

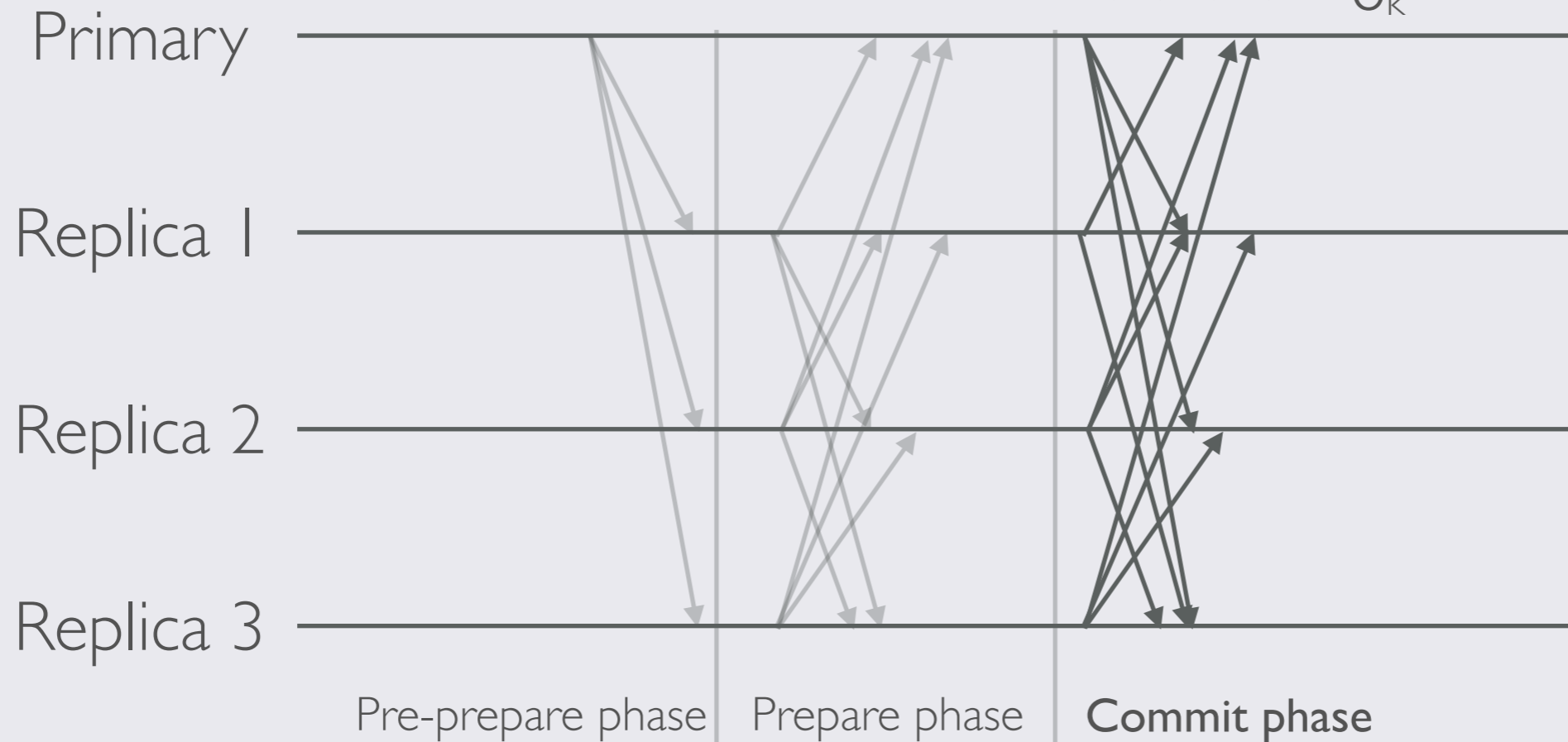
EECS 591

DISTRIBUTED SYSTEMS

Manos Kapritsos
Fall 2021

COMMIT

After collecting a P-Certificate, replica k sends $\langle \text{COMMIT}, v, n, d, k \rangle_{\sigma_k}$ to all replicas

