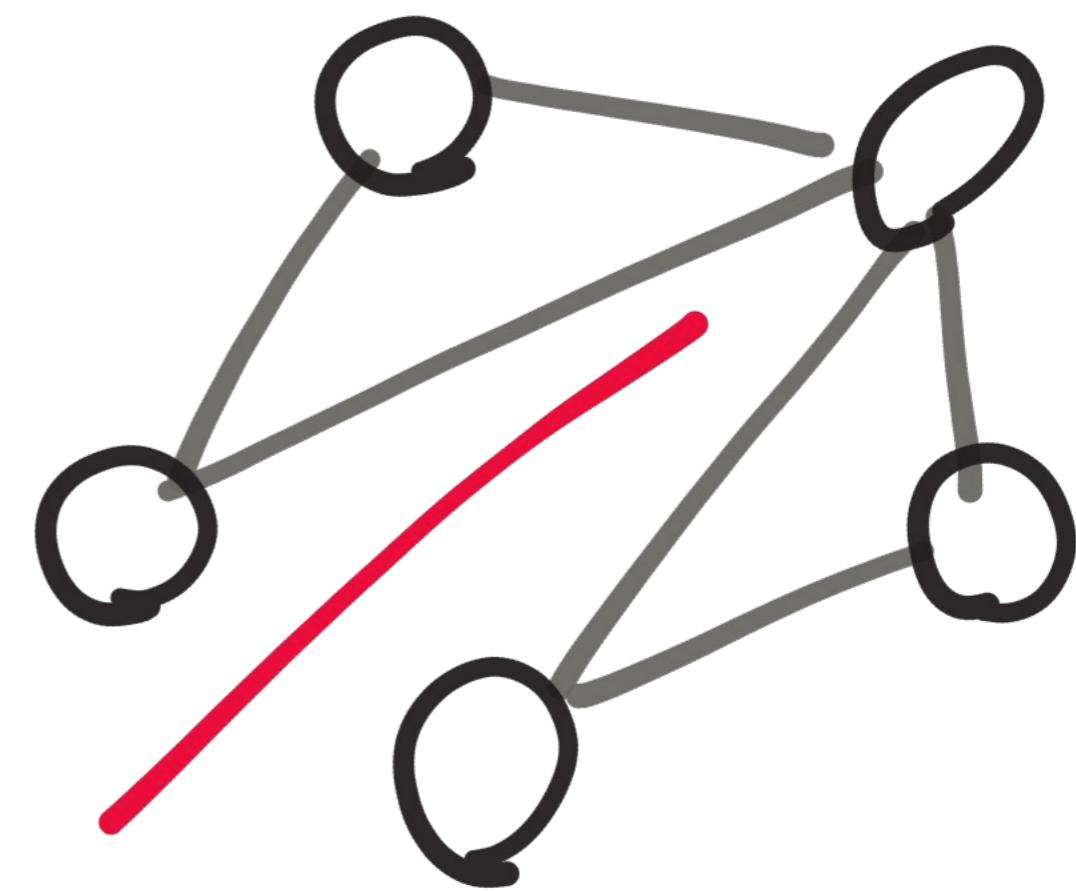
@falican re-read that article, most of the complaints in it are pre 1.0 btw, given we're probably not far off 1.5, grains of salt needed ;)



