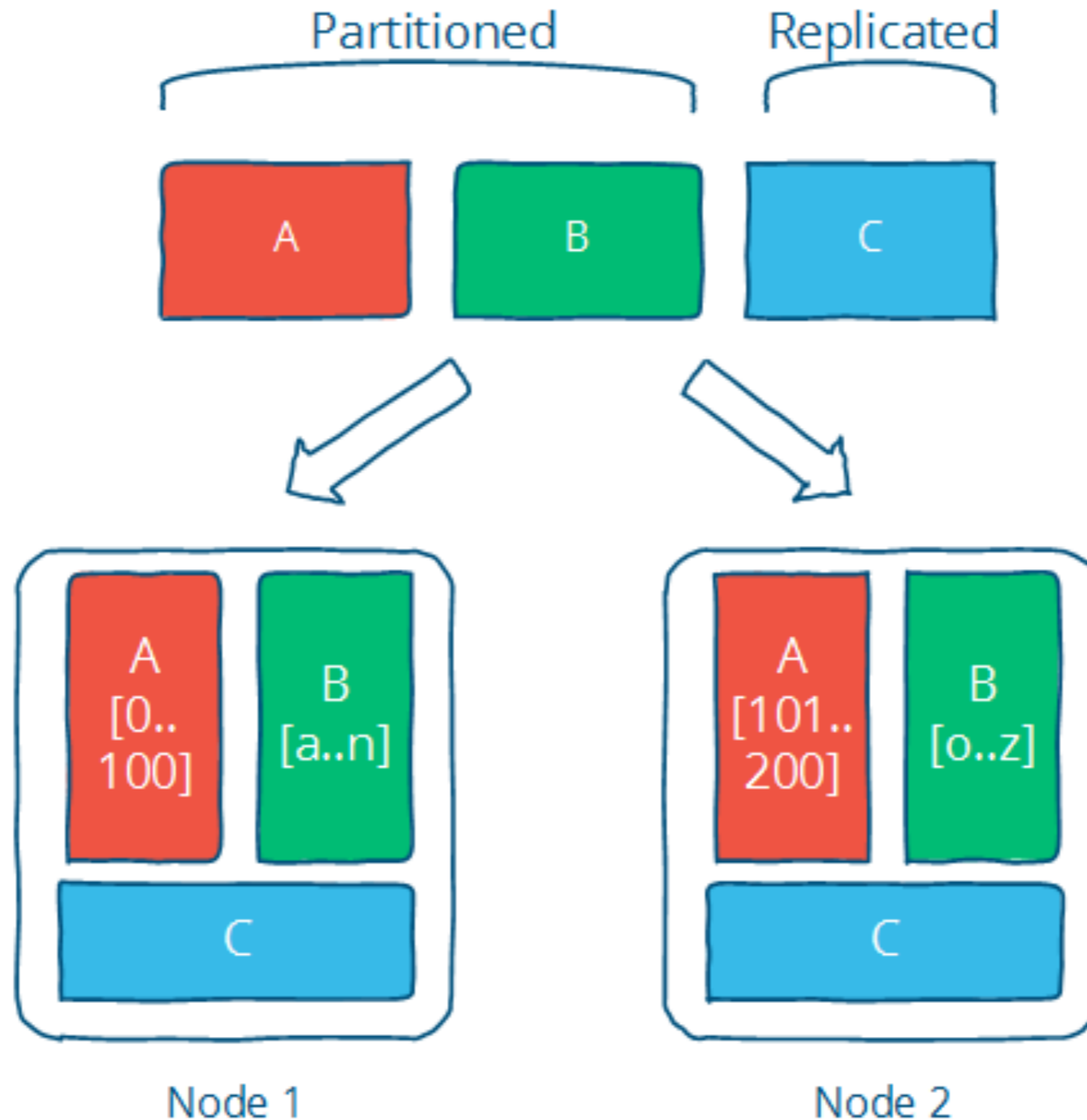# Data Design Techniques

# Partitioning

- <u>Partitioning</u> is dividing the dataset into smaller distinct independent sets

- Partitioning is used to reduce the impact of dataset growth since each partition is a subset of the data.

- Partitioning improves performance by limiting the amount of data to be examined and by locating related data in the same partition.

- Partitioning improves availability by allowing partitions to fail independently, increasing the number of nodes that need to fail before availability is sacrificed.

- Partitioning is also very much application-specific, so it is hard to say much about it without knowing the specifics.
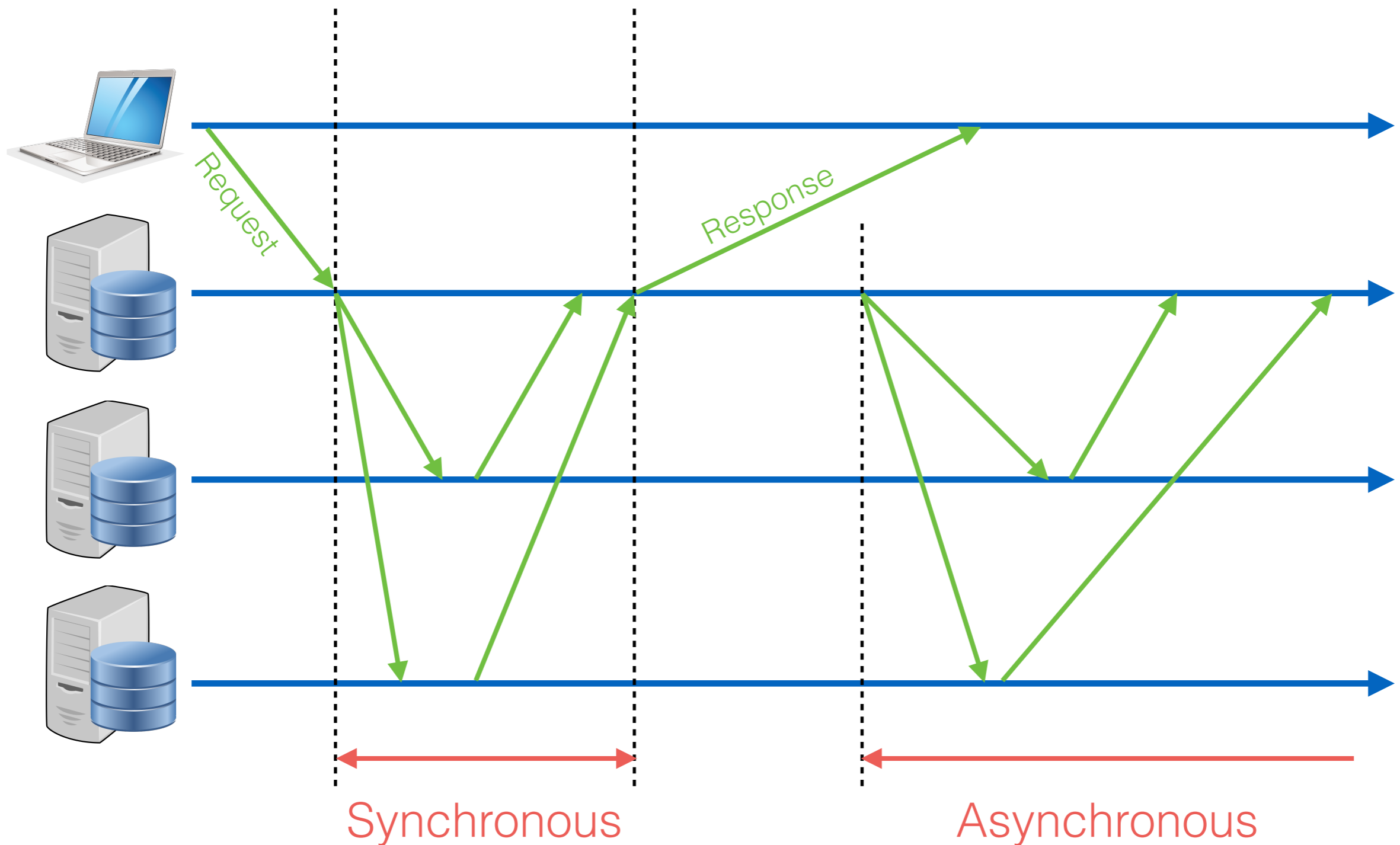
# Replication

- <u>Replication</u> is making copies of the same data on multiple machines; this allows more servers to take part in the computation.

- Replication - copying or reproducing something - is the primary way in which we can fight latency.

- Replication improves performance by making additional computing power and bandwidth applicable to a new copy of the data.