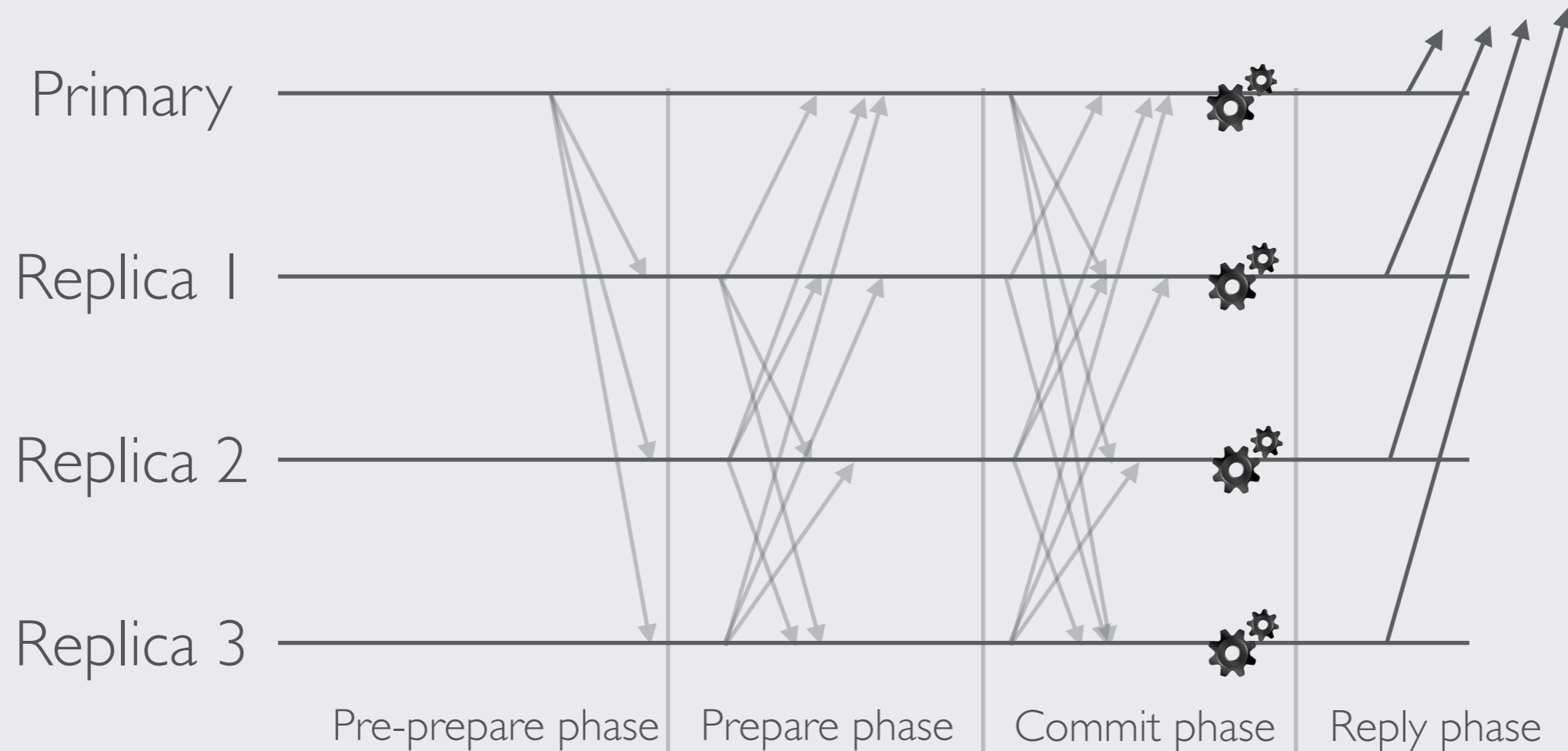


COMMIT CERTIFICATE

- C-Certificates ensure consistent order of requests **across** views
 - **Cannot miss** a P-Certificate during view change
- A replica has a C-Certificate(m, v, n) iff:
 - it had a P-Certificate(m, v, n)
 - its log contains $2f + 1$ matching COMMIT messages from distinct replicas (including itself)
- A replica executes a request when:
 - it gets a C-Certificate for it
 - it has executed all requests with smaller sequence numbers

REPLY

After executing a request, replica k replies to the client with $\langle \text{REPLY}, v, t, c, k, r \rangle_{\sigma_k}$



How many matching requests must the client wait for?

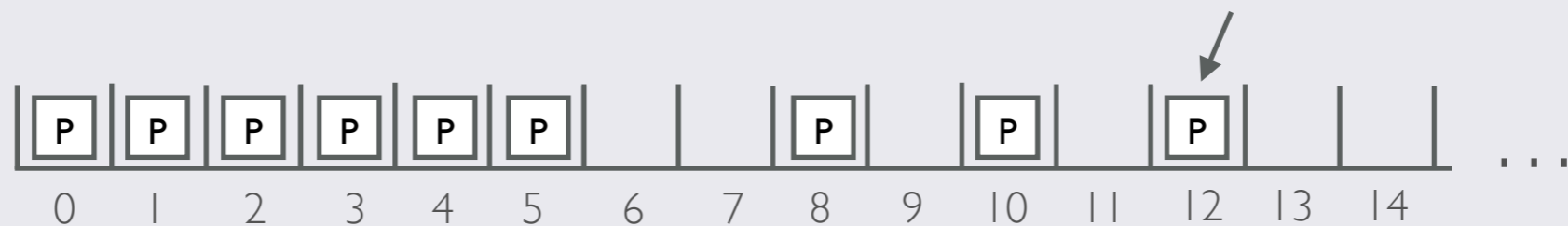
TO ARMS, REPLICAS!!

- A disgruntled replica mutinies:
 - Stops accepting messages (except for VIEW-CHANGE and NEW-VIEW messages)
 - sends $\langle \text{VIEW-CHANGE}, \mathbf{v} + 1, \mathcal{P} \rangle_{\sigma_k}$
 - \mathcal{P} contains all P-Certificates known to replica \mathbf{k}
- A replica joins mutiny after seeing $f + 1$ distinct VIEW-CHANGE messages
- Mutiny succeeds if the new primary collects a new-view certificate \mathcal{V} , indicating support from $2f + 1$ distinct replicas (including itself)

ON TO VIEW $v+1$: THE NEW PRIMARY

- The “primary-elect” p' (replica $v+1 \bmod N$) extracts from the new-view certificate \mathcal{V} :

- the highest sequence number h of any message for which \mathcal{V} contains a P-Certificate \boxed{P}



- two sets \mathcal{O} and \mathcal{N} :

- ▶ if there is a P-certificate for n, m in \mathcal{V} , where $n \leq h$
add $\langle \text{PRE-PREPARE}, v+1, n, m \rangle_{\sigma_p}$ to \mathcal{O}