Retesting with

ES 1.50

Intersecting
partitions still cause
data loss



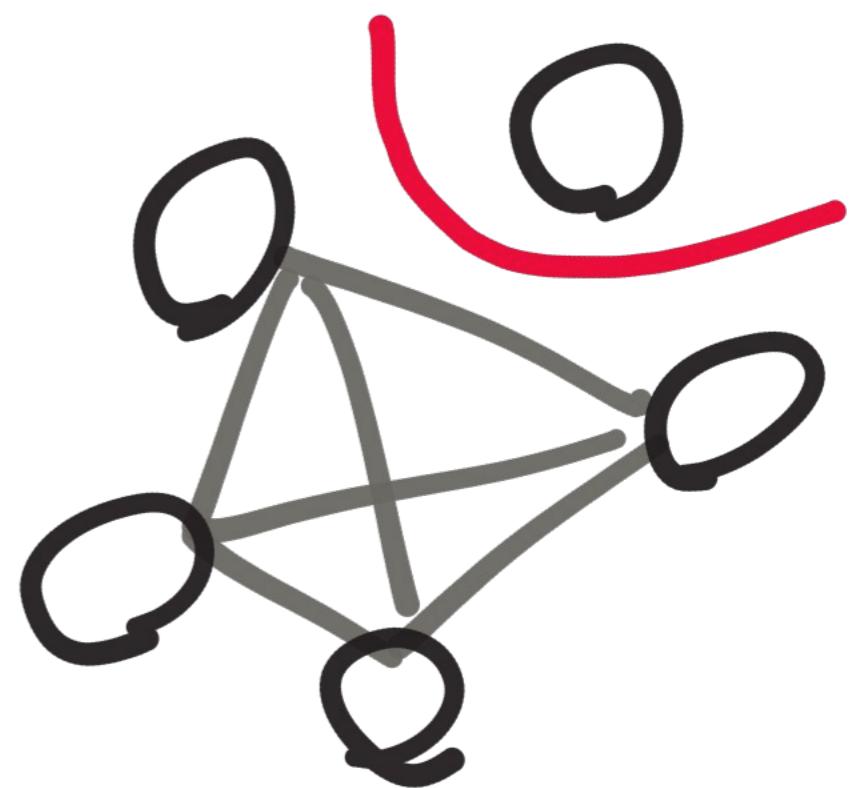
Long-lived split-

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
- 



AEROSPIKE



RELIABILITY

100% Uptime
with strong consistency (ACID)

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

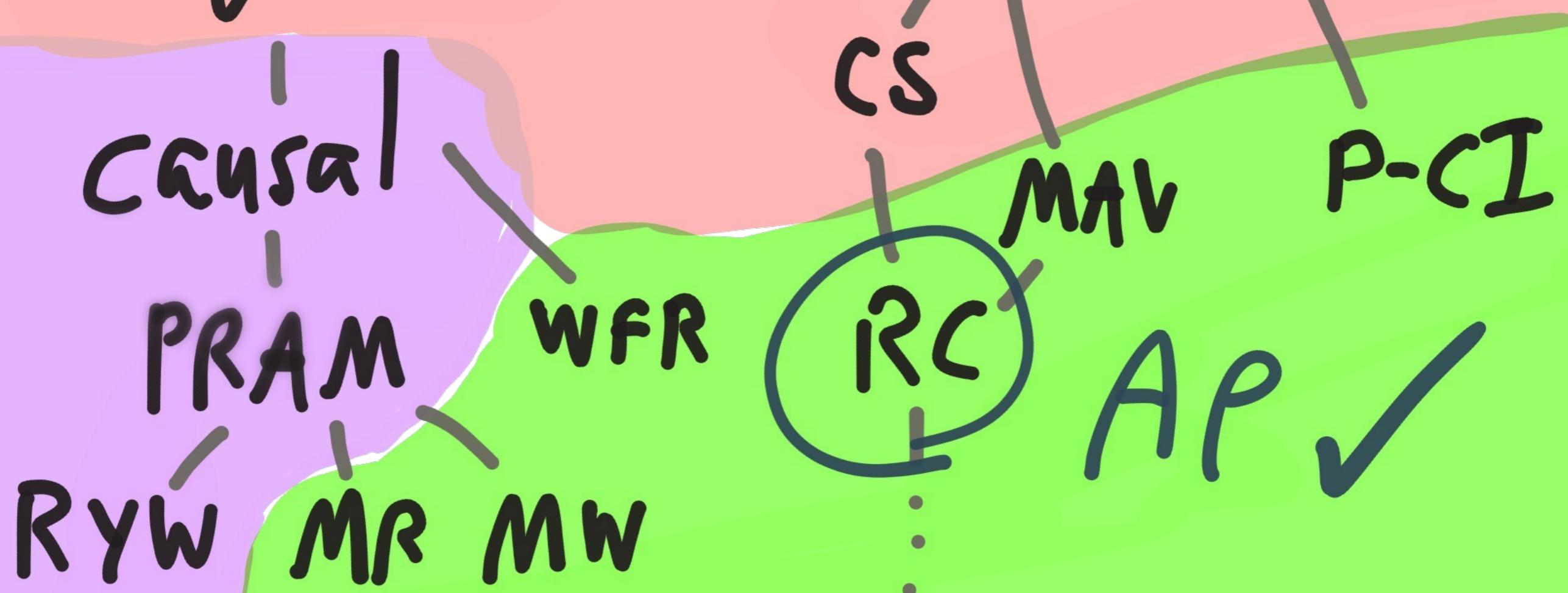
"AS is by and large an
AP system that provides
high consistency"

"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

WFR

MR MW

RR

CS

SI

MAV

RC'

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"



Search



DISCOVER

DOWNLOAD

DEVELOP

DEPLOY

DOCUMENTATION

COMMUNITY

Get Started!

