# Dynamo concepts in depth.
## Pavlo Baron, codecentric AG

Pavlo Baron

pavlo.baron@codecentric.de
@pavlobaron

So Dynamo isn't about speed.
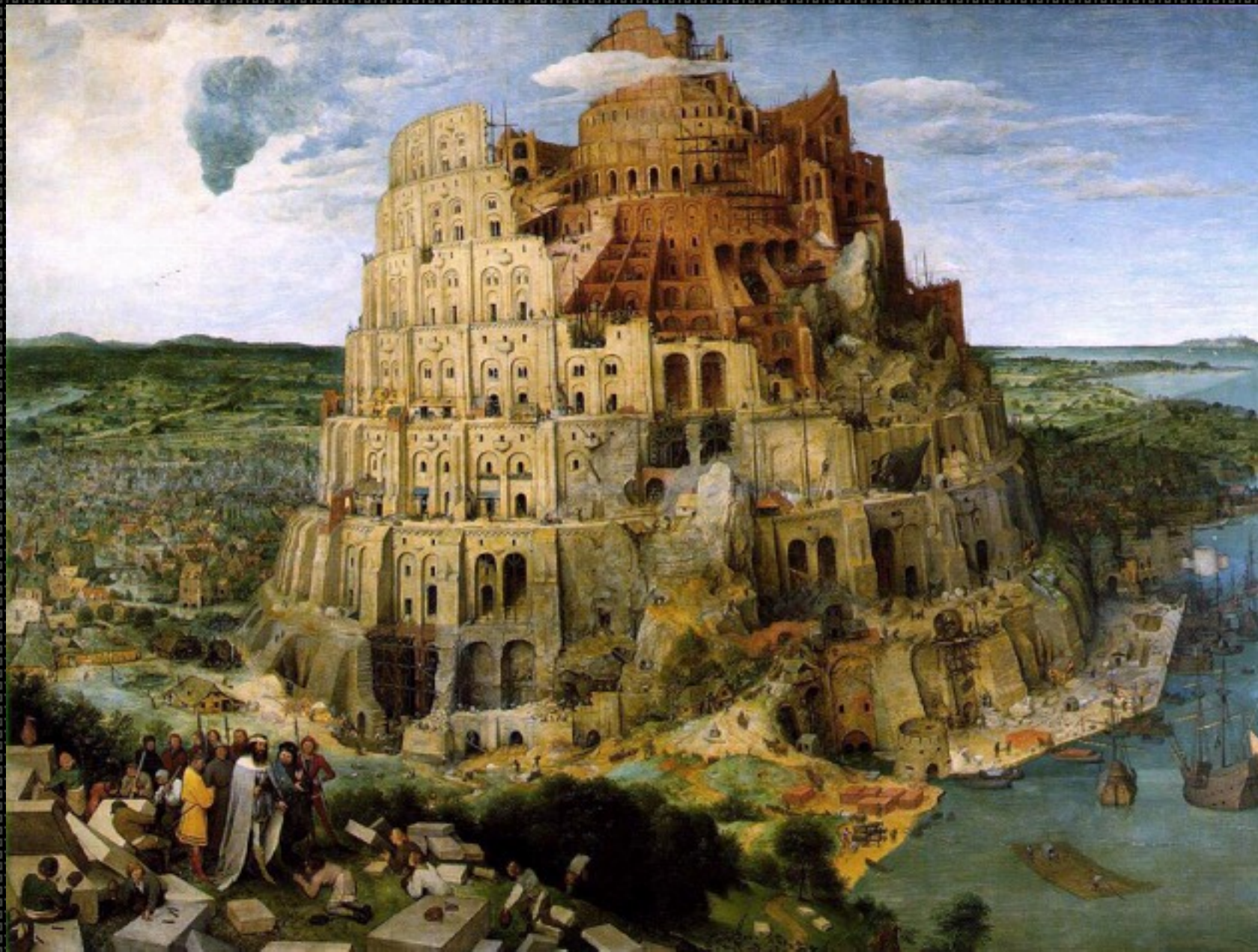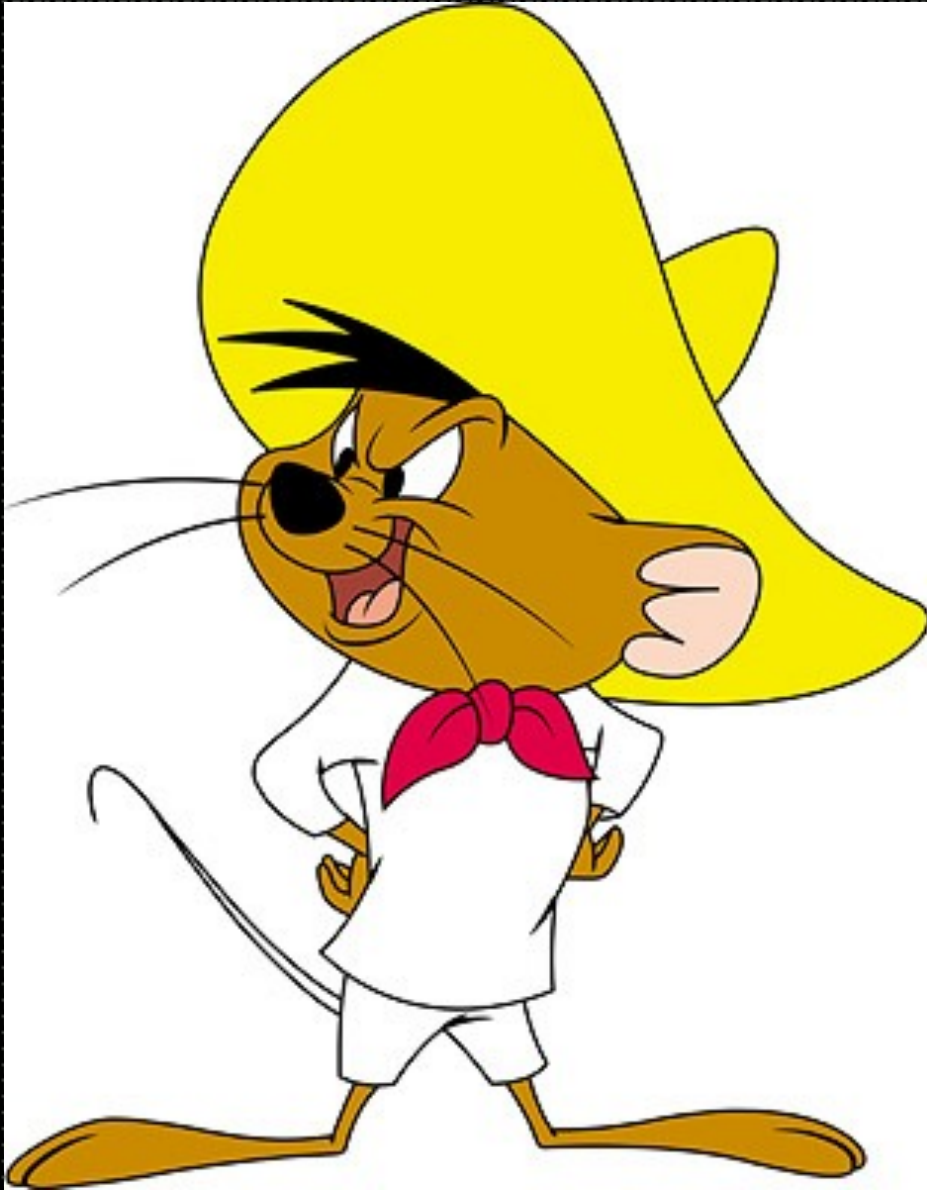
It's about immediate, reliable writes.

It's about operation relaxation.