- Replication improves availability by creating additional copies of the data, increasing the number of nodes that need to fail before availability is sacrificed.
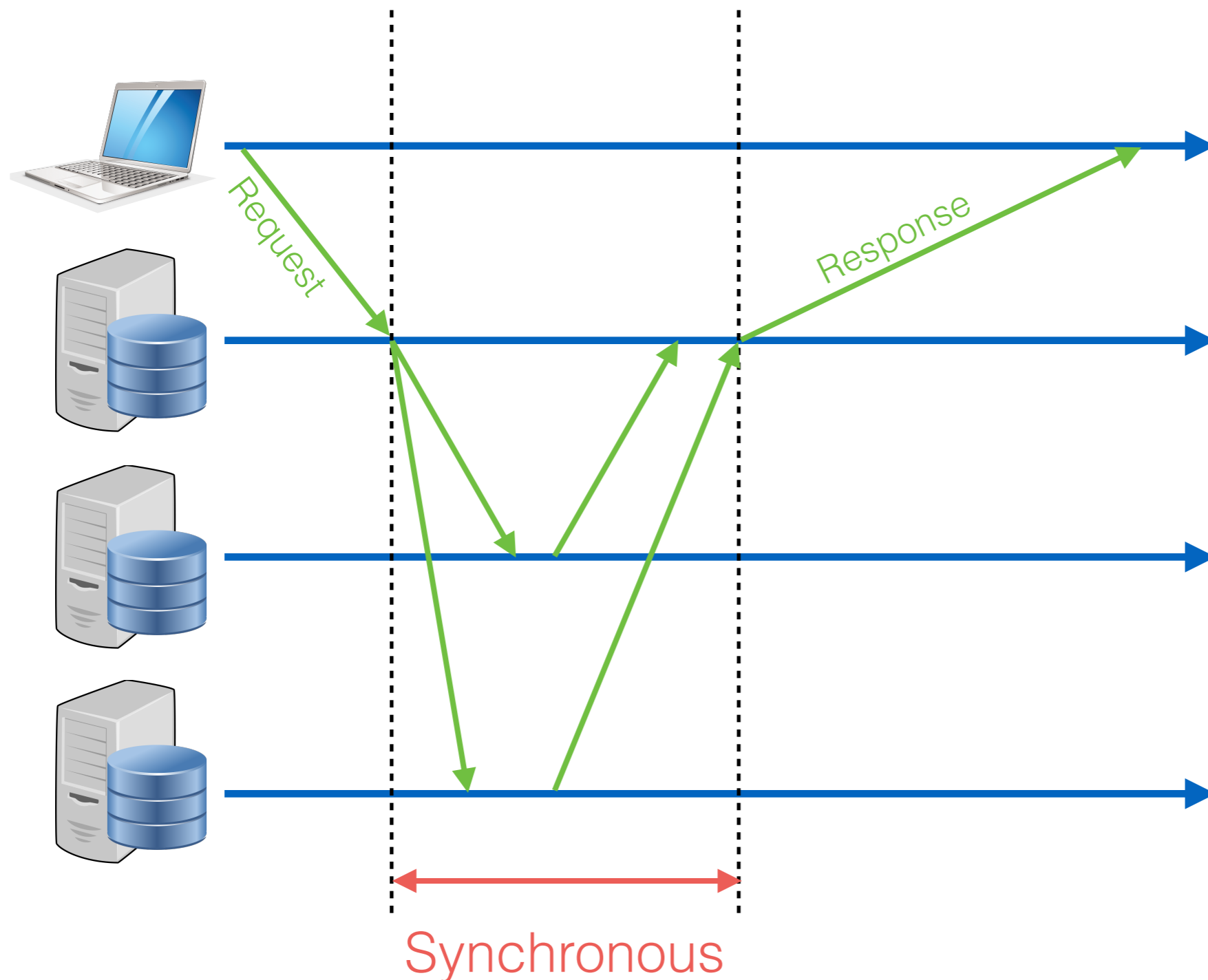
# Replication & Consistency

- one of many problems in distributed systems

- it is often the part that people are most interested in

- provides a context for many subproblems, such as leader election, failure detection, consensus and atomic broadcast

- References:

  - Coulouris, Dollimore, Kindberg, Blair, Distributed Systems - Concepts and Design, ch 15

  - http://book.mixu.net/distsys

  - http://research.microsoft.com/en-us/people/philbe/chapter7.pdf

  - http://research.microsoft.com/en-us/um/people/lamport/pubs/paxos-simple.pdf

  - http://static.googleusercontent.com/media/research.google.com/it//archive/paxos_made_live.pdf

# Communication Patterns



Request

Response

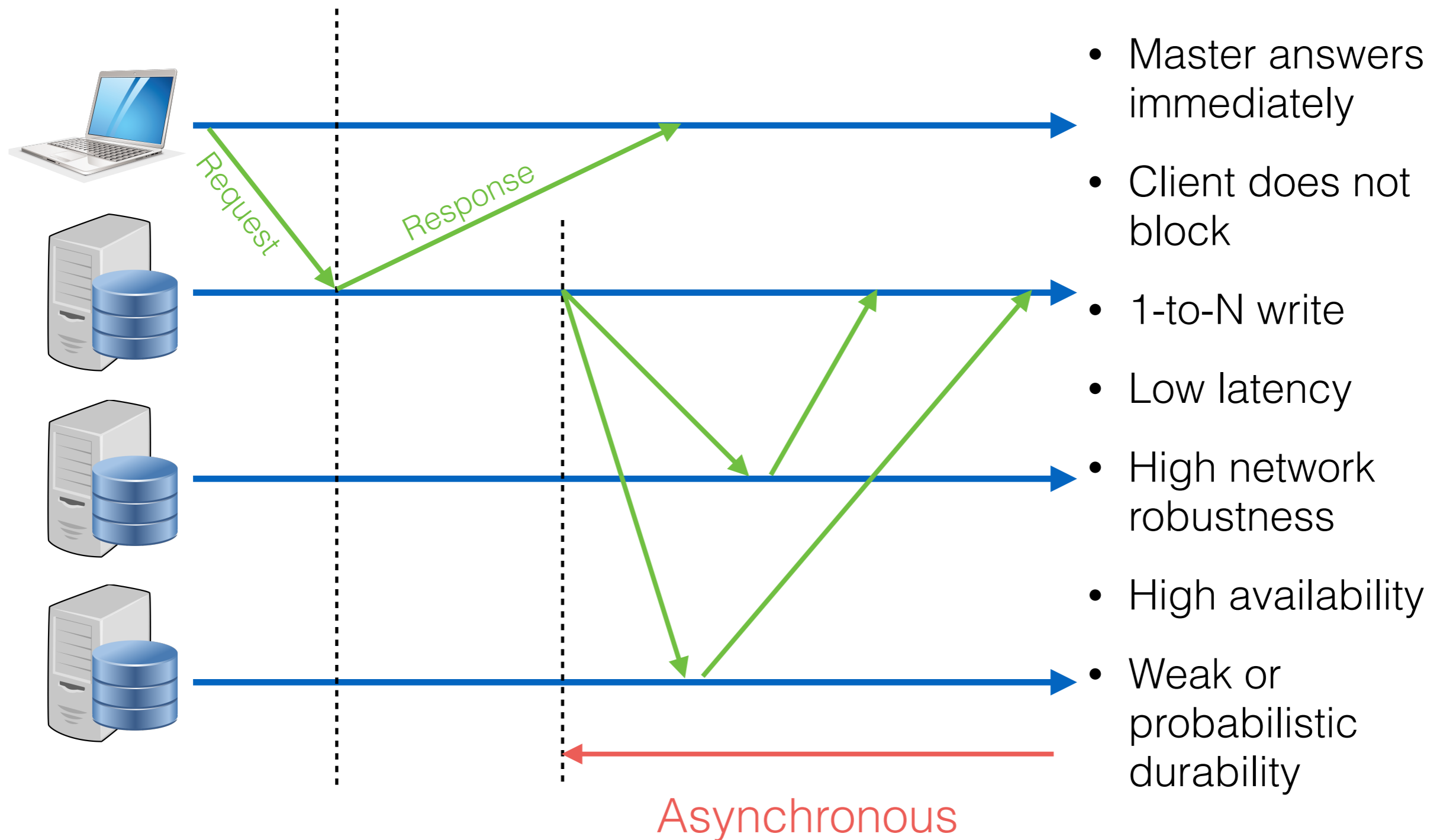Synchronous

Asynchronous

# Synchronous Replication

or active, or eager, or push, or pessimistic



Request

Response

Synchronous

- Three stages
- Client blocks
- N-to-N write
- **Transactions**
- High latency
- Low network robustness
- Low availability
- High durability

# Asynchronous Replication

or passive, or pull, or lazy



- Master answers immediately
- Client does not block
- 1-to-N write
- Low latency
- High network robustness
- High availability
- Weak or probabilistic durability

Request

Response

Asynchronous

# Designing Replication Strategies

|  | Master (primary) | Group (everywhere) |
|---|---|---|
| **Asynch (lazy)** | | |
| **Synch (eager)** | | |