DEPLOY

TRAINING

RESOURCES

**Aerospike is now available with Click-to-Deploy on Google Compute Engine**

You can quickly spin up an Aerospike cluster for development or production.

**Deploy Aerospike clusters easily in Amazon EC2**

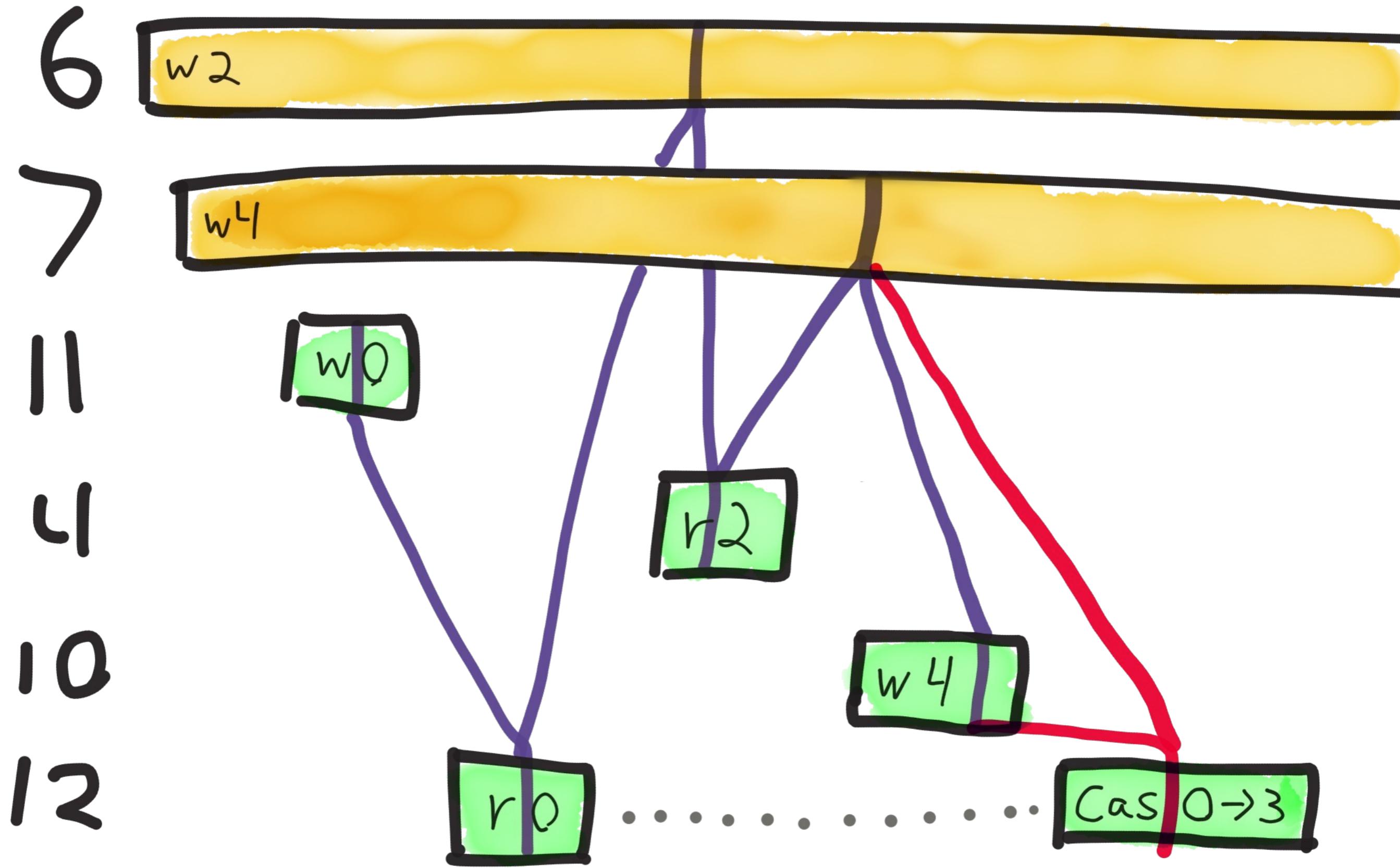
Learn how to install on EC2 instances, get tips on capacity planning, configuration, benchmarking and using bcache for faster persistent storage.