It's about distribution and fault tolerance.
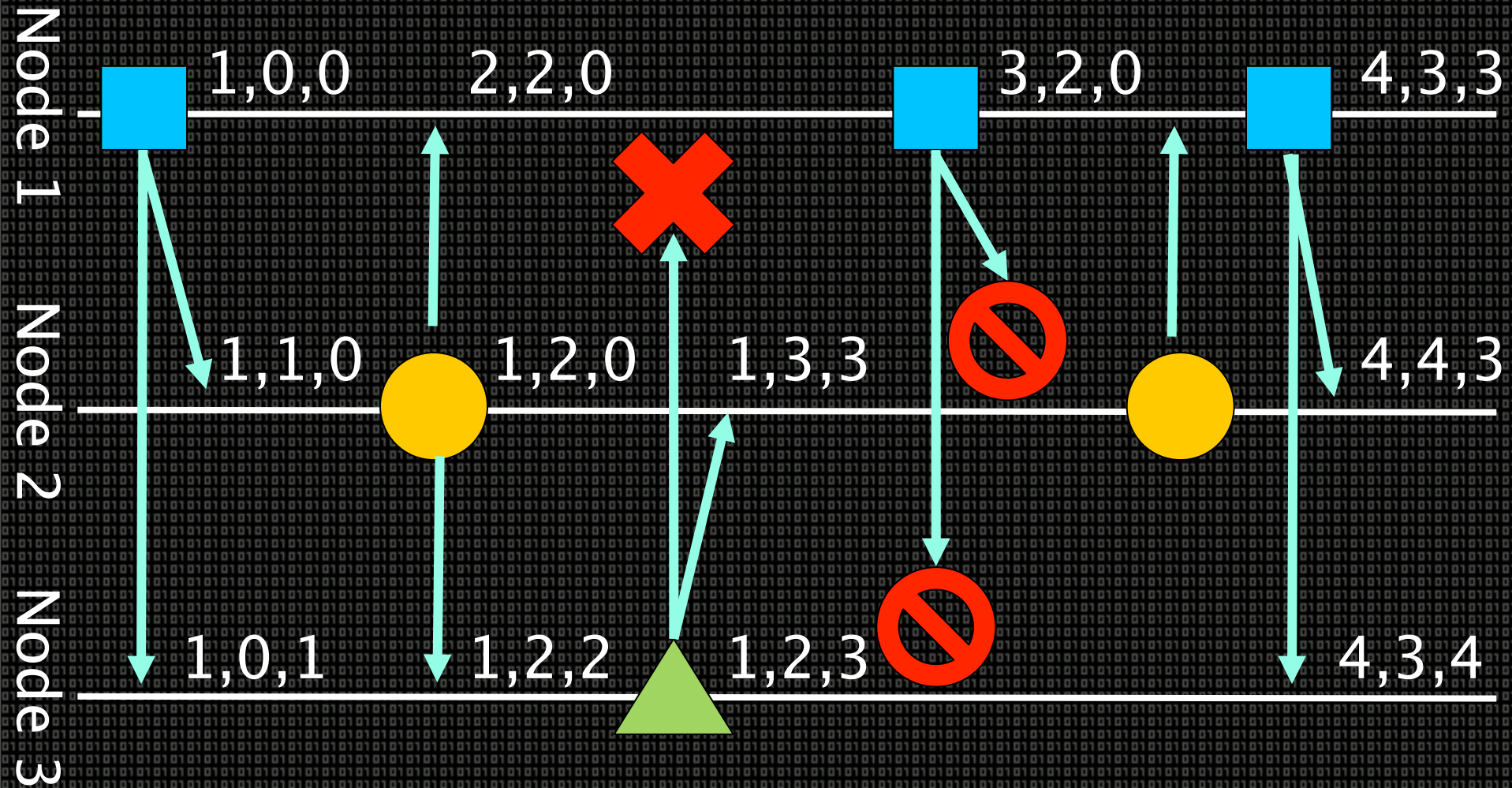
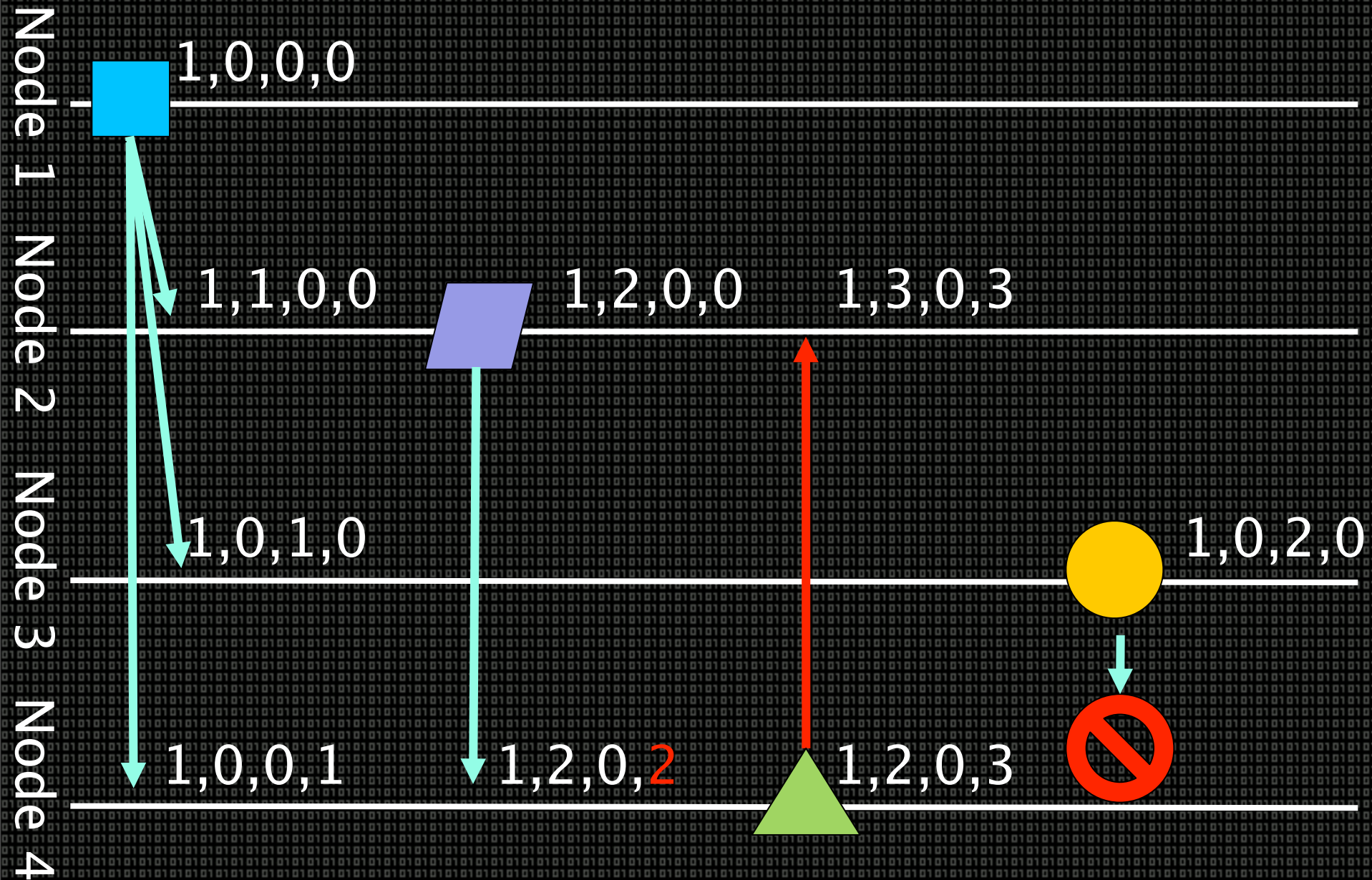It's about almost linear scalability.