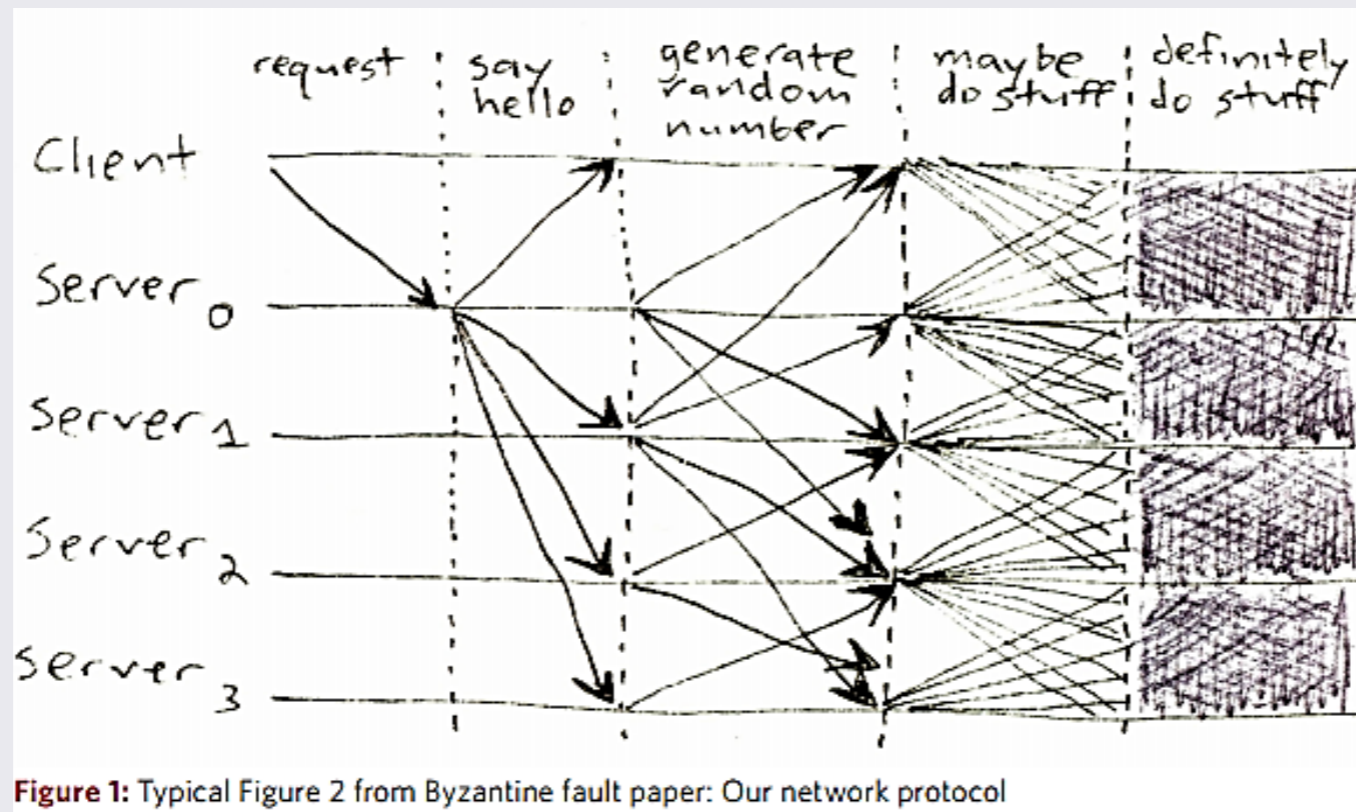
- ▶ otherwise, if $n \leq h$ but there is no P-Certificate
add $\langle \text{PRE-PREPARE}, v+1, n, \text{null} \rangle_{\sigma_p}$ to \mathcal{N}

- p' sends $\langle \text{NEW-VIEW}, v+1, \mathcal{V}, \mathcal{O}, \mathcal{N} \rangle_{\sigma_p}$ to all replicas

ON TO VIEW $v+1$: THE REPLICA

- A replica accepts a NEW-VIEW message for $v+1$ if
 - it is signed properly
 - it contains in \mathcal{V} valid VIEW-CHANGE messages for $v+1$
 - it can verify locally that \mathcal{O} is correct (repeating the primary's computation)
- Adds all entries in \mathcal{O} to its log (as did p')
- Sends a PREPARE to all replicas for each message in \mathcal{O}
- Adds all PREPARE messages to its log and enters new view

BFT: A PERSPECTIVE



[Mickens 2013]

On the other hand:
Google has used BFT in its datacenters and
so do many blockchain approaches

ADMINISTRIVIA

Midterm

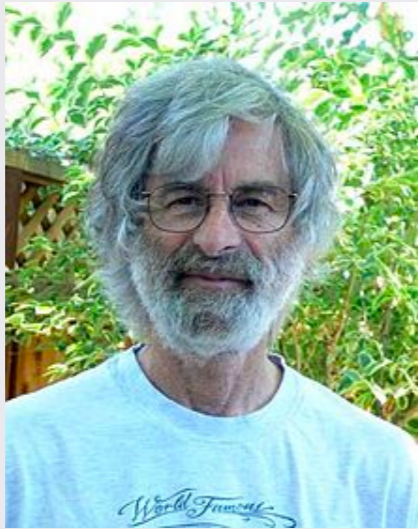
- Next Wednesday 10/27, 12-1:20pm, during class
 - You can use any material listed on the course website

No class on Monday 10/25

- Conflict with SOSF

Research part

- Starts on Monday 11/1
 - You should read both papers and you can review either one