**1M TPS on 1 Server**

This can be achieved on a \$5k commodity server. Don't believe it? Follow this easy, 4 step

So for a CAS

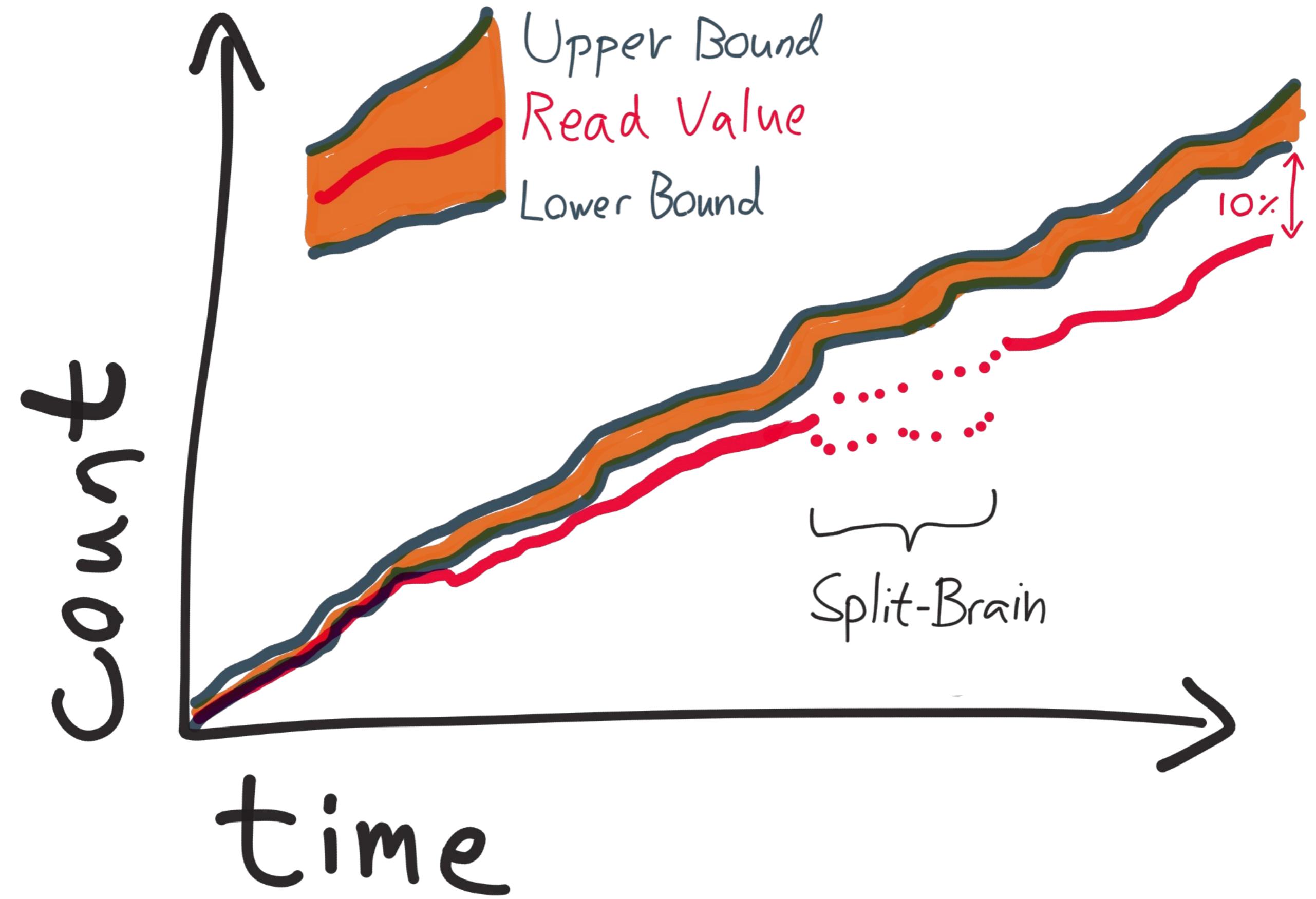
Register...



Counters will

drop increments

during a partition



"Deployment modes

AP vs. CP"



CP is vaporware

Application Merge!