V(i), V(j): competing

Conflict resolution:

1: siblings, client
2: merge, system
3: voting, system

Vector clocks

Node 1    1,0,0    2,2,0    3,2,0    4,3,3

Node 2    1,1,0    1,2,0    1,3,3    4,4,3

Node 3    1,0,1    1,2,2    1,2,3    4,3,4

# Vector clocks

Node 1    1,0,0,0

Node 2    1,1,0,0        1,2,0,0        1,3,0,3

Node 3    1,0,1,0                                    1,0,2,0

Node 4    1,0,0,1        1,2,0,2        1,2,0,3

# O(1) for data lookups / delta tracking

#

N, M: nodes
HT(N), HT(M): hash trees

M needs update:
    obtain HT(N)
    calc delta(HT(M), HT(N))
    pull keys(delta)

a

ab    ac

abc   abd   acb   acc

abe   abd   ada   adb

ab   ad

a

Node a.2

# "Equal" nodes based decentralized distribution

# Consensus, agreement, voting, quorum

X bit integer space
$$0 <= N <= 2 \wedge X$$

or: 2 x Pi
$$0 <= A <= 2 \times Pi$$
$$x(N) = \cos(A)$$
$$y(N) = \sin(A)$$

V: vnodes holding a key
W: write quorum
R: read quorum
DW: durable write quorum
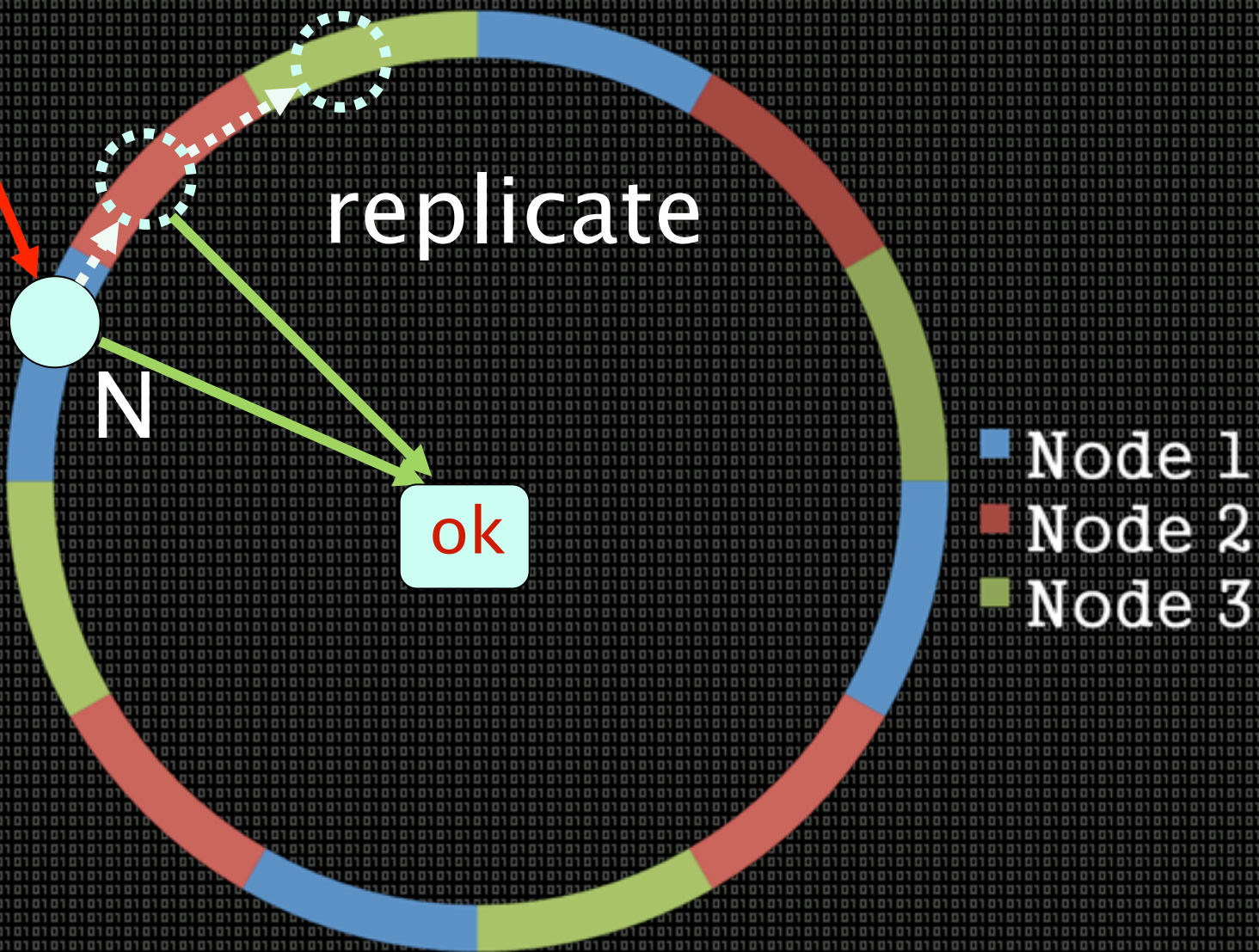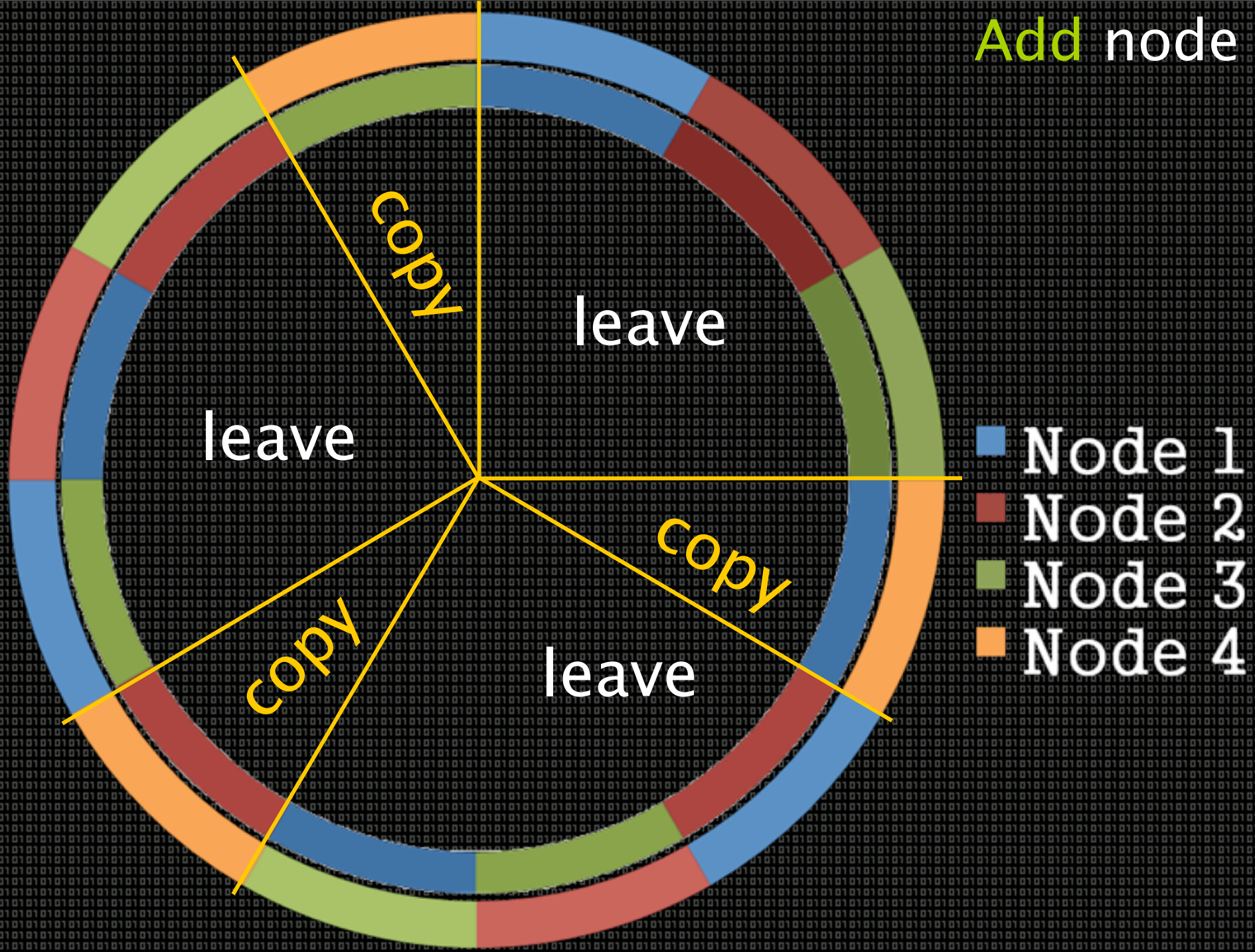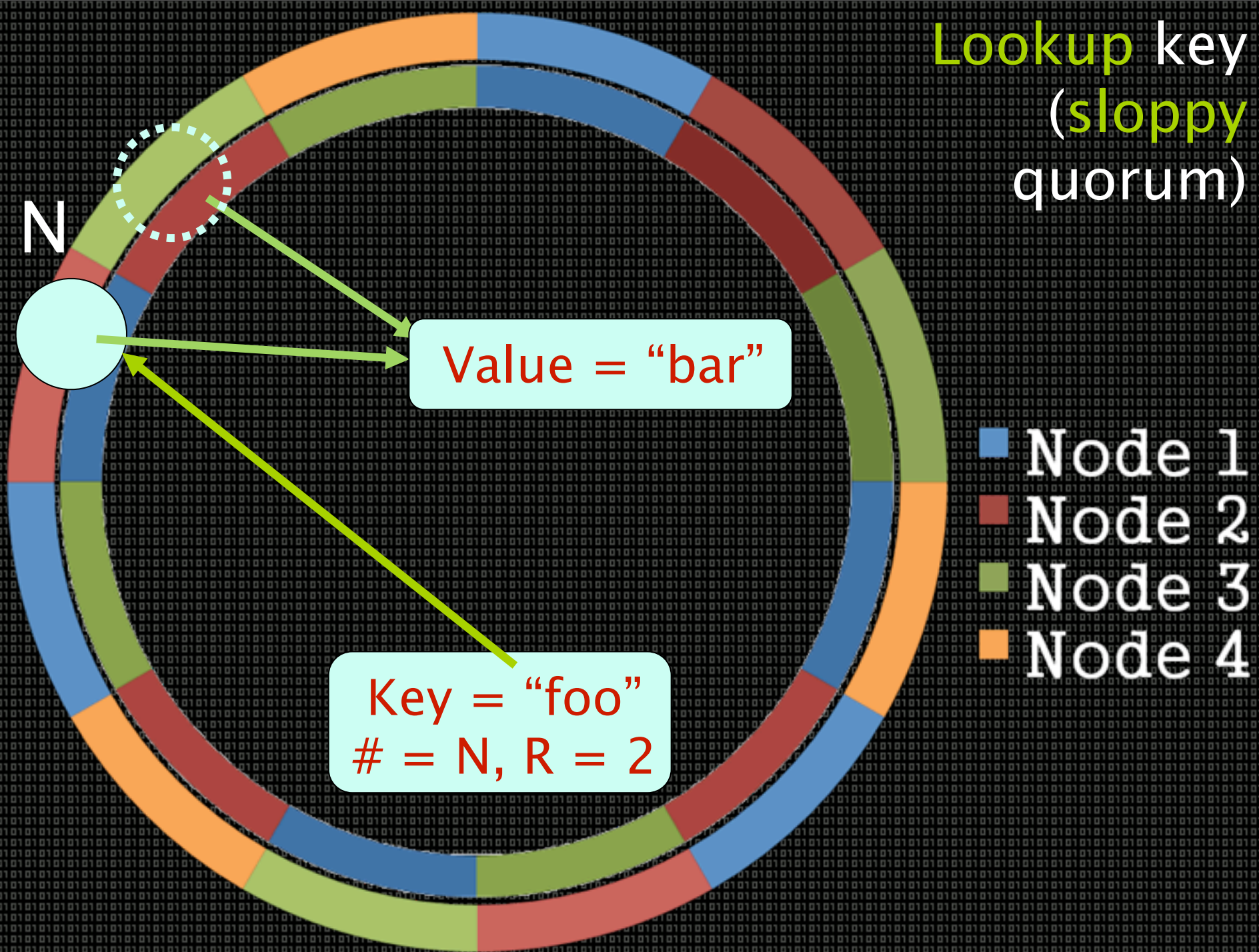
$$W > 0.5 * V$$
$$R + W > V$$