Leslie
Lamport



Barbara
Liskov

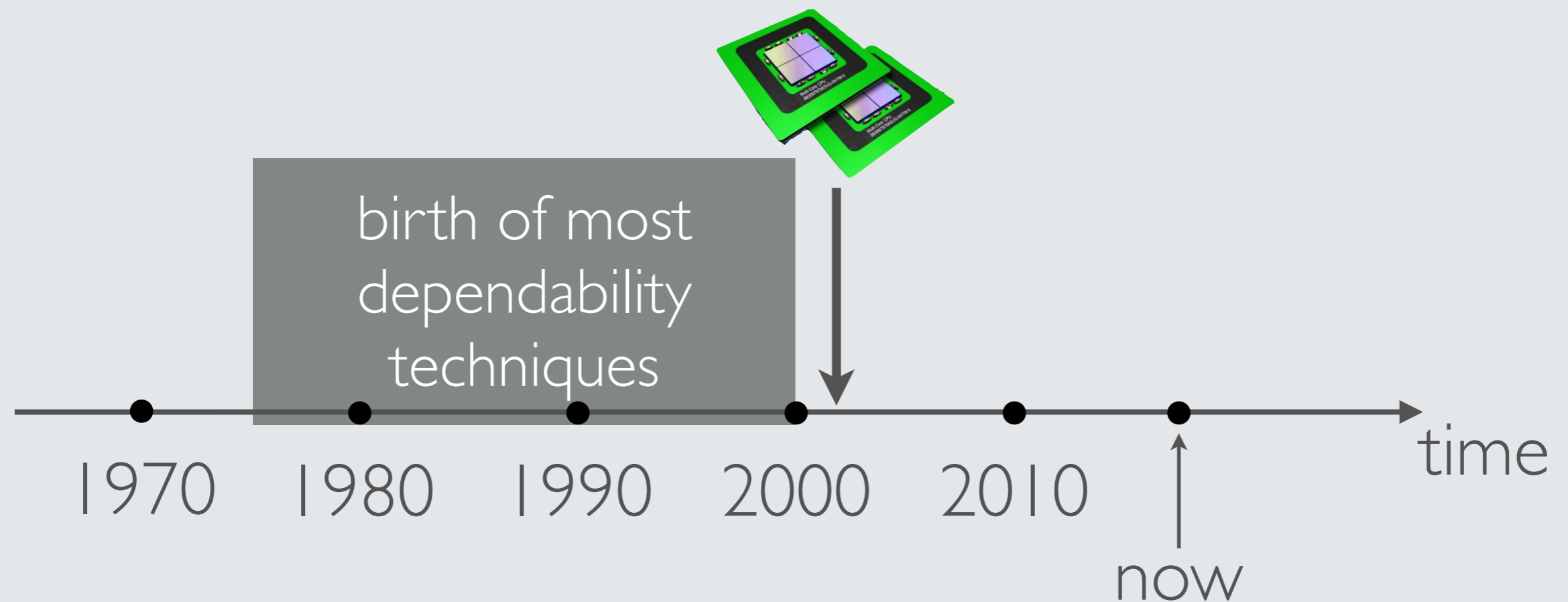


Lorenzo
Alvisi

EVE: REPLICATING MULTITHREADED SERVERS

Kapritsos, Wang, Quema, Clement, Alvisi, Dahlin

THE ACHILLES' HEEL OF REPLICATION



Challenge: scale to **multithreaded execution**

How do we build dependable
~~multithreaded~~ services?

Answer:

State Machine Replication

STATE MACHINE REPLICATION

Ingredients: a server

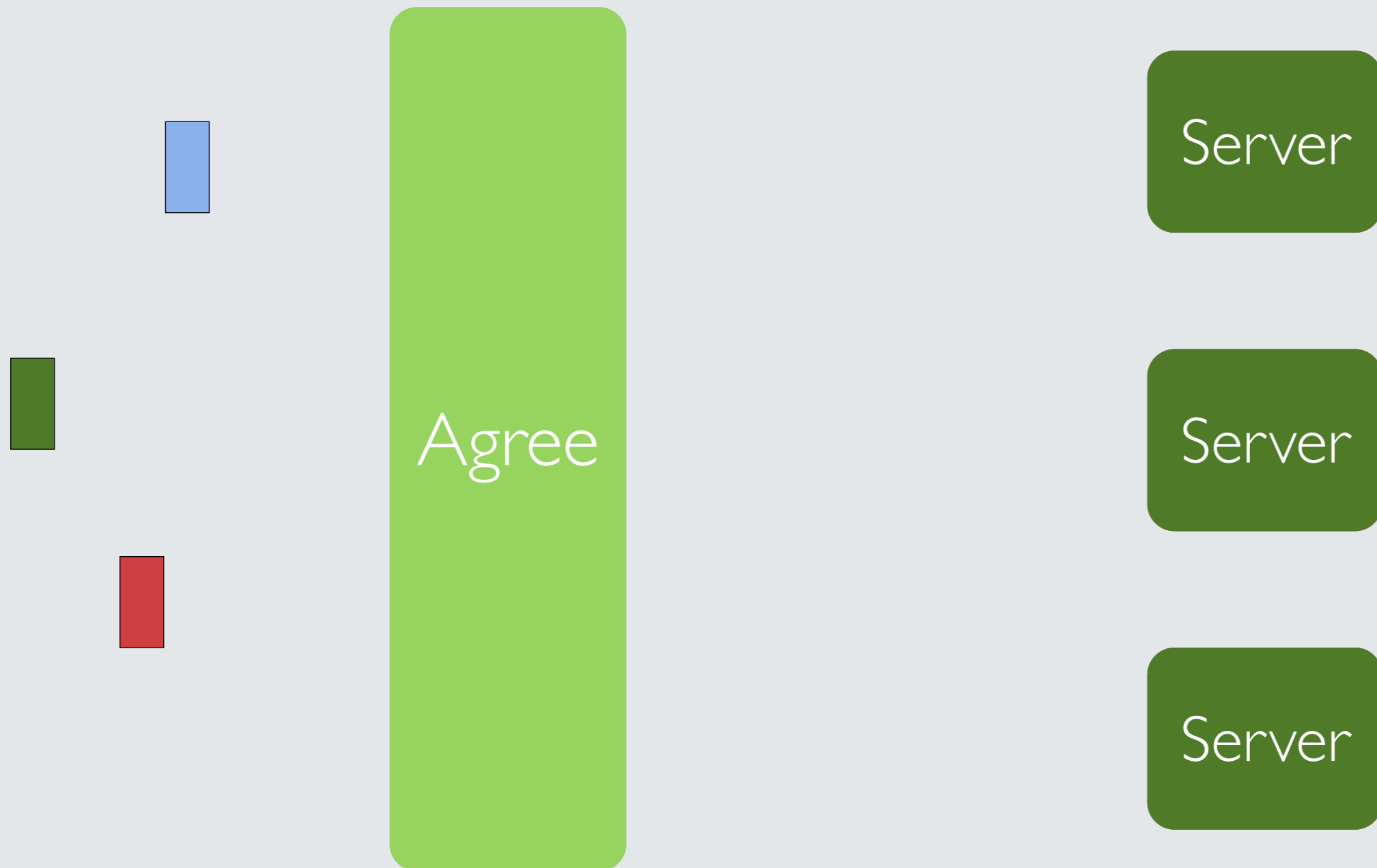
- 1. Make server deterministic (state machine)*
- 2. Replicate server*
- 3. Provide all replicas with the same input*