# Pros & Cons

Synchronous

+ No inconsistencies

+ Local copies always updated

+ Atomic changes

- Global single transaction

- Performance

Group

+ Anyone transaction
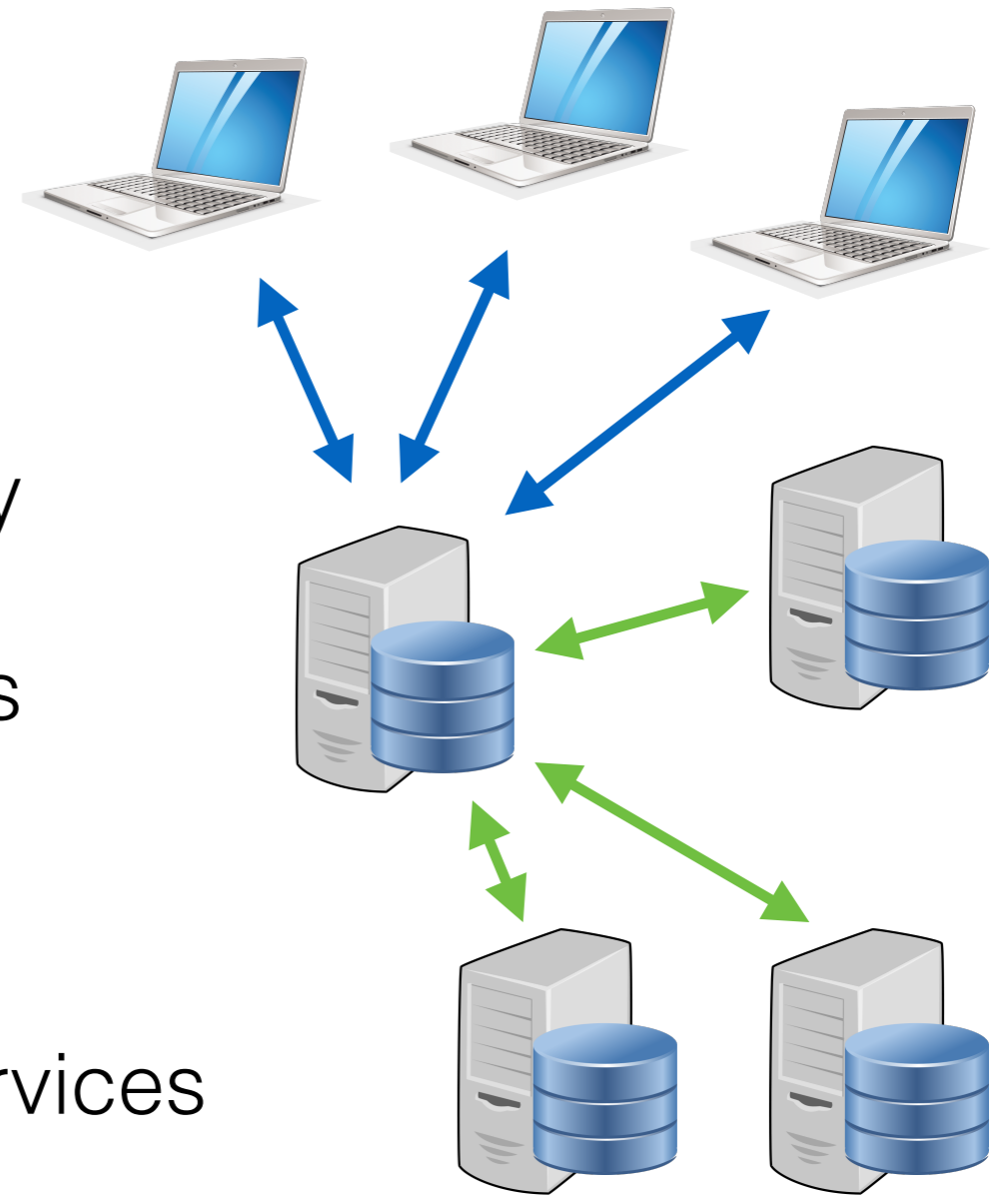
+ Load balancing

- Synchronized copies

Asynchronous

+ Local transaction

+ Performance

- Inconsistencies

- Local copies not updated

- Replication not transparent

Primary

+ Simple

+ Centralized updates

- Load unbalancing

- Single POF
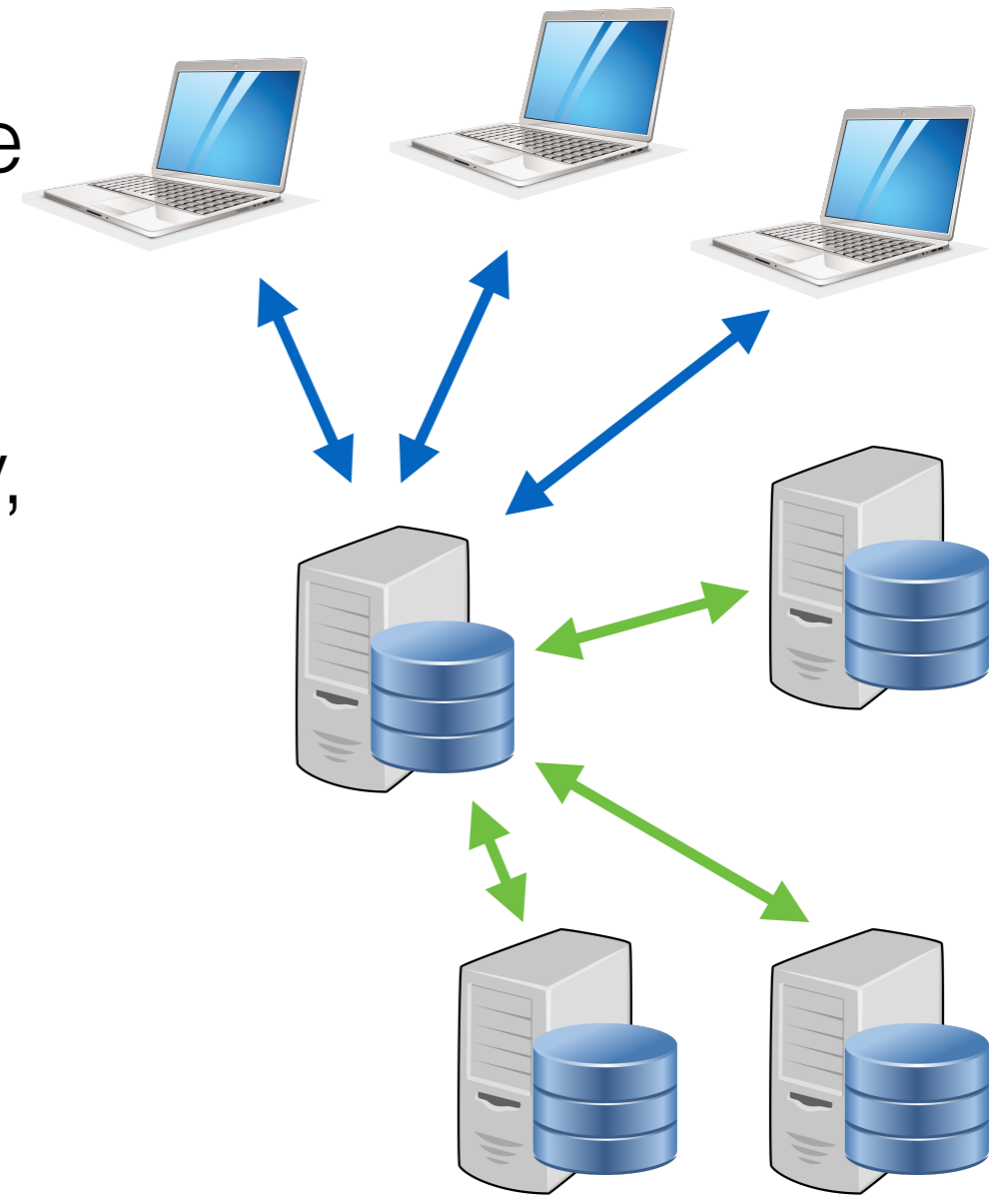
- Local copies not updated

# Asynch + Primary

- Primary Copy Replication

  - Simplest approach:

    - All requests go to primary copy

    - Other copies serve as backups

  - Widely used in practice

    - Keep backups of important services

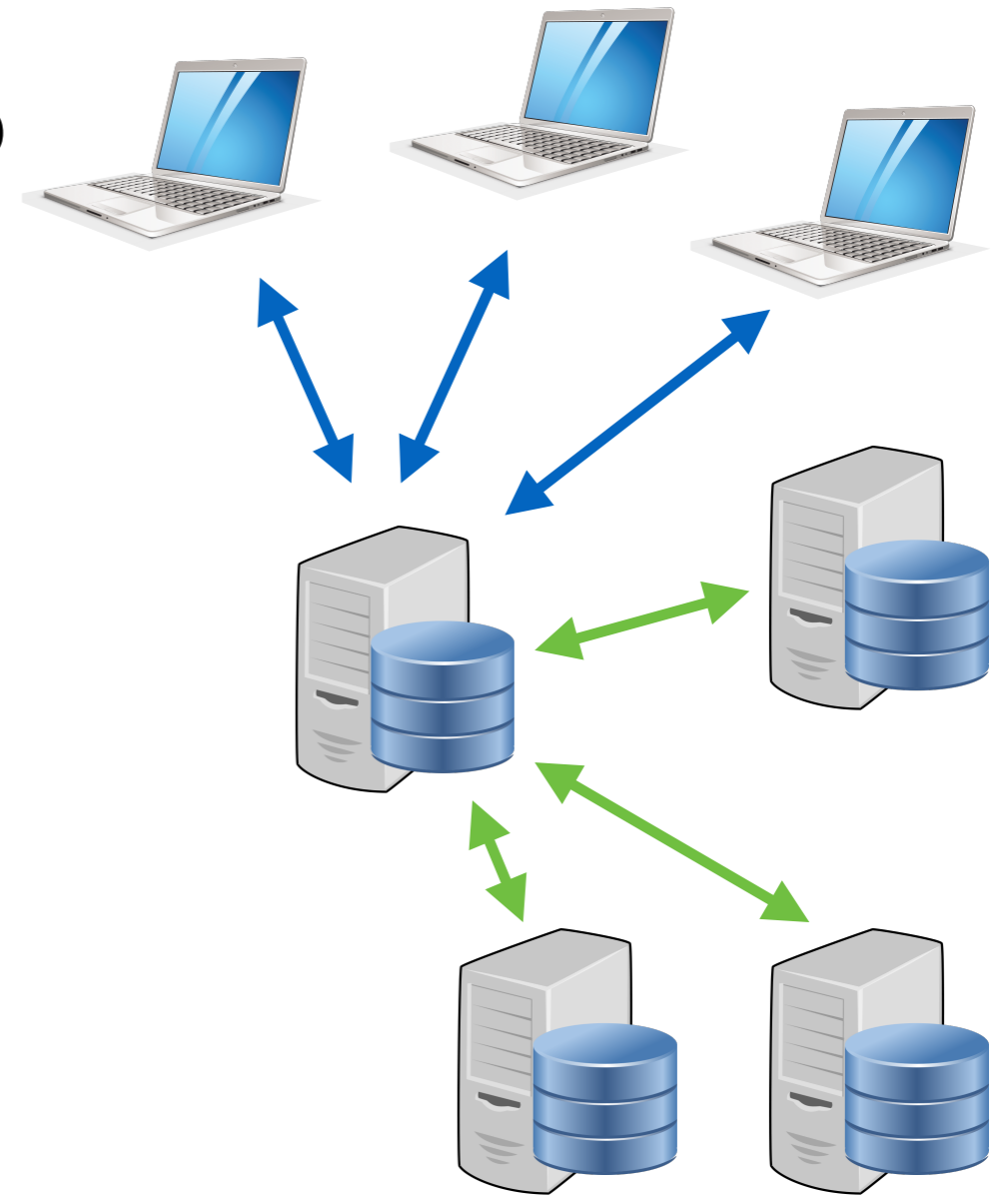    - Backup can go "live" when primary fails

# Primary Copy Replication

- All updates are performed on the primary

- log of operations (or alternatively, changes) is shipped across the network to the backup replicas

- Asynch and synch variants
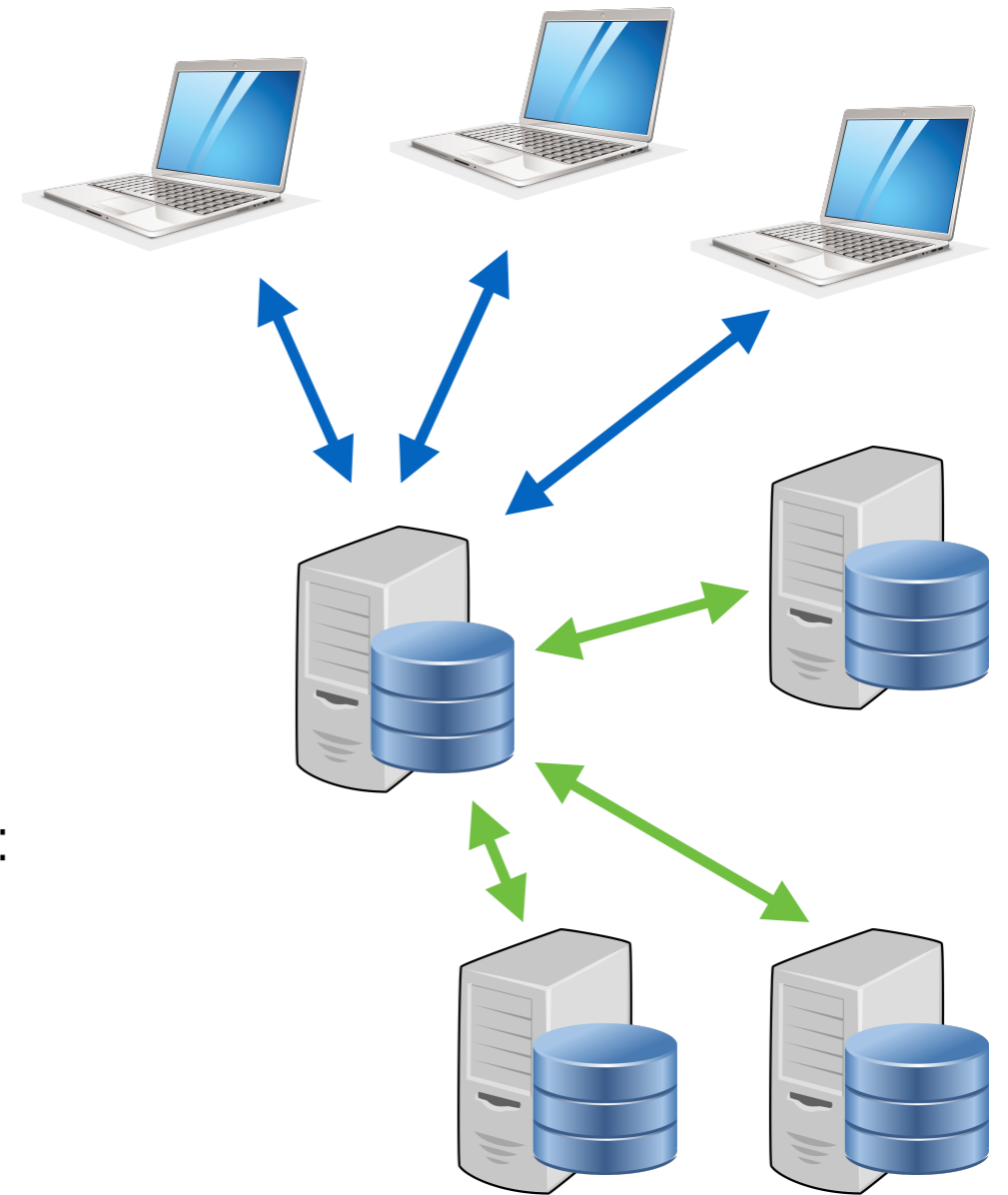
- Asynch used by MySQL and MongoDB

# Primary Copy Fault Tolerance

- Failures can cause violations to consistency guarantees

- Site failures, network partitions, undelivered messages

- We need an <u>atomic commitment protocol</u> (ACP)

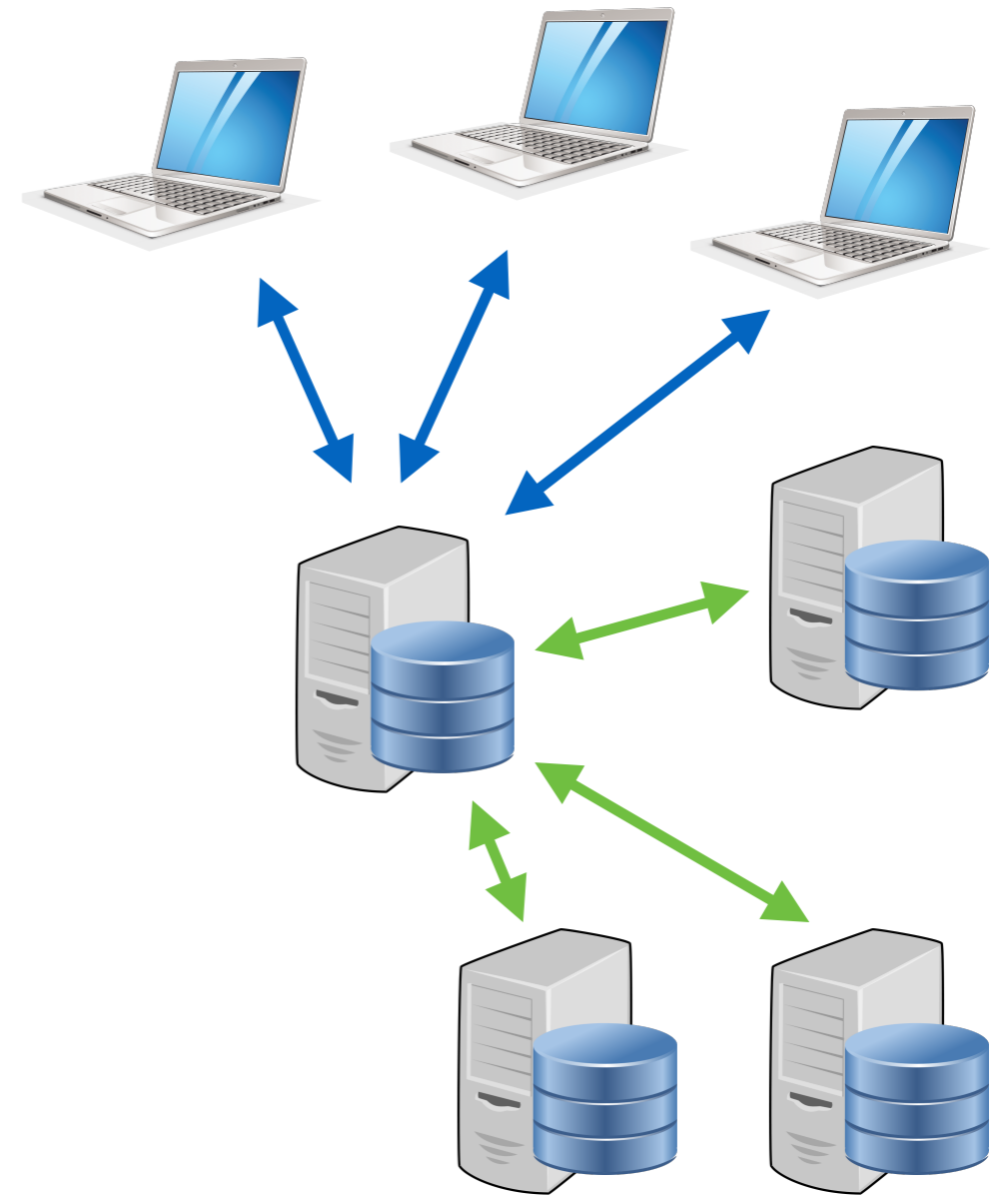- Commonly used in database systems (i.e., transactions)