Insert key
(sloppy quorum)

Key = "foo"
# = N, W = 2

replicate

N

ok

Node 1
Node 2
Node 3

Lookup key (sloppy quorum)

N

Value = "bar"

Key = "foo"
# = N, R = 2

Node 1
Node 2
Node 3
Node 4

Remove node

copy

leave

Node 1
Node 2
Node 3

Gossip – node down/up

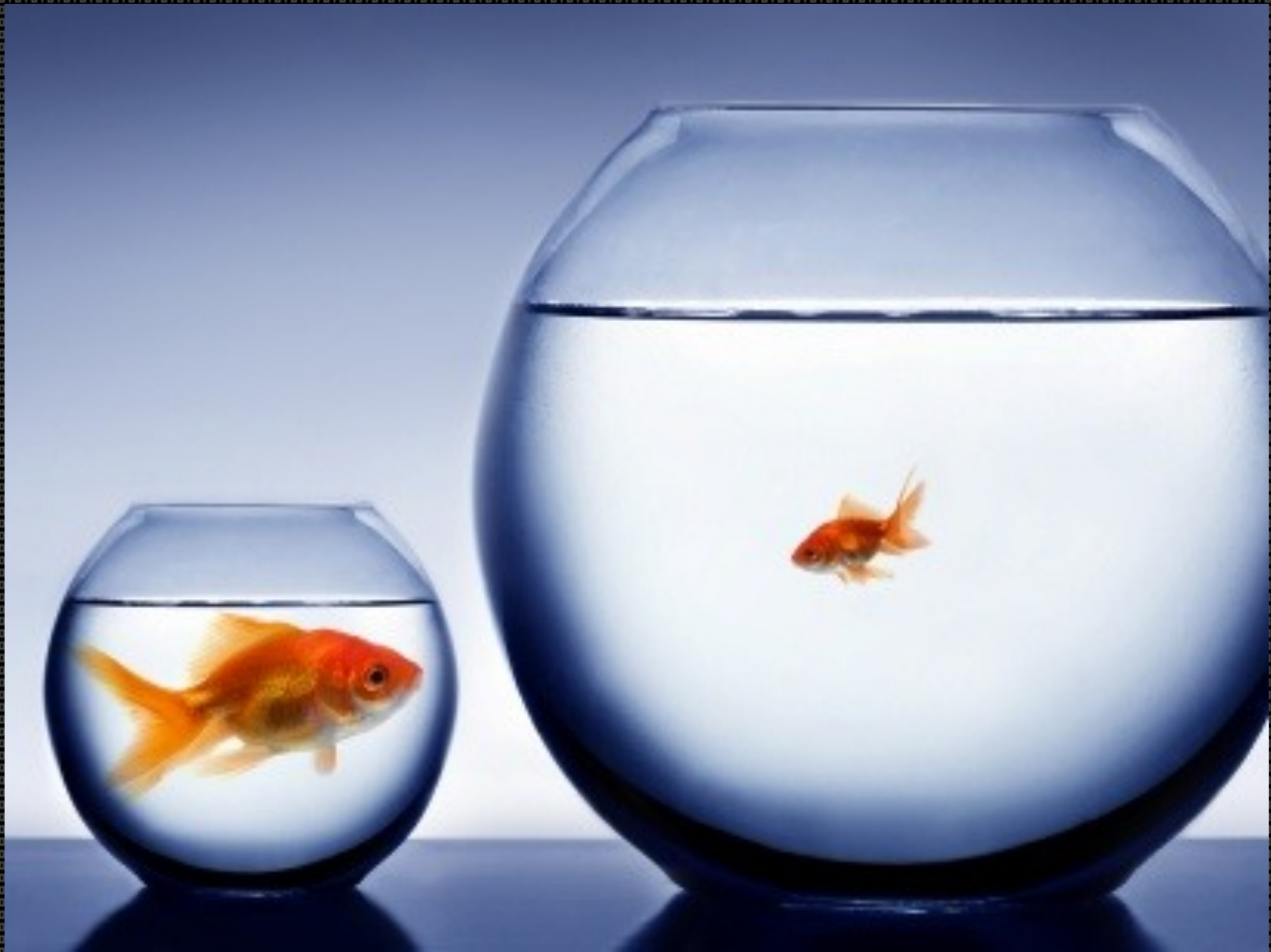Node 1
Node 2
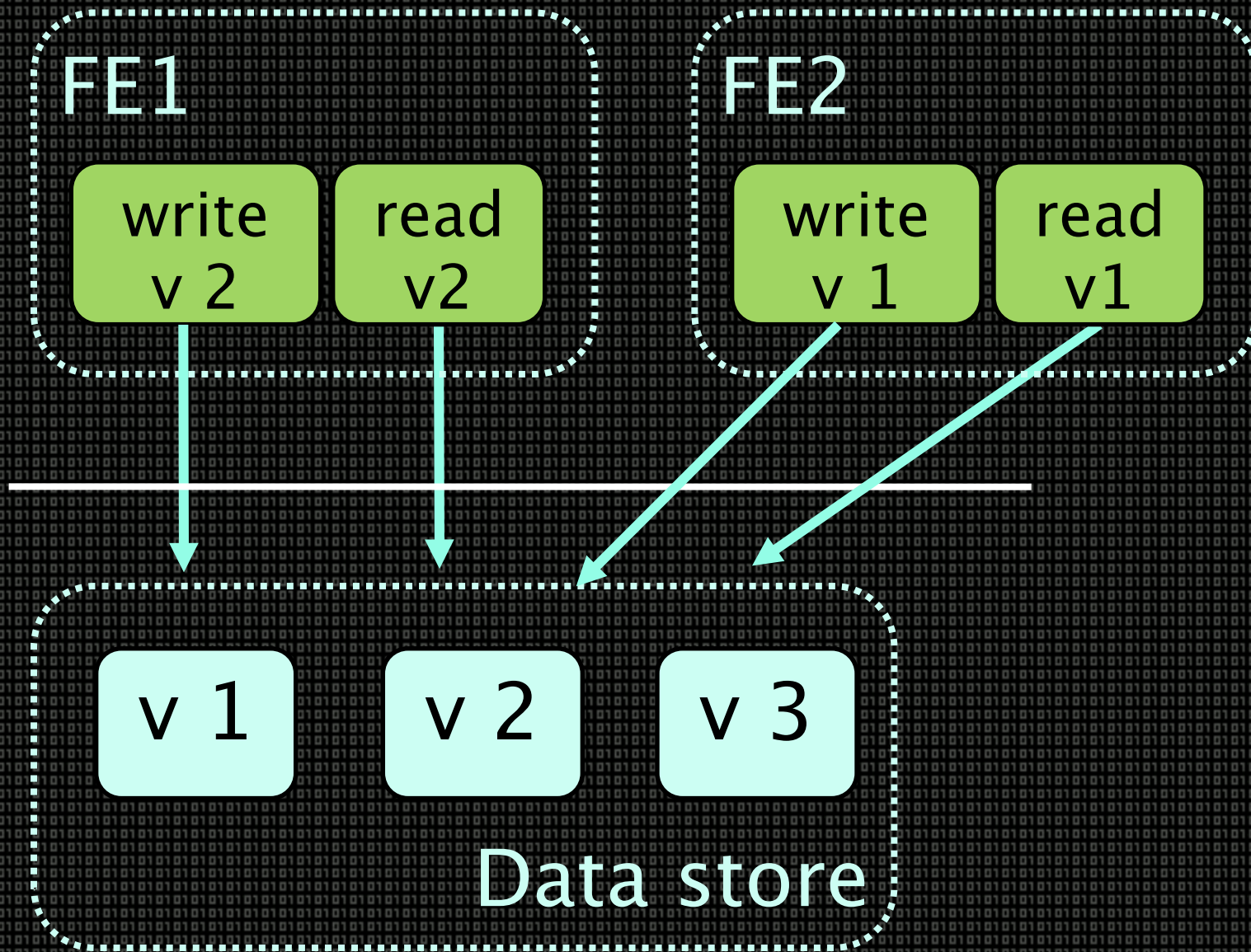Node 3
Node 4