"Read... will return
both versions, allowing
the app. to resolve"

Like “CP mode”

this feature does

not exist.

No CRDTs for you!

UPDATES

~~strong serializable~~

linearizable

~~sequential~~

PRAM

RYW MR MW

~~serializable~~

RR

CS

1

20

4

May



P-CI

Recap



Mongo & EJ

have made

improvements!



Still Surprises

Lurking in the



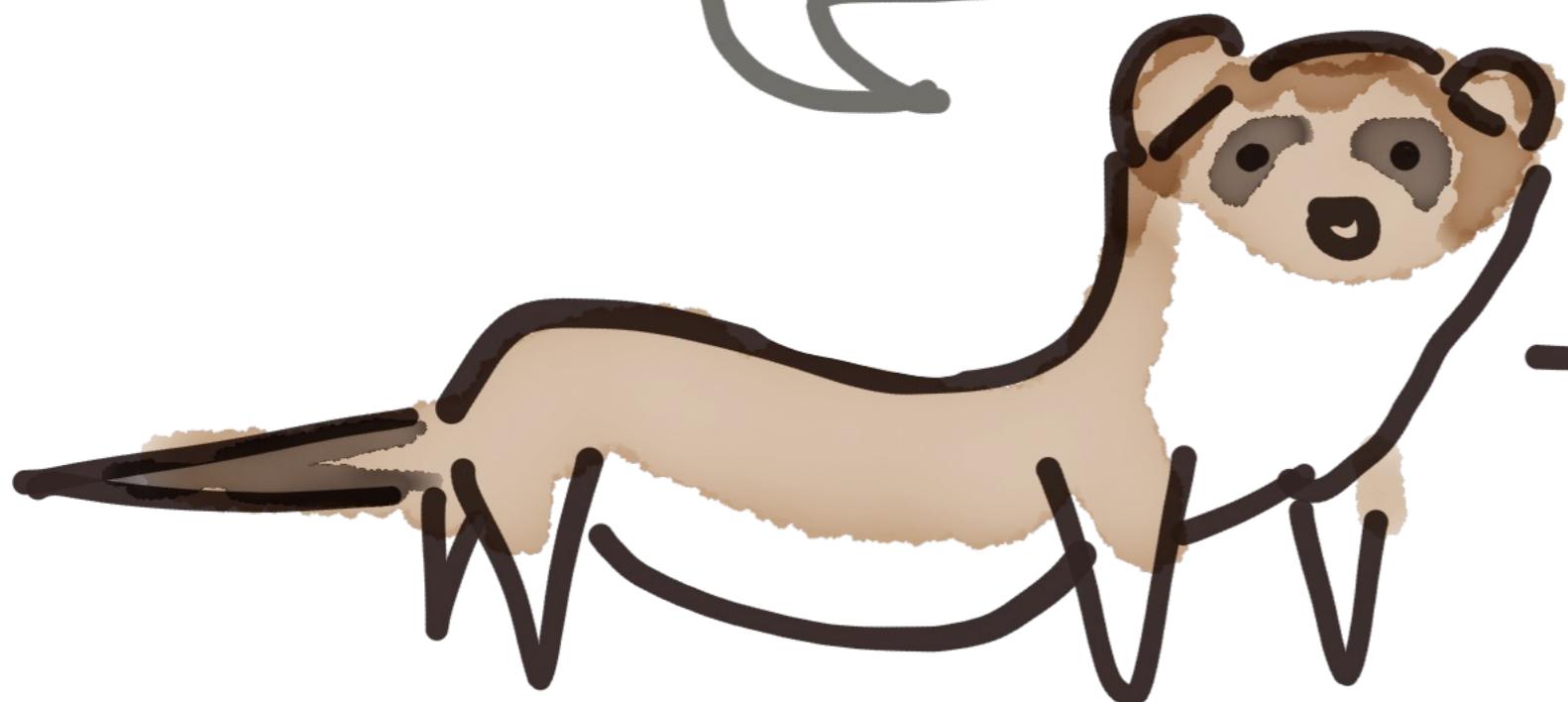
Read the docs!



"Strict"

"ACID"

"Strong"



- WORDS

*Figure out the
invariants
your system needs*

How much



can you tolerate?

Then test it

— for —

YOURSELF

Consider
your
failure modes

Process Crash

#Kill -9 1234

Node failure

- AWS terminate
- Physical power switch

Clock Skew

date 1028000

faketime ...