# Atomic Commitment Protocols

- 1 coordinator, N participants

- Coordinator knows the names of all participants

- Participants know the name of the coordinator, not necessarily each other

- Each site contains a distributed transaction log surviving failures

- Each process may cast exactly one of two votes: **Yes** or **No**

- Each process can reach exactly one of two decisions: **Commit** or **Abort**
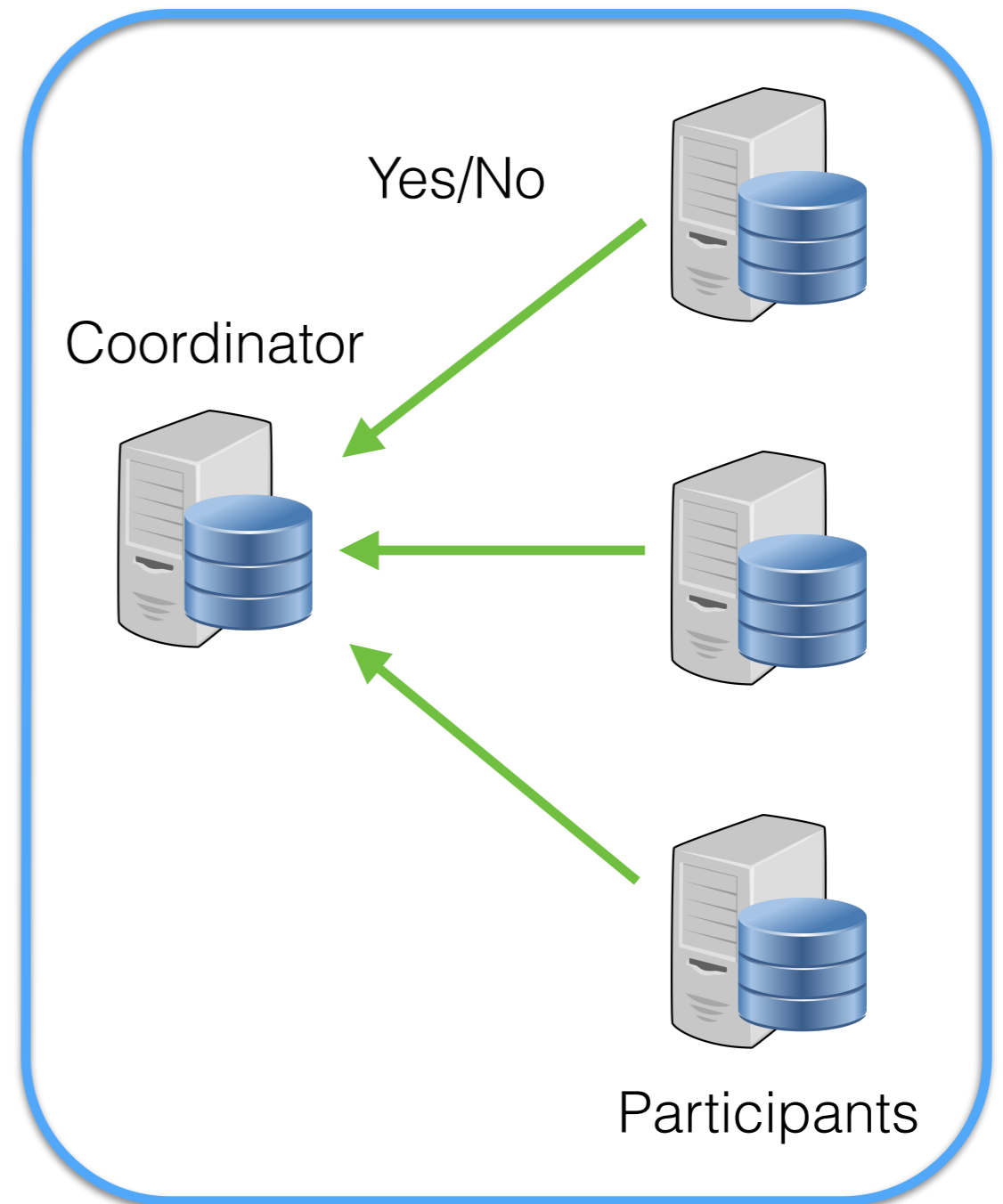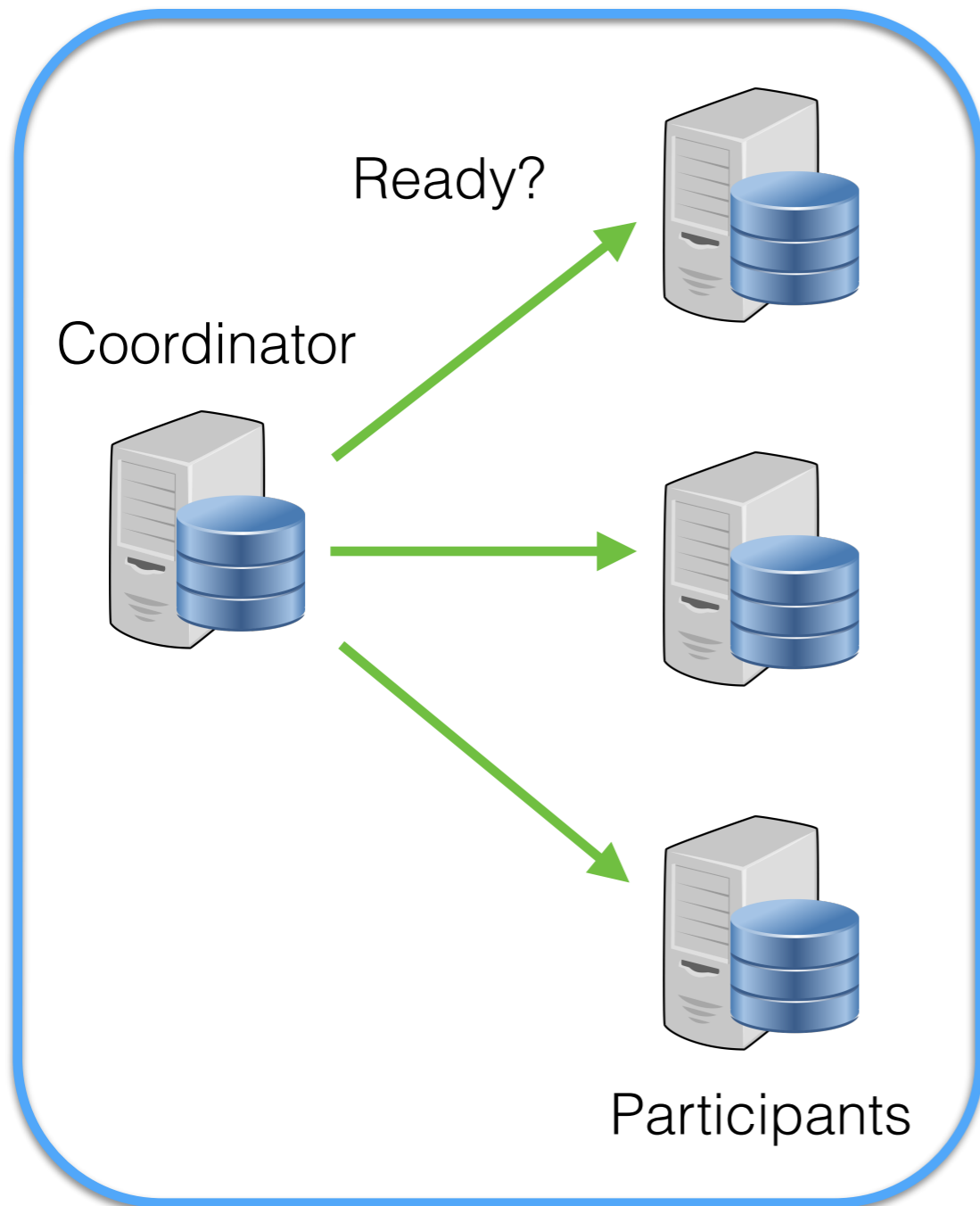
- Consensus Problem!!!!
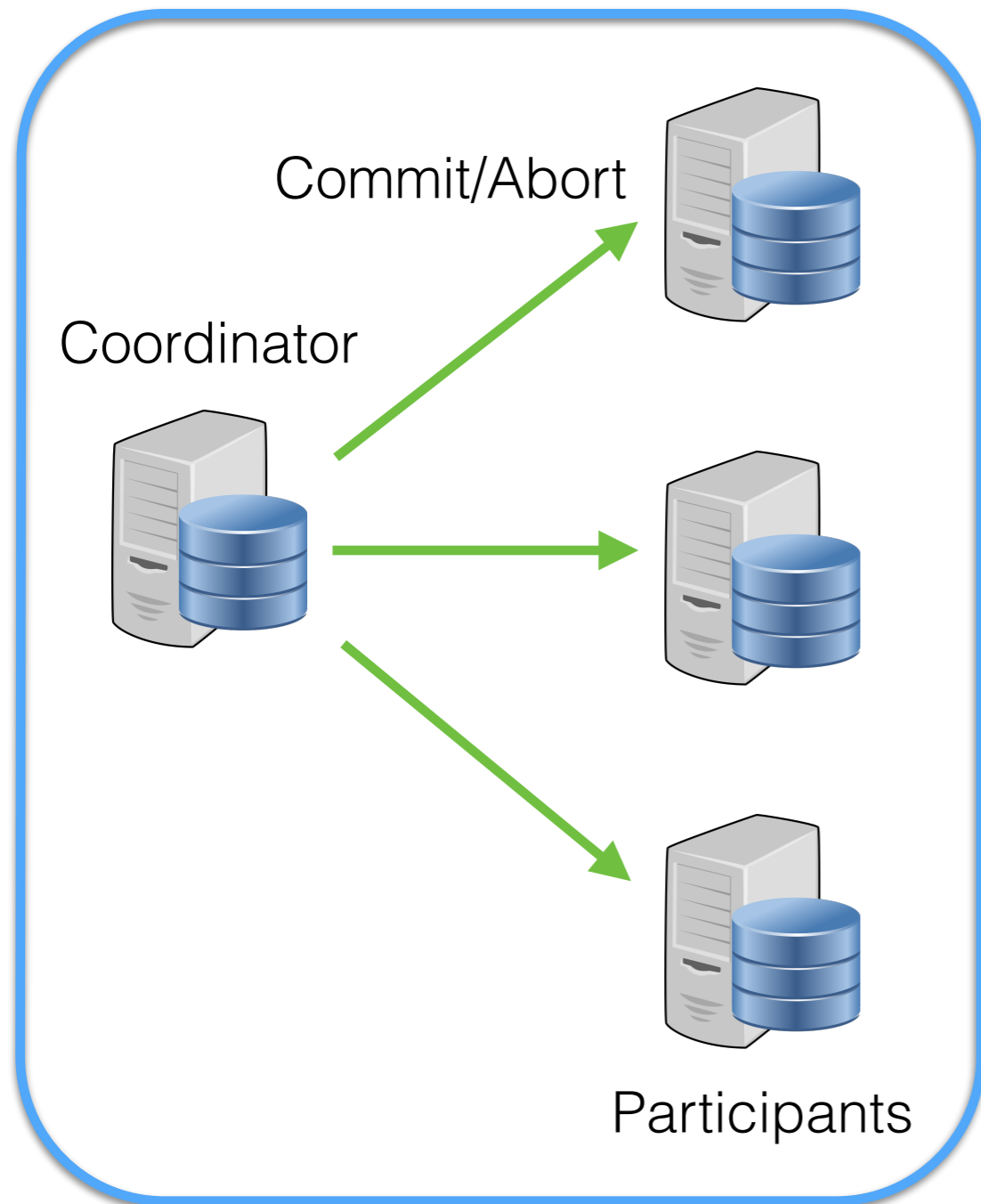
# ACP Conditions