update
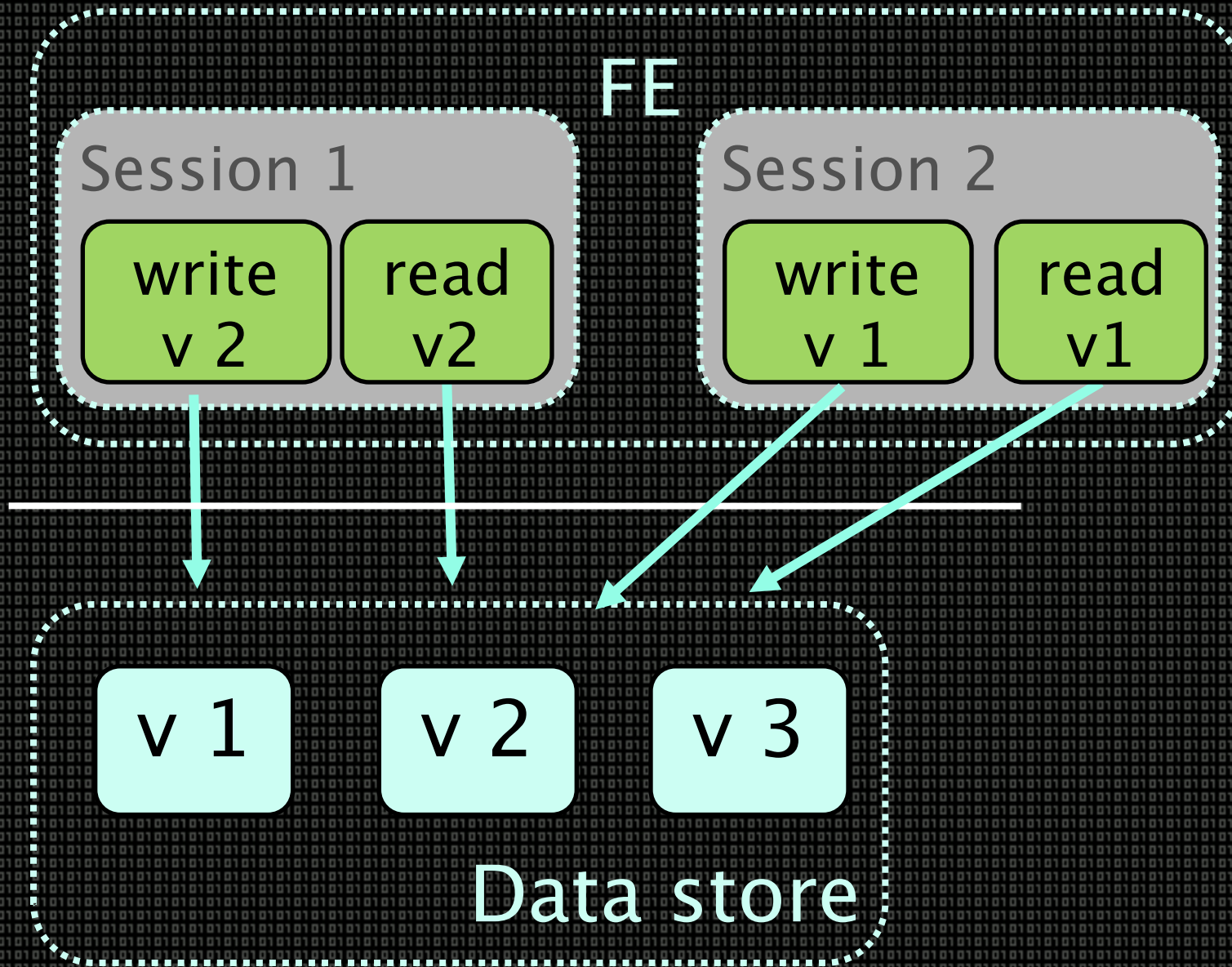
update, 4 down

update

update

read, 4 up

read

**B**asically **A**vailable,
**S**oft-state,
**E**ventually consistent

Opposite to ACID

Read your write consistency

Monotonic read consistency

FE1
read v 2
read v2
read v3

FE2
read v 3
read v4

Data store
v 1
v 2
v 3
v 4

N: node, G: group including N

node(N) is unavailable
      replicate to G or
      store data(N) locally
      hint handoff for later
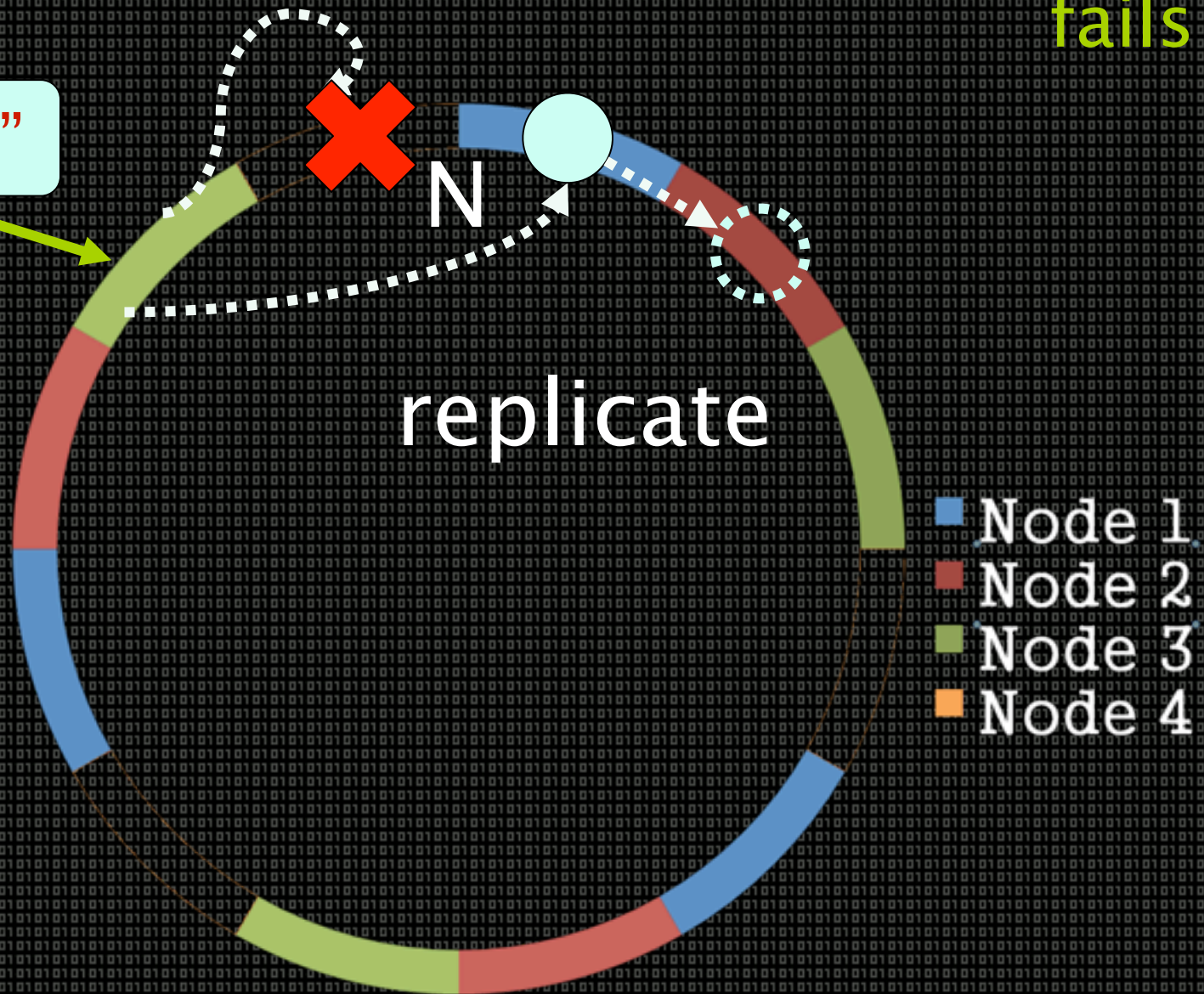node(N) is alive
      handoff data to node(N)