input

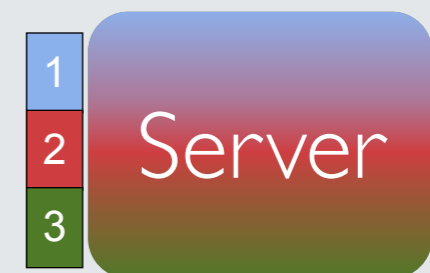
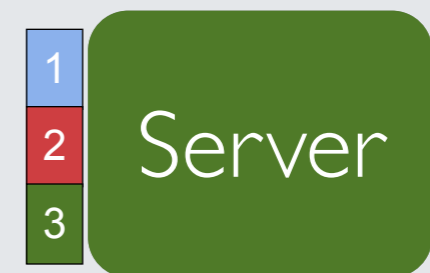
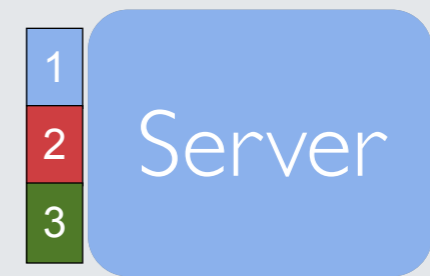
Server

Guarantee: correct replicas will produce the same output

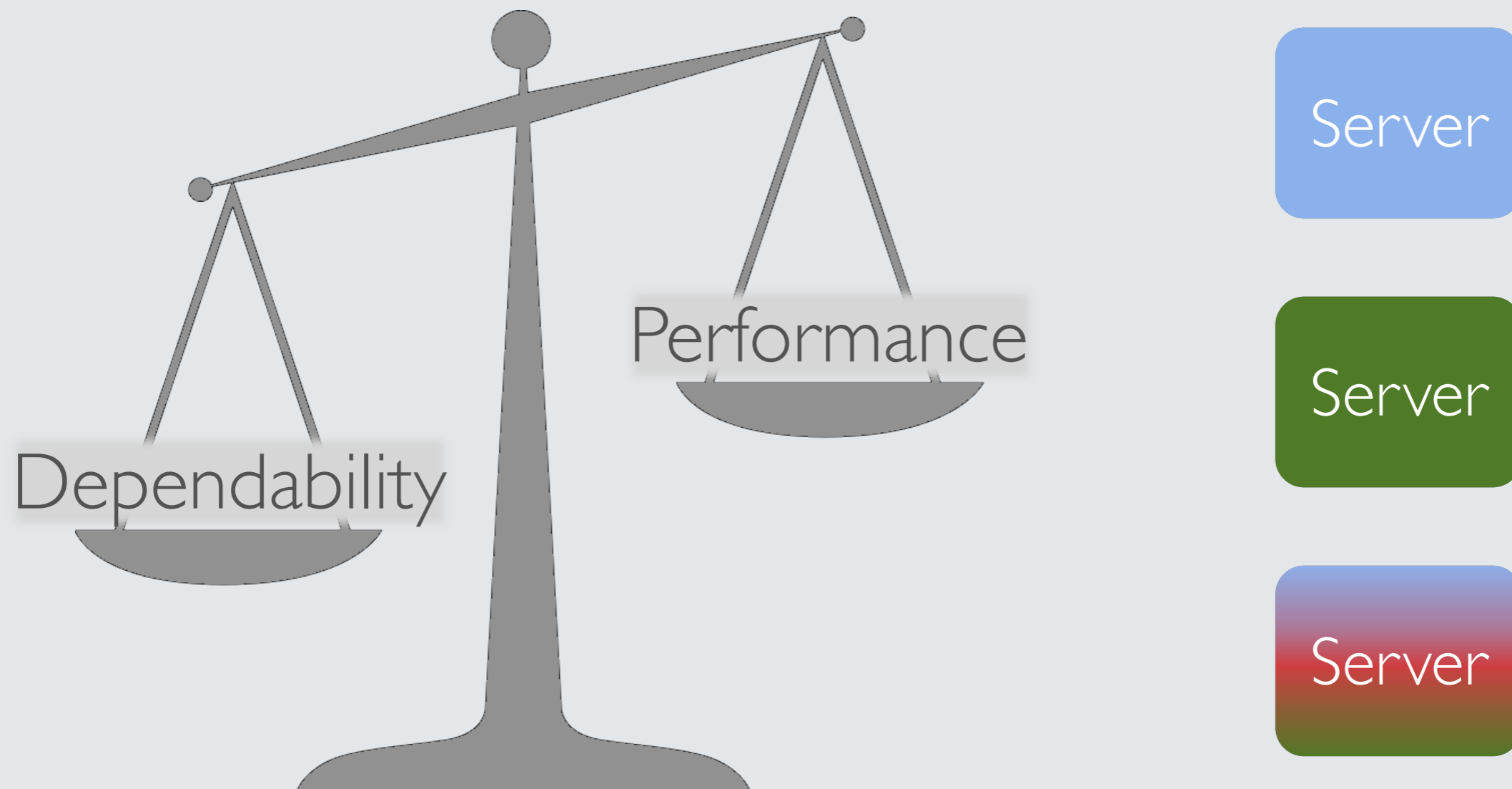
SMR IMPLEMENTATION



How do we build dependable
~~multithreaded~~ services?