1. All processes <u>that reach a decision</u> reach the same one.

2. A process cannot reverse its decision after it has reached one.

3. The **Commit** decision can only be reached if <u>all</u> processes voted **Yes**.

4. If there are no failures **and** all processes voted **Yes**, then the decision will be to **Commit**.

5. Consider any execution containing only failures that the algorithm is designed to tolerate. At any point in this execution, if all existing failures are repaired and no new failures occur for sufficiently long, then all processes will eventually reach a decision.

# Two Phase Commit (2PC)

# Two Phase Commit (2PC)



Commit/Abort

Coordinator

Participants

- Terminates if no faults

- Satisfy conditions 1-4

- Does not satisfy condition 5 (as expected)

- *3n* messages for *n* nodes

- What happens with faults?

# 2PC - Phase 1

1: coordinator sends **ready** to all participants

1: if a participant receives **ready** from the coordinator:

2:      if it is ready to commit:

3:          send **yes** to the coordinator

4:      else:

5:          send **no** to the coordinator

# 2PC - Phase 2

1:    if the coordinator receives only **yes** messages:

2:          send **commit** to all nodes

3:    else

4:          send **abort** to all nodes


1:    if a node receives **commit** from the coordinator:

2:          commit the transaction

3:    else (**abort** received)

4:          abort the transaction

# 2PC and Faults

- <u>Fail-stop</u> model: We assume that a failed node does not re-emerge

- Failures are detected (instantly)

  - <u>Timeouts</u> are used in practical systems to detect failures (FLP theorem does not apply there!)

- When a waiting process is interrupted by a timeout, the process must take special action, called a <u>timeout action</u>.
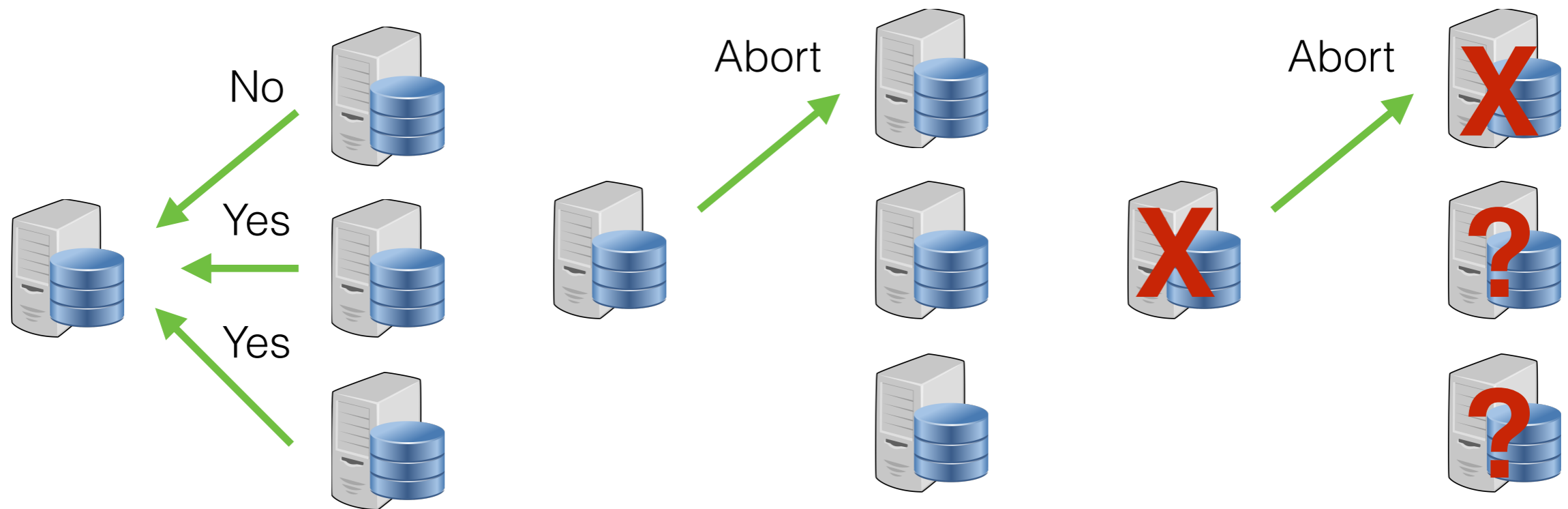
# Timeout Actions

- Coordinator waiting: no decision taken yet, can decide to send **abort** or **commit** to participants

- Participant waiting for ready: can unilaterally decide **abort** before it sends **yes**

- Participant waiting for **commit/abort**: has already voted **yes**, so it can not unilaterally decide. 2PC blocks.

- We need a <u>safety mechanism</u> to deal with block (not part of 2PC).

# Safety Mechanism

- All participants must know each other

- Coordinator failures is detected

- A new coordinator is selected

- It must ask the other participants if a participant has already received a **commit**

- A participant that has received a **commit** replies **yes**, otherwise it sends **no** and *promises not to accept a commit that may arrive from the old coordinator*

- If some participant replied **yes**, the new coordinator broadcasts **commit**

- What if we have multiple failures???