# handoff

Replica
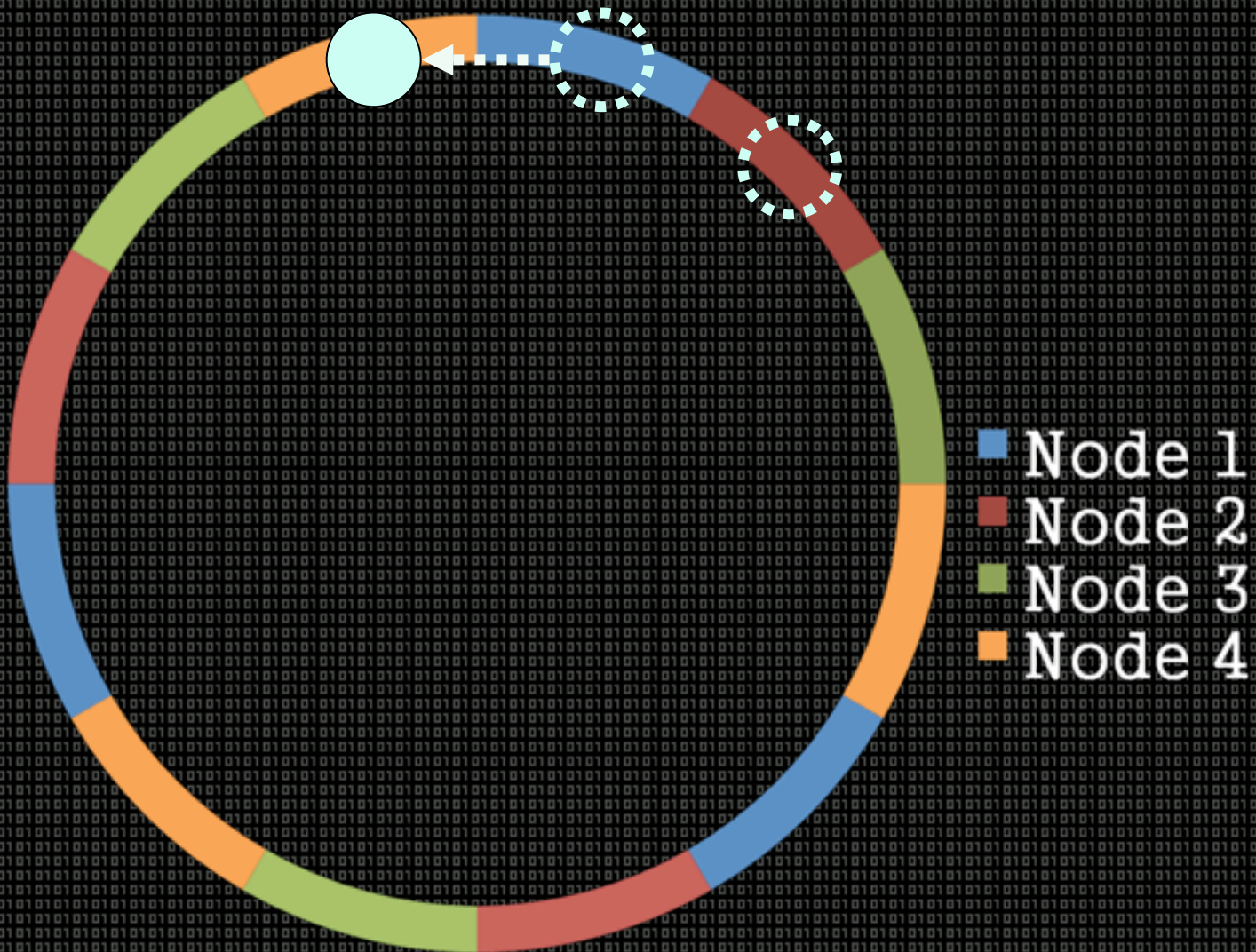recovers



Node 1
Node 2
Node 3
Node 4

Key = "foo", # = N –> handoff hint = true
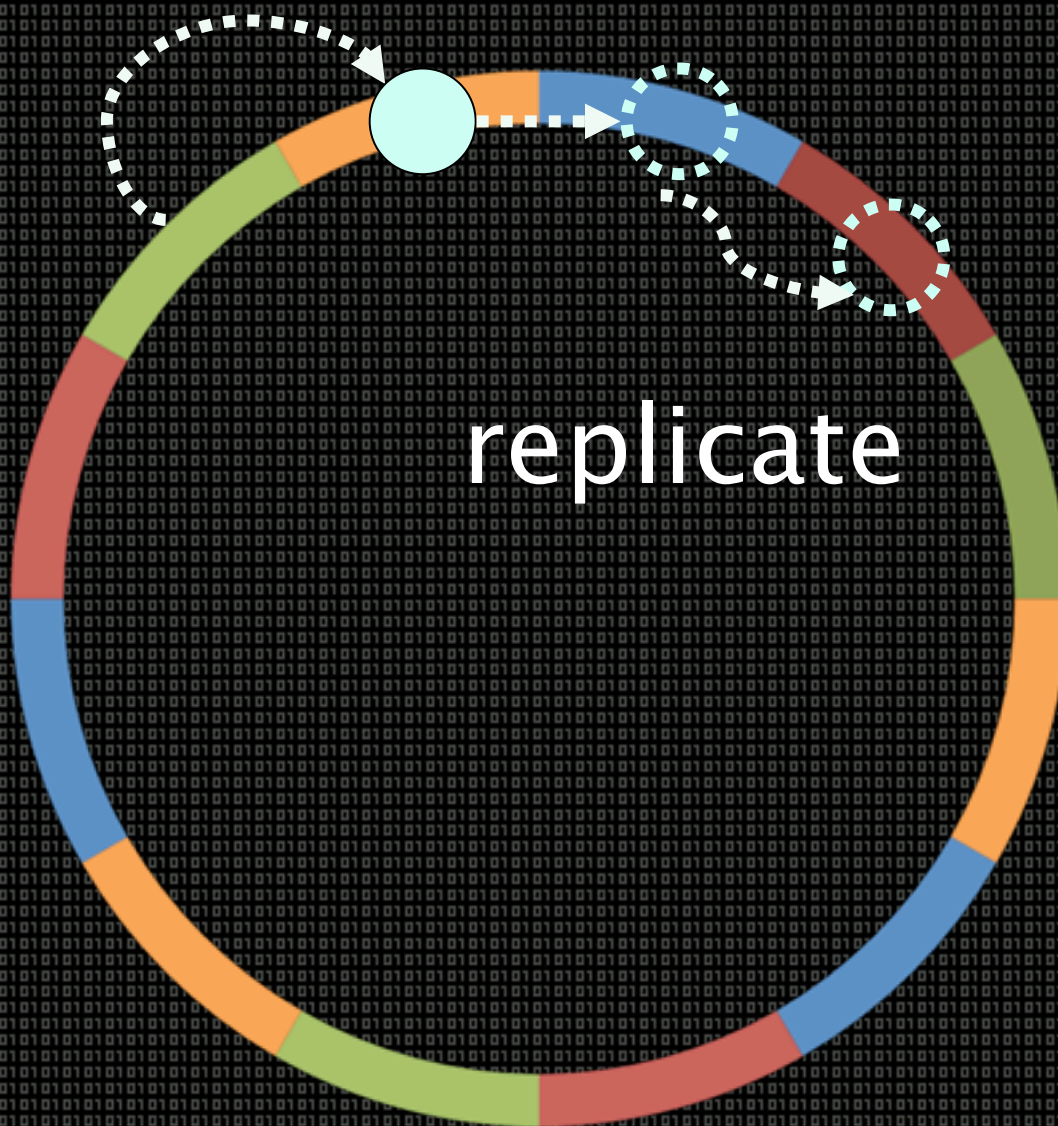
All replicas fail

N

Node 1
Node 2
Node 3
Node 4

All
replicas
recover

handoff

replicate

Node 1
Node 2
Node 3
Node 4

# Latency is an adjustment screw

# Availability is an adjustment screw

CA – irrelevant
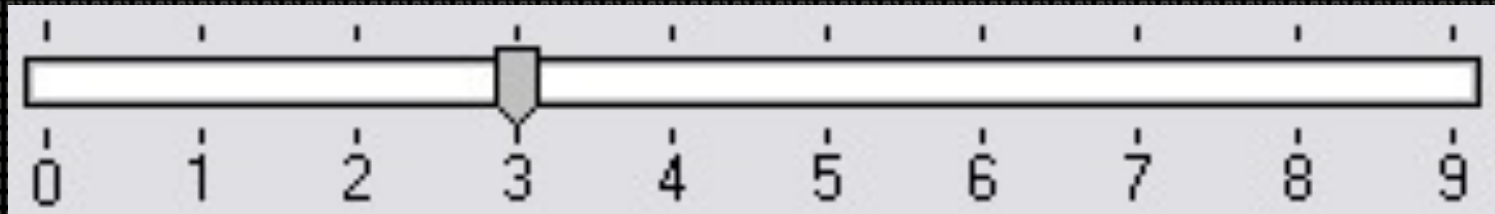
CP – eventually unavailable offering maximum consistency

AP – eventually inconsistent offering maximum availability

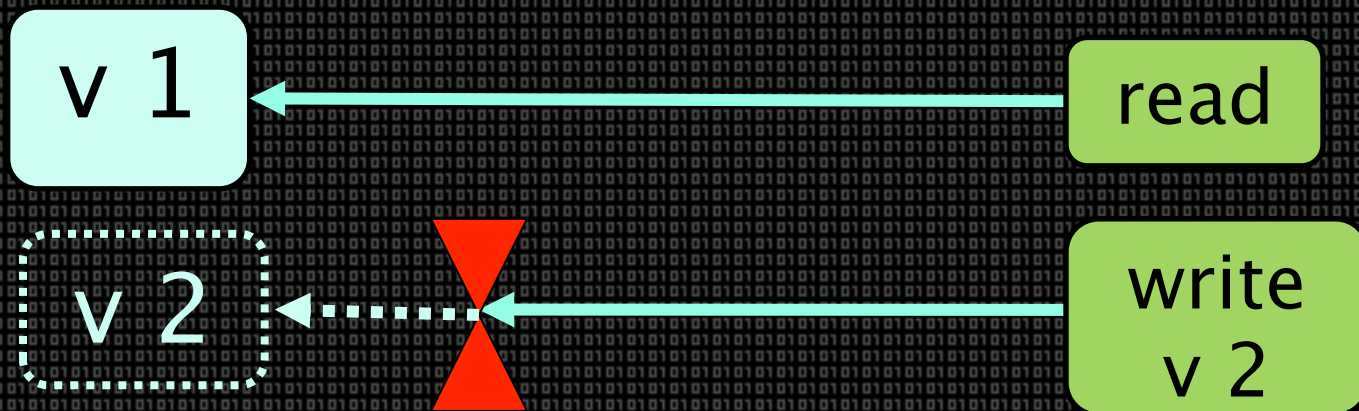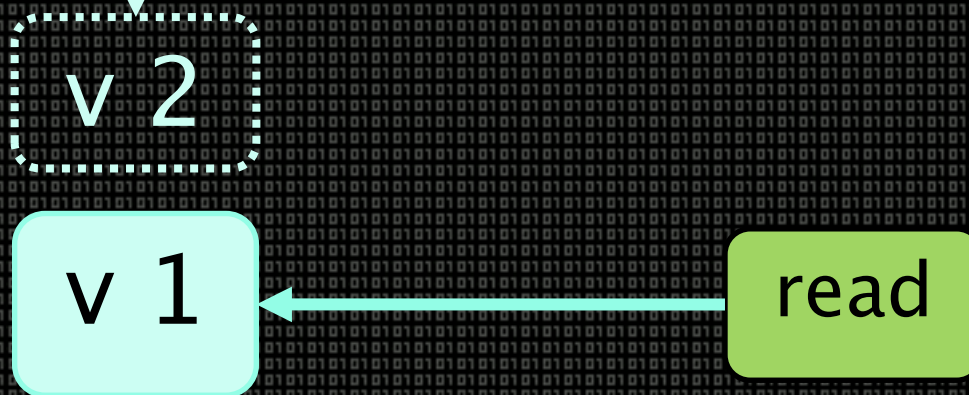A                                                                                    C