How do we build dependable multithreaded services?



Eve (OSDI '12)

Scaling replication to multithreaded execution

SMR requires replica convergence

Agree

Execute

Agree-Execute enforces
sequential execution

EXECUTE-VERIFY



Execute

First execute...

(multithreaded and without
agreeing on the order)

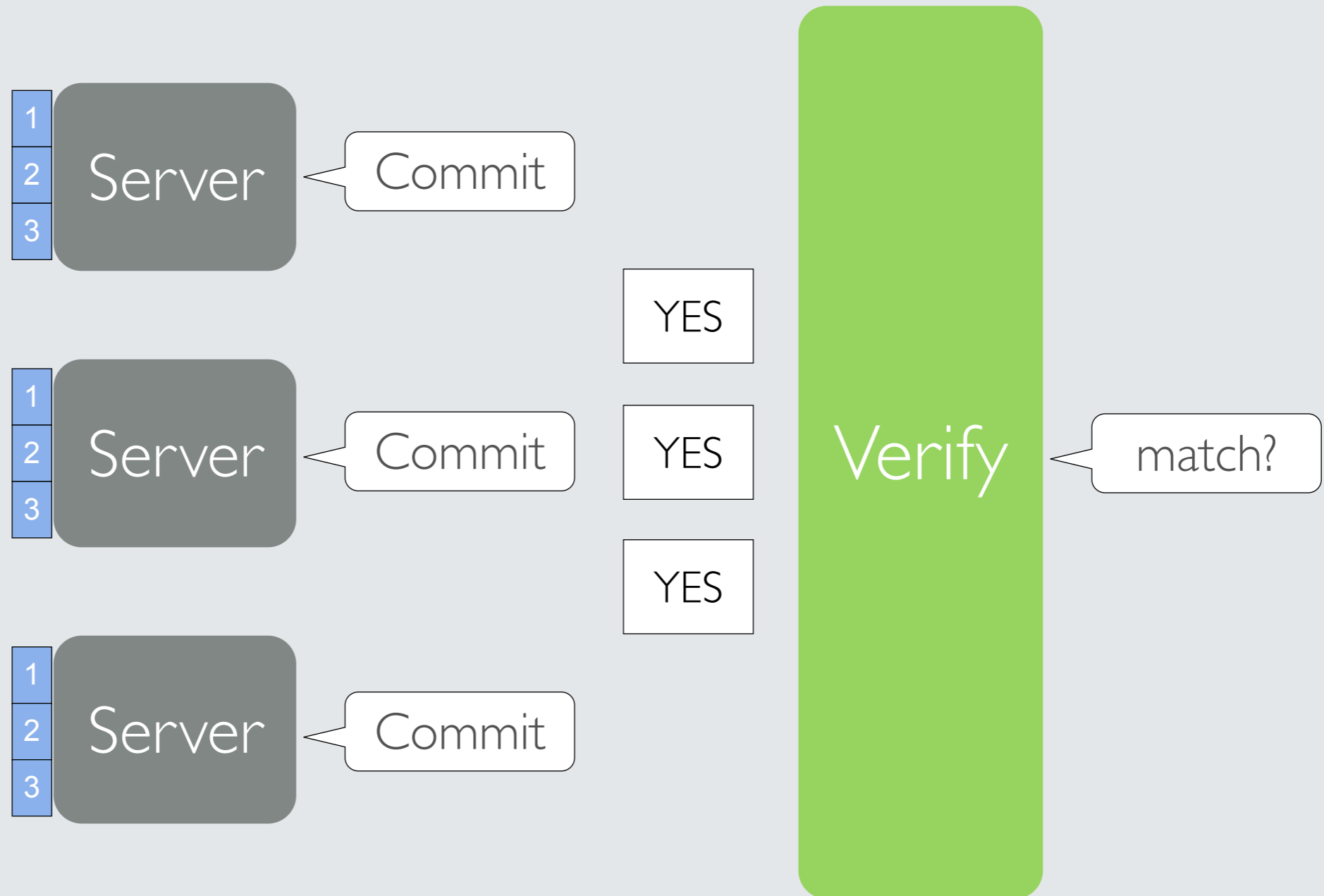


Verify

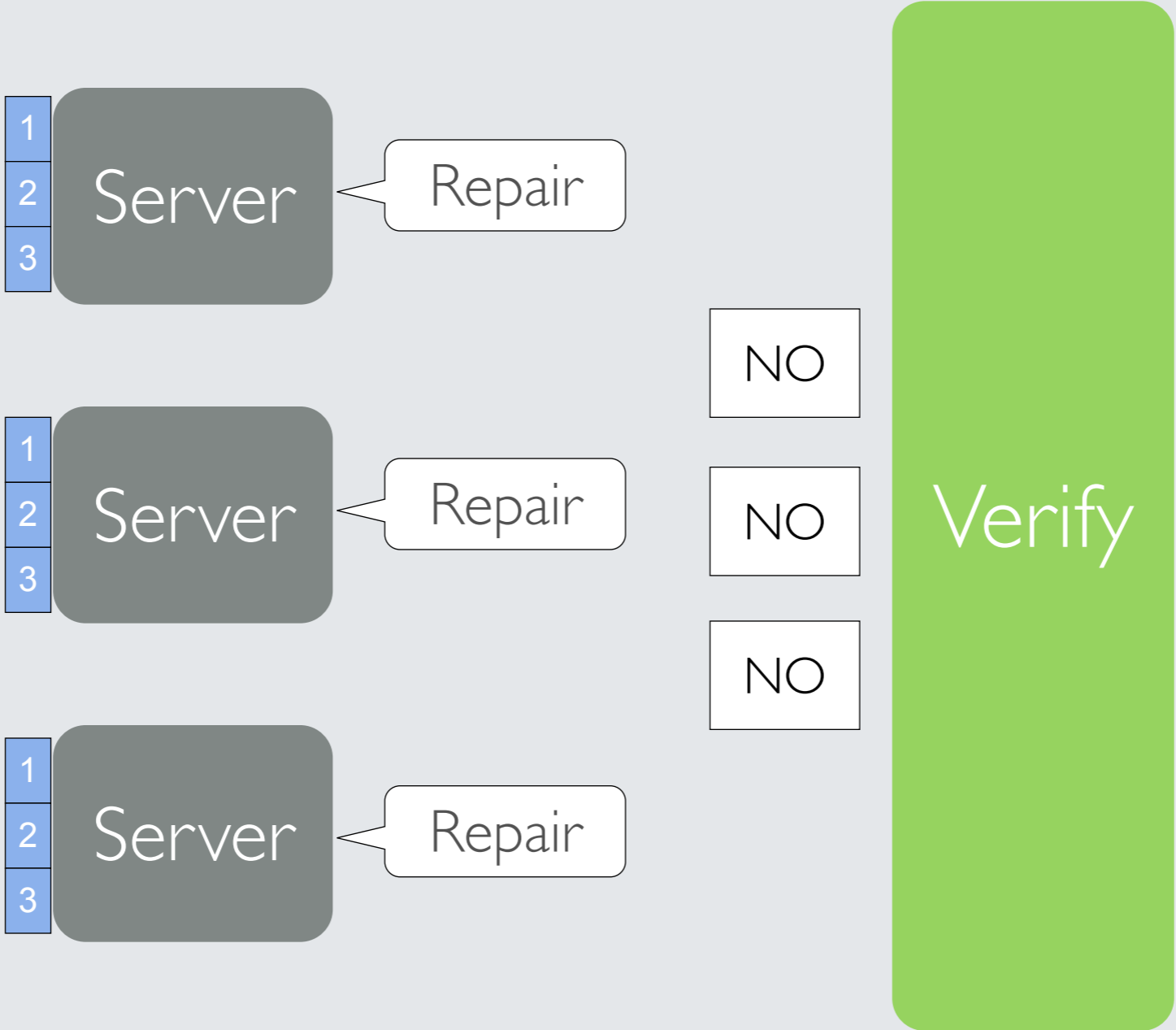
...then verify

(that replicas agree
on the outcome)

ON CONVERGENCE



ON DIVERGENCE



Repair: rollback and re-execute sequentially

Eve's logic at a glance

Frequent

```
if (converged)
  commit
else
  repair divergence
```

Uncommon

1. Make divergence uncommon

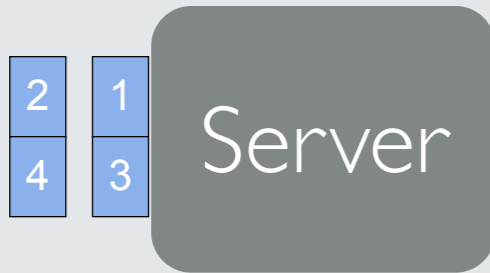
2. Detect divergence efficiently

3. Repair divergence efficiently

MAKING DIVERGENCE UNCOMMON



```
if (converged)