# 2PC and multiple failures

- Multiple participant failures:

  - No problem as long as the coordinator lives
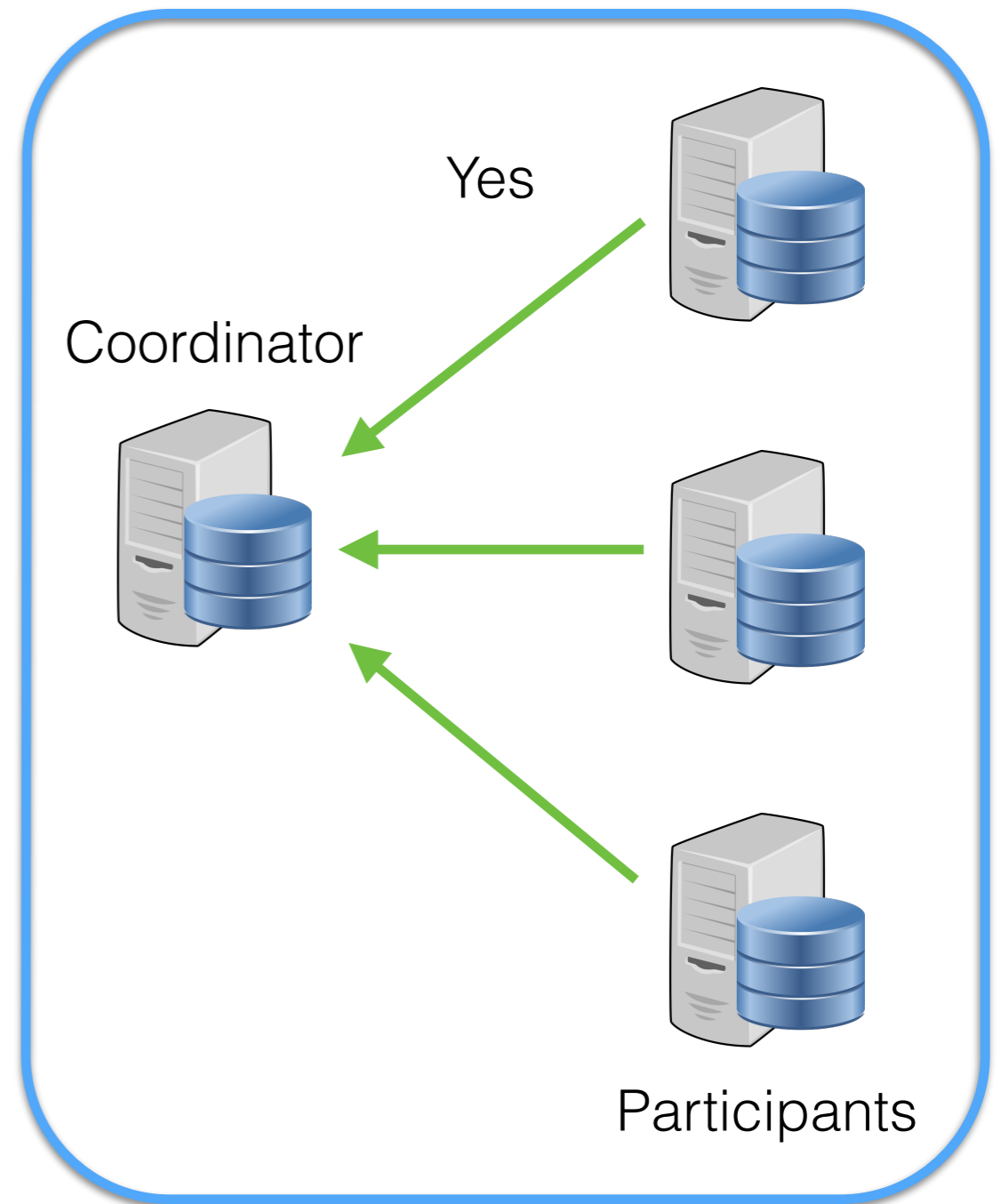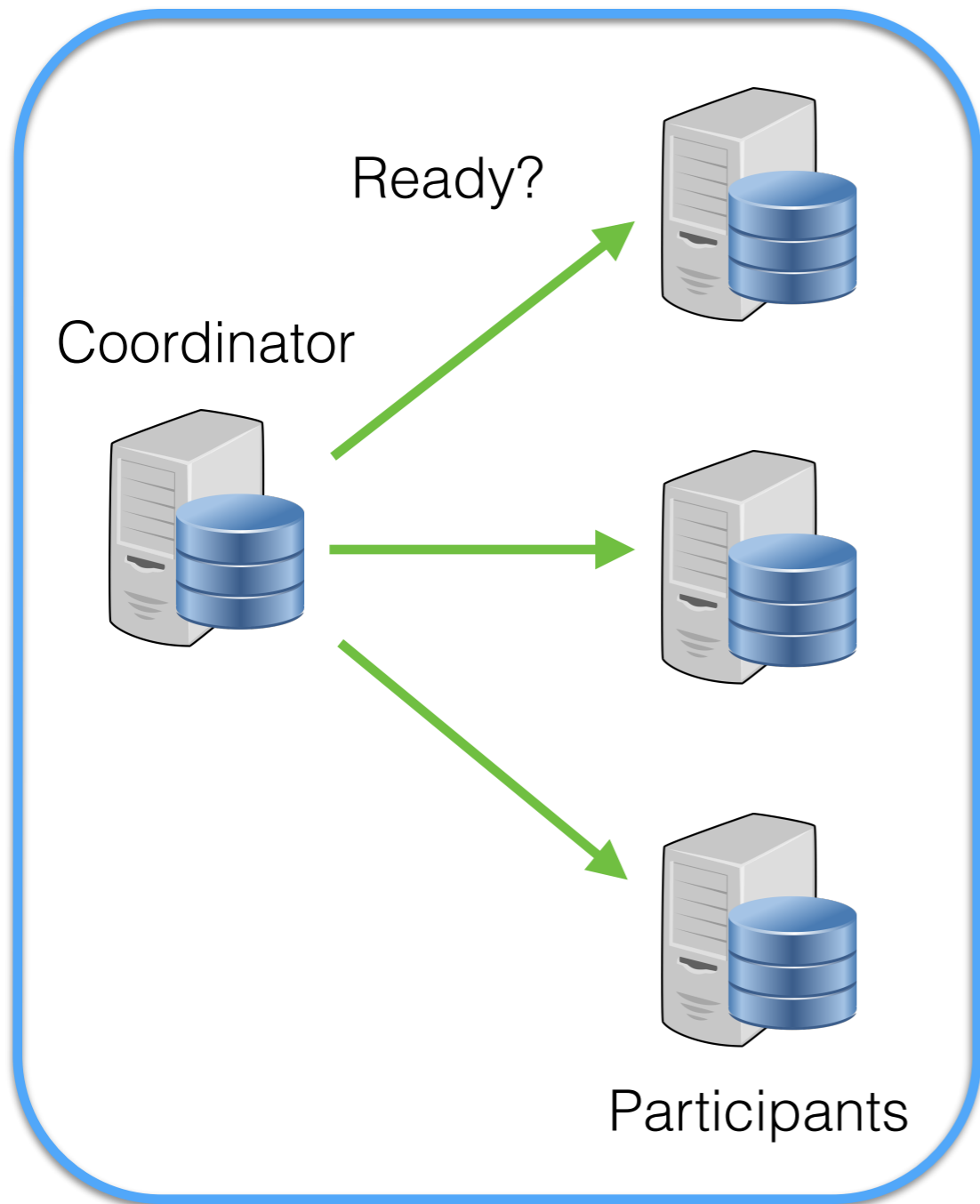
- Coordinator and participant failures:



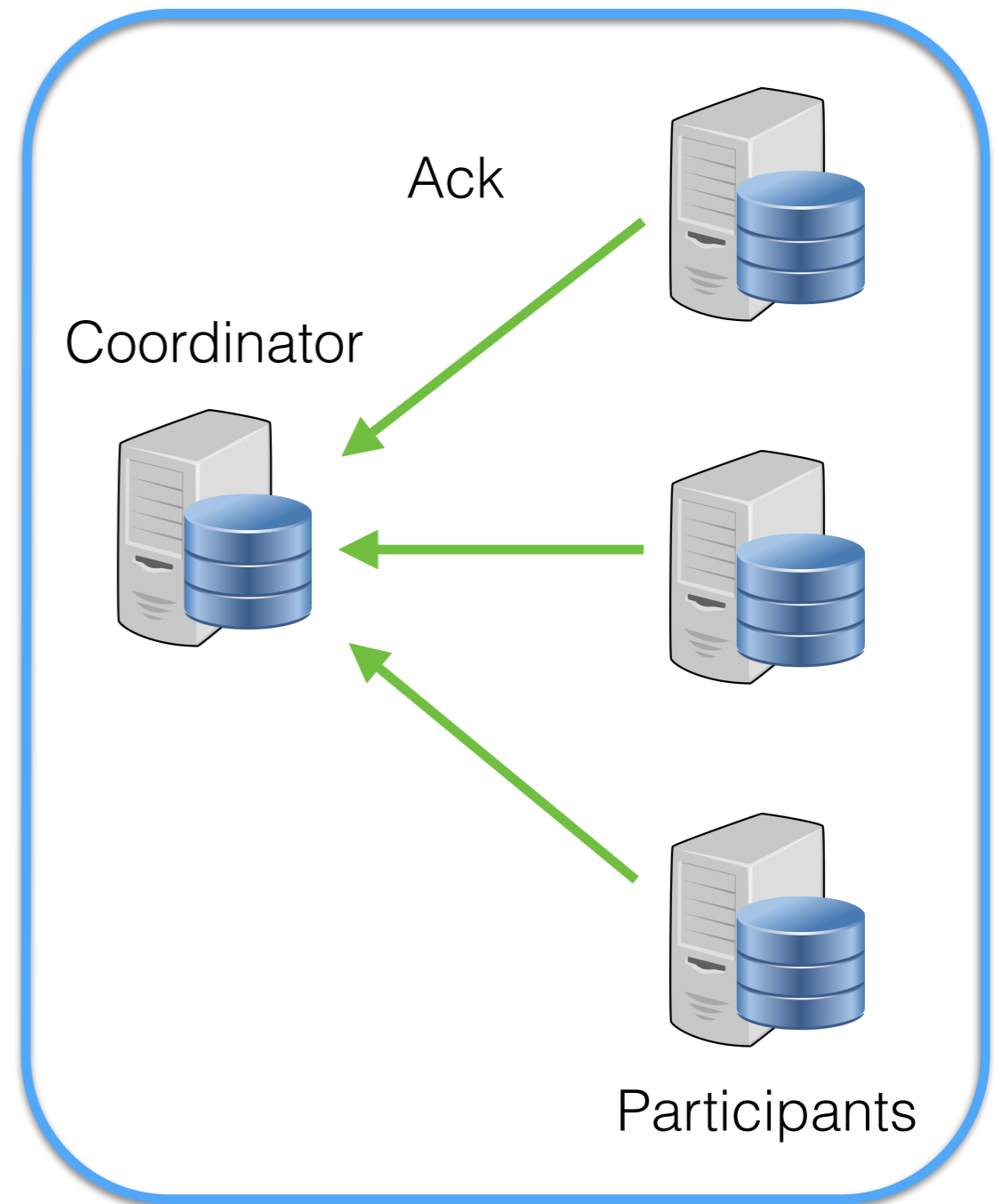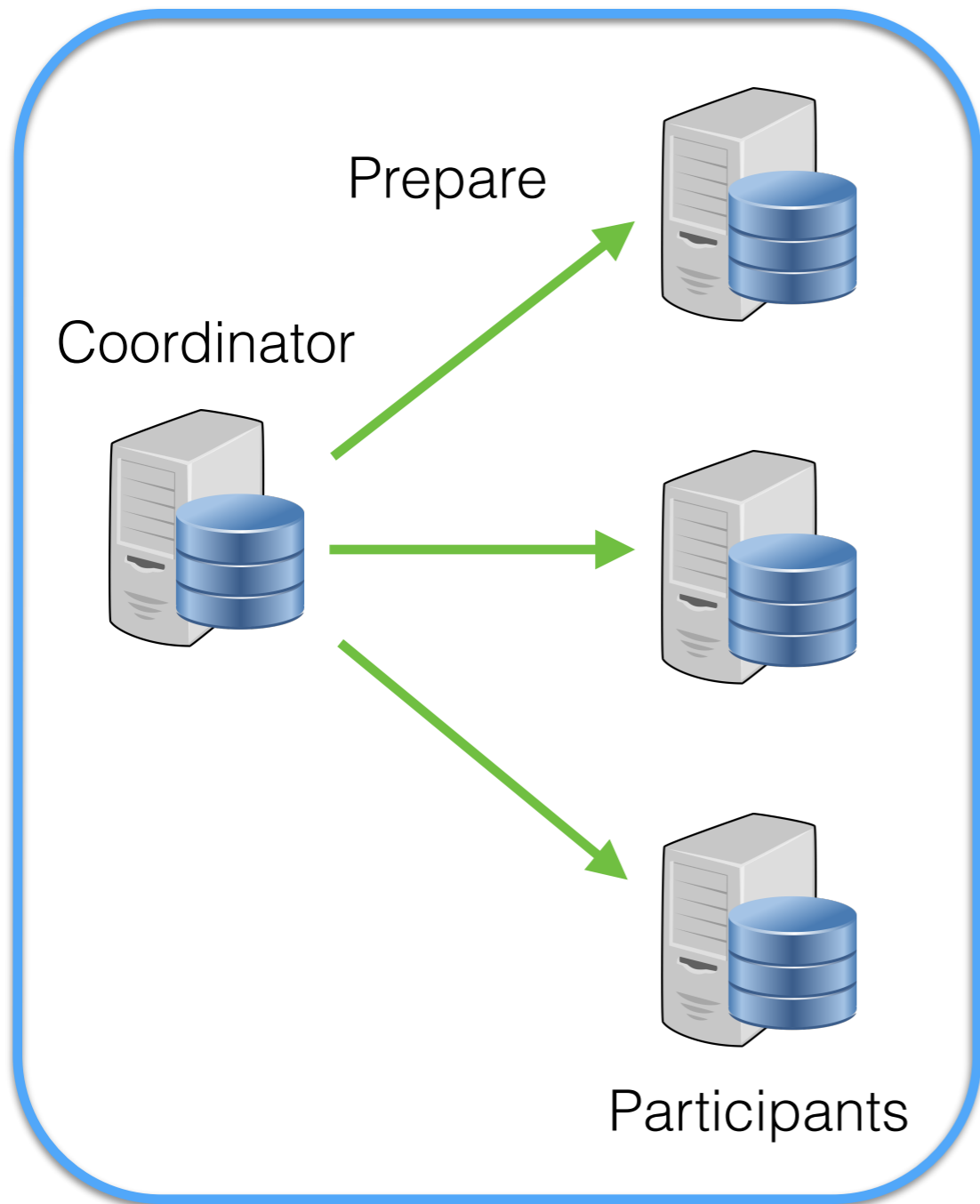- The remaining participants cannot take a decision

# 2PC and multiple failures

- Solution: Add another phase to the protocol!

  - The new phase precedes the commit phase

  - The goal is to inform all participants that all are ready to commit or abort

  - At the end of this phase, every participant knows whether or not all participants want to commit before any participant has actually committed or aborted!

- This protocol is called the three-phase commit protocol (3PC)

# Three Phase Commit (3PC)
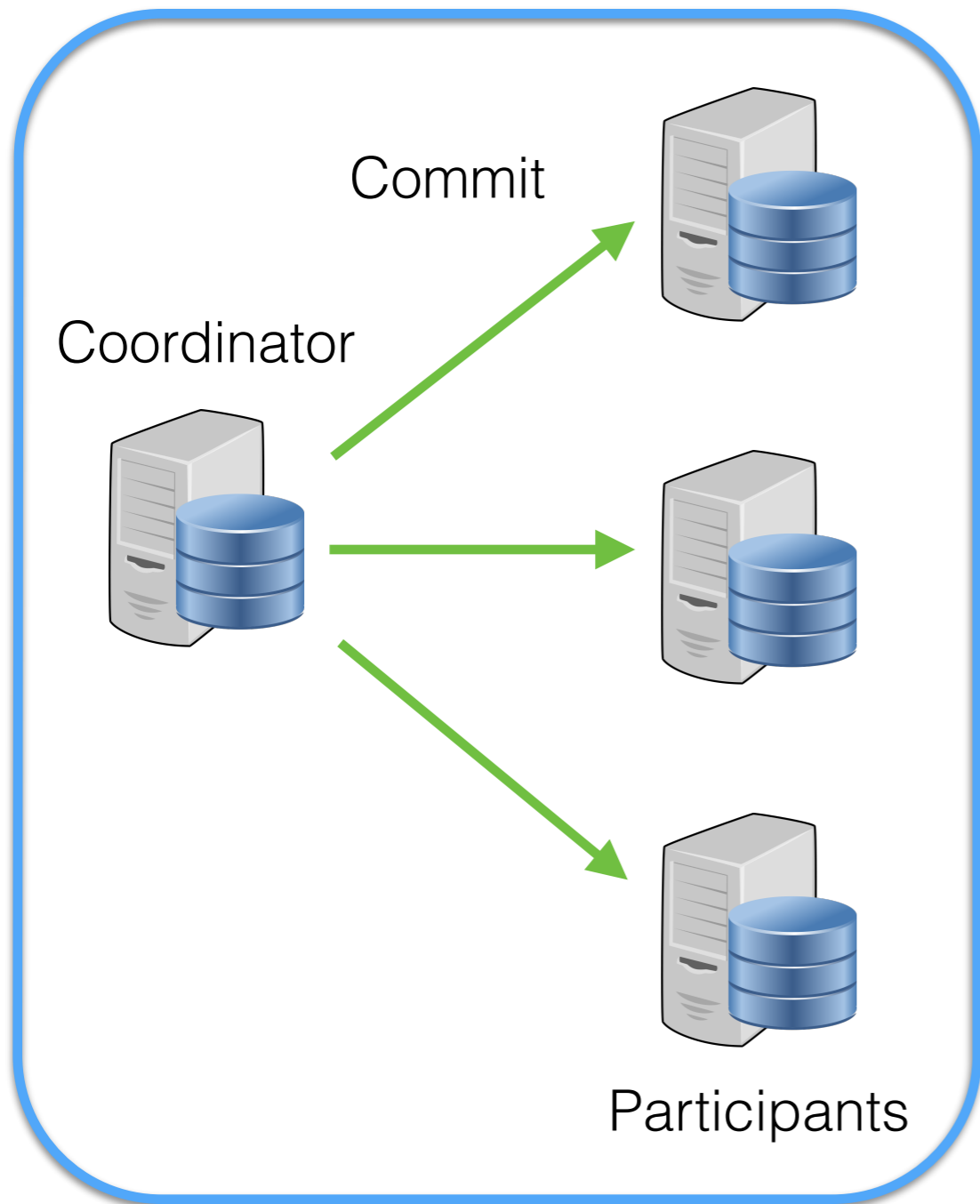