GC/TQ Pause

killall -s STOP foo

killall -s CONT foo

Network Partition

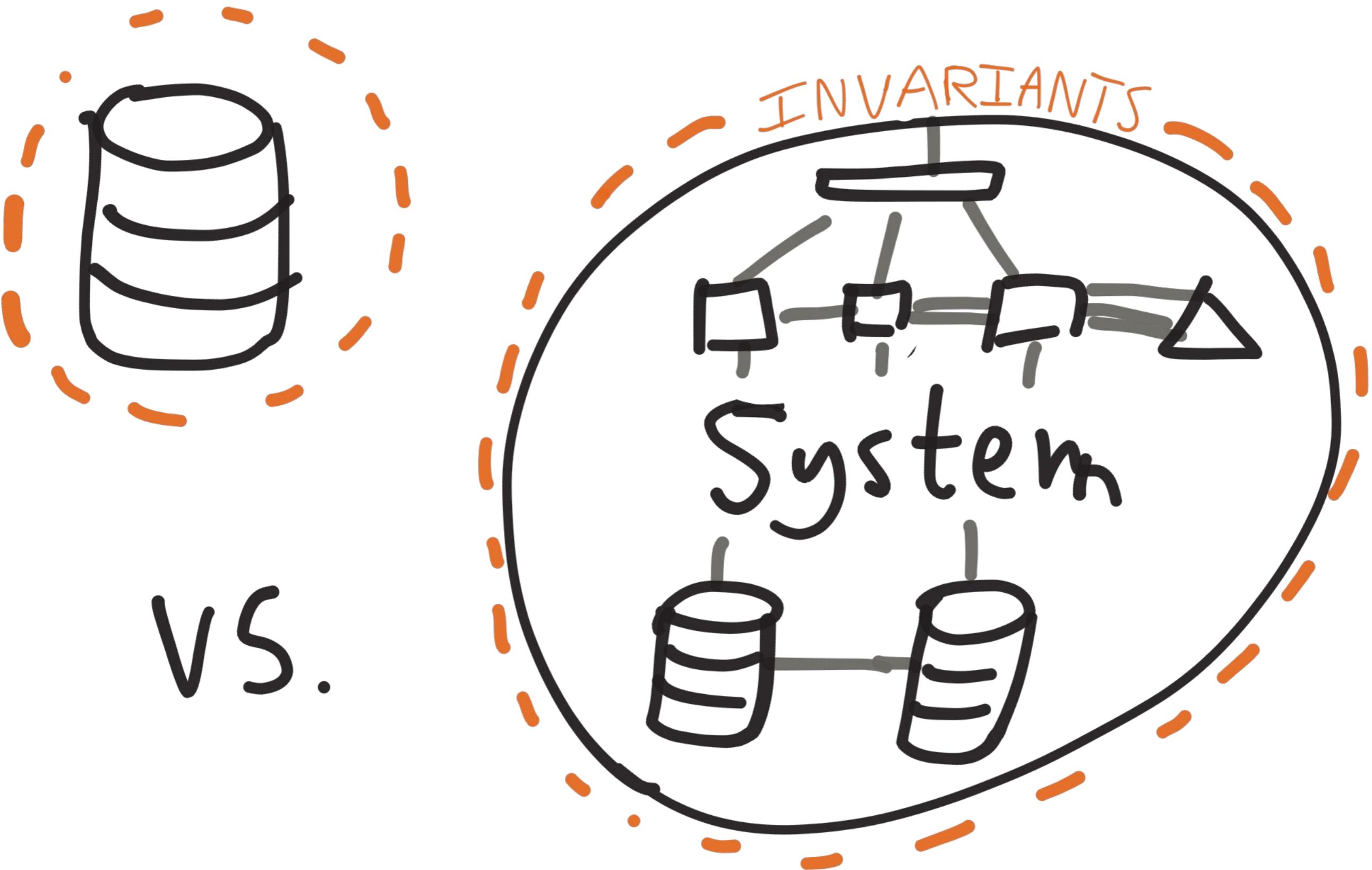
iptables -j DROP

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

Test your

Systems

end-to-end



vs.

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

Thanks!

Kevin Porter

Lucien Valmar

Boaz Leskes

Lee Hinman

Matt Kangas

Stripe

Evan Broder

Asya Kamsky

Kelly Stirman

Marc Hedlund

Jepsen

github.com/aphyr/jepsen