  commit
else
  repair divergence
```



Idea: identify commutative requests

Mixer: group together commutative requests

- Execute requests within a group in parallel



Mixer is a hint, not an oracle

EXAMPLE: TPC-W MIXER

Transaction	Read tables	Write tables
getBestSellers	item, author, order_line	
doCart	item	shopping_cart_line, shopping_cart
doBuyConfirm	customer, address	order_line, item, cc_xacts, shopping_cart_line

3 frequent transactions of the TPC-W browsing workload

EFFICIENT DIVERGENCE DETECTION

Need to compare application states & responses frequently

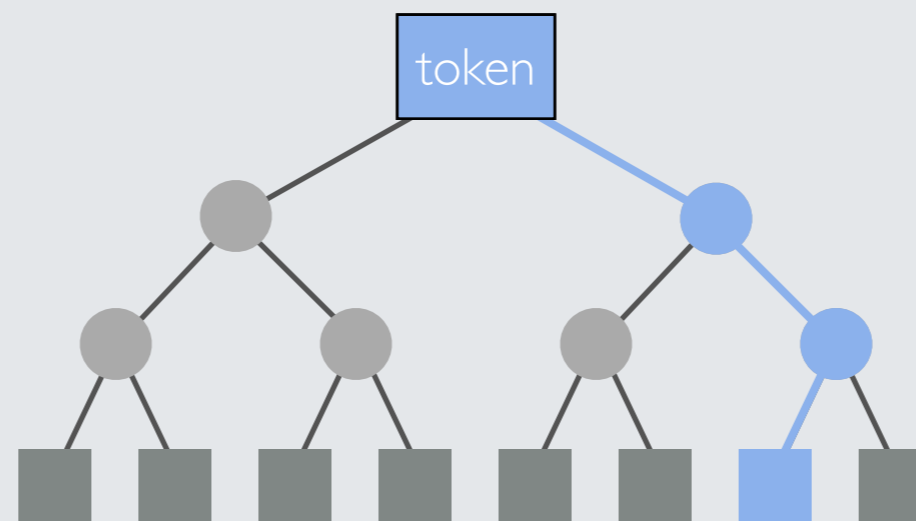
if (converged)

commit

else

repair divergence

Application
state



Merkle tree

EFFICIENT DIVERGENCE REPAIR

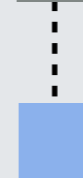
```
if (converged)
  commit
else
  repair divergence
```

Need to rollback application states after every divergence

Application
state



Rollback



Copy-on-write

```
if (converged)
  commit
else
  repair divergence
```

1. Make divergence uncommon

Mixer

2. Detect divergence efficiently

Merkle tree

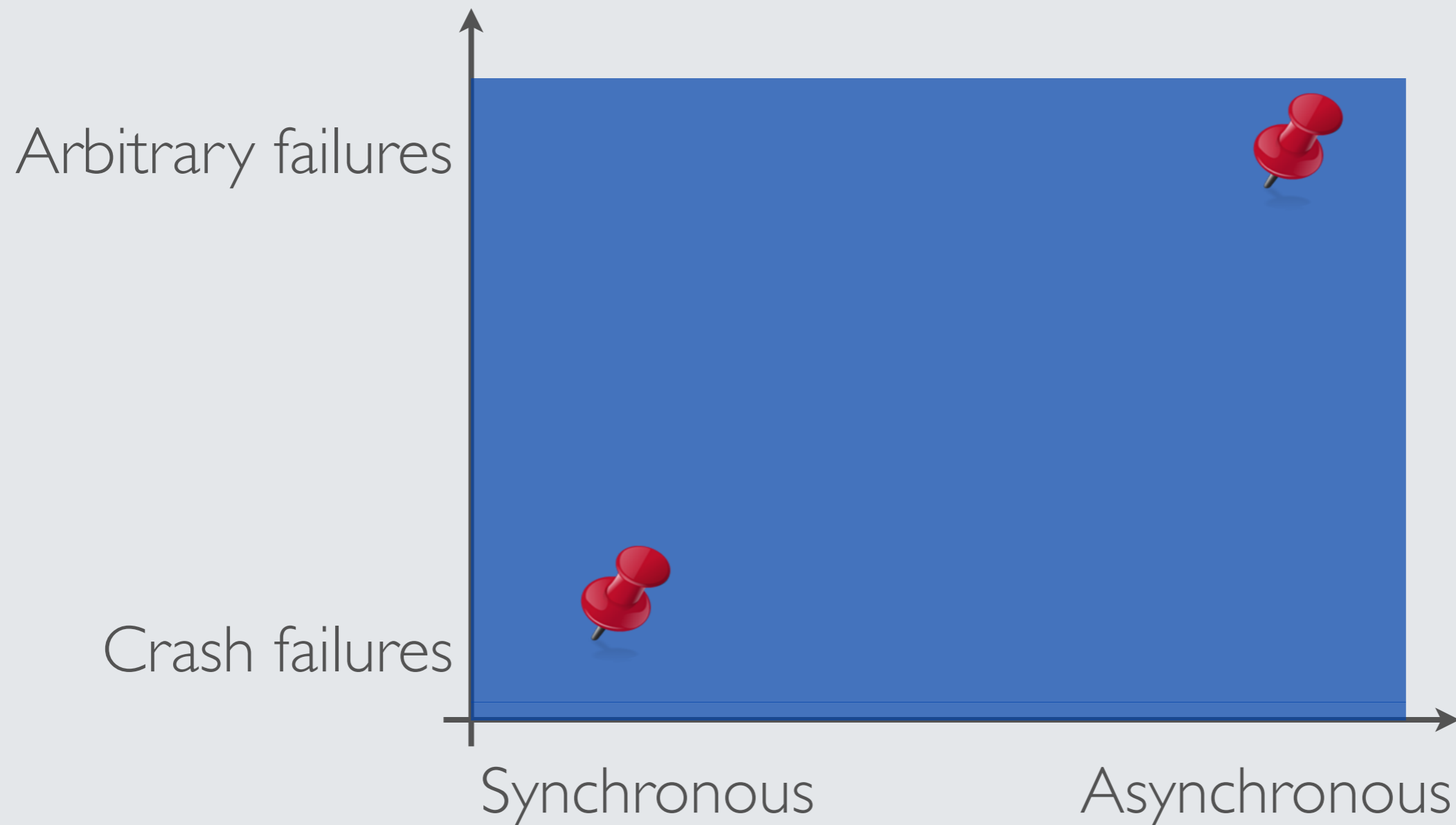
3. Repair divergence efficiently

Copy-on-Write

MASKING CONCURRENCY BUGS



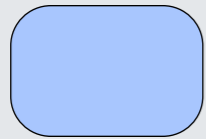
EXECUTE-VERIFY: AN ARCHITECTURAL CHANGE



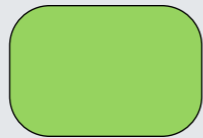
CONFIGURATIONS

Asynchronous BFT

Execution



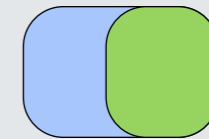
Verification



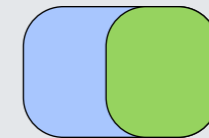
Tolerates 1 arbitrary fault

Synchronous primary-backup

Primary



Backup



Tolerates 1 omission fault

EVALUATION

What is the performance benefit of Eve compared to traditional SMR systems?

Application: H2 Database Engine

Workload: TPC-W (browsing)