Replica 1

CP

v 1 ← read

v 2 ← write v 2

v 2

v 1 ← read

Replica 2

Replica 1

CP (partition)

v 1 ← read

v 2 ← write v 2

Replica 2

v 1 ← read

Friday, August 31, 12
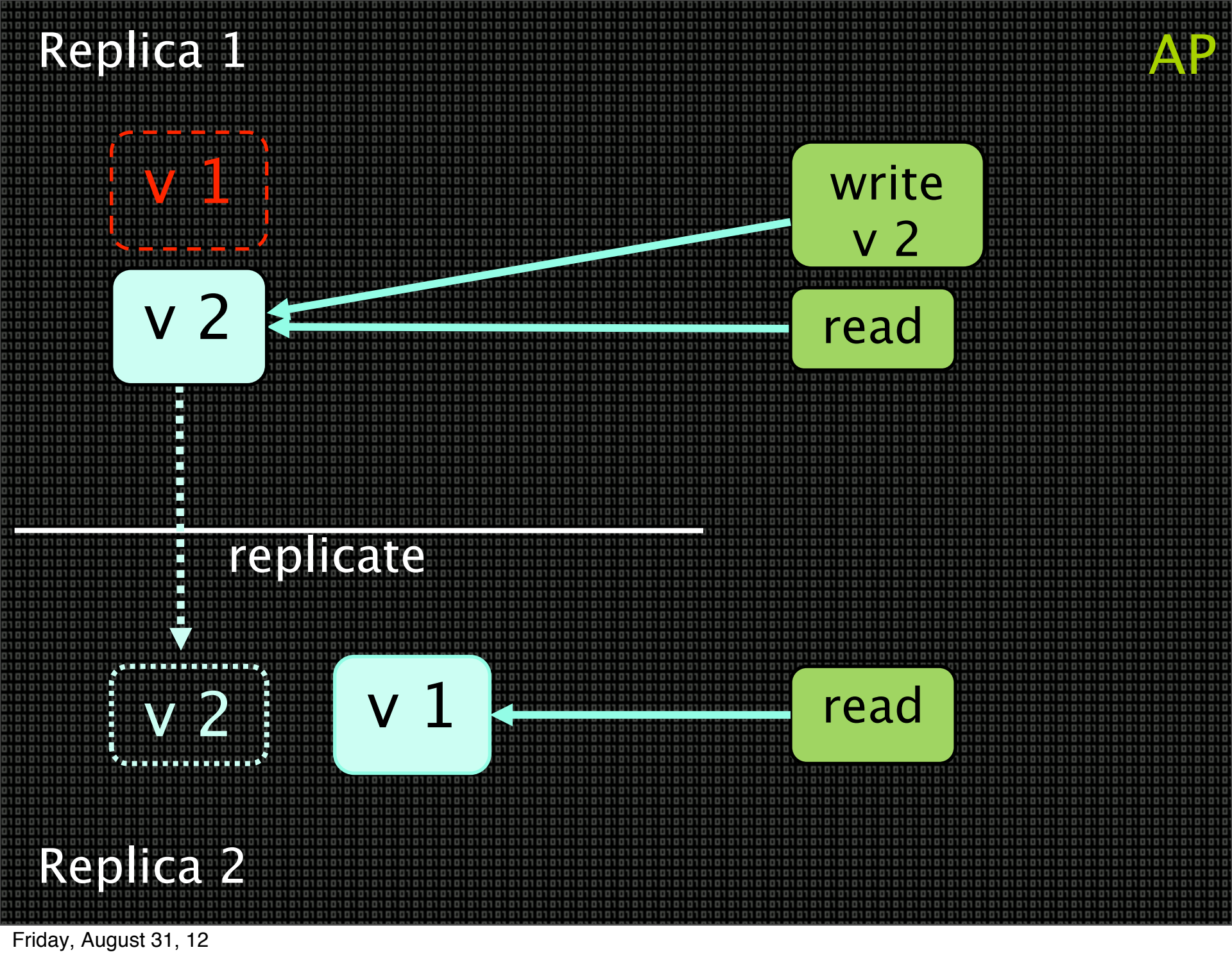
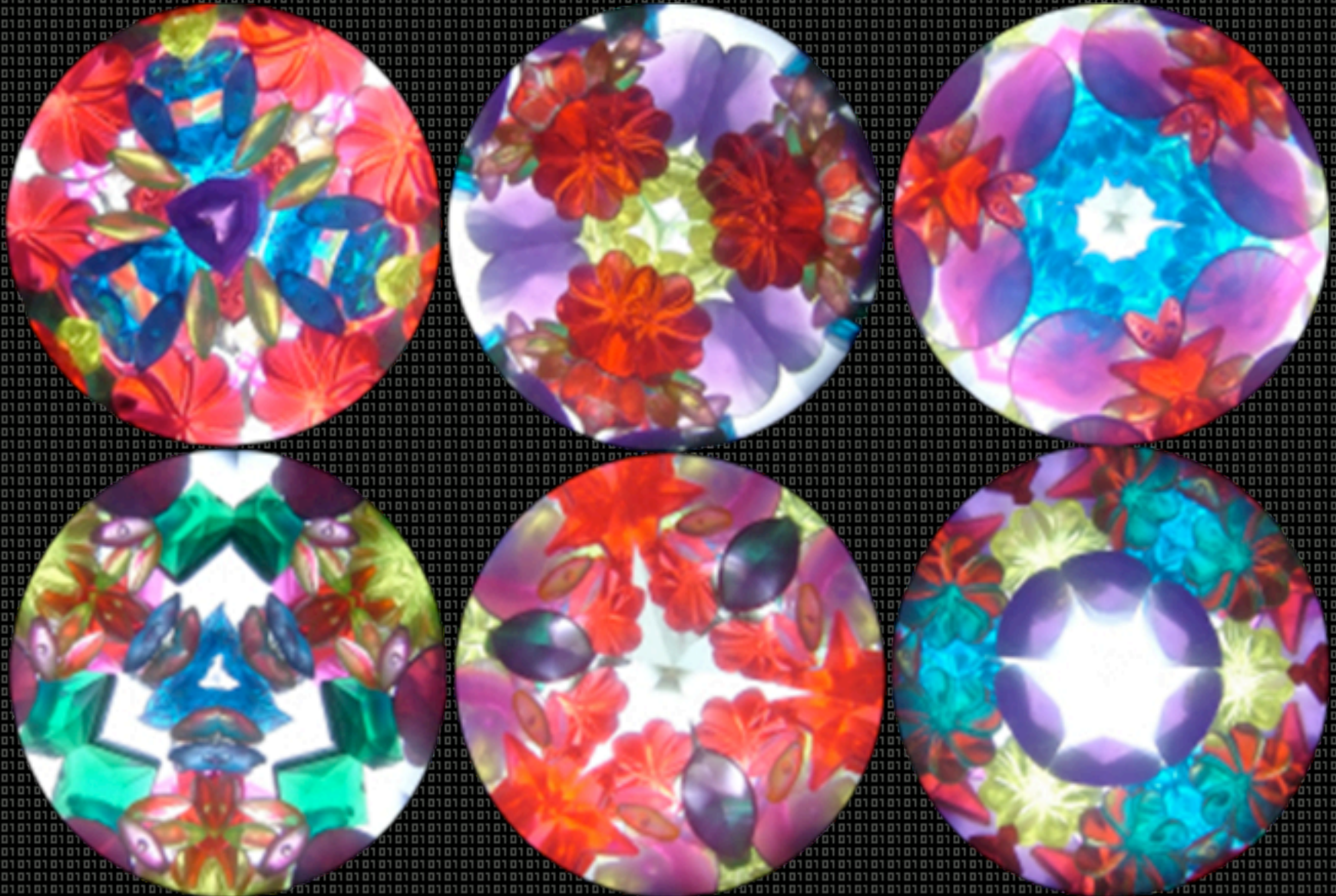Replica 1　　　　　　　AP (partition)

v 1

write
v 2

v 2

read

hint
handoff

v 2

v 1

read

Replica 2

# Frequent structure changes

Thank you

Many graphics I've created myself

Some images originate from istockphoto.com

except few ones taken from Wikipedia and product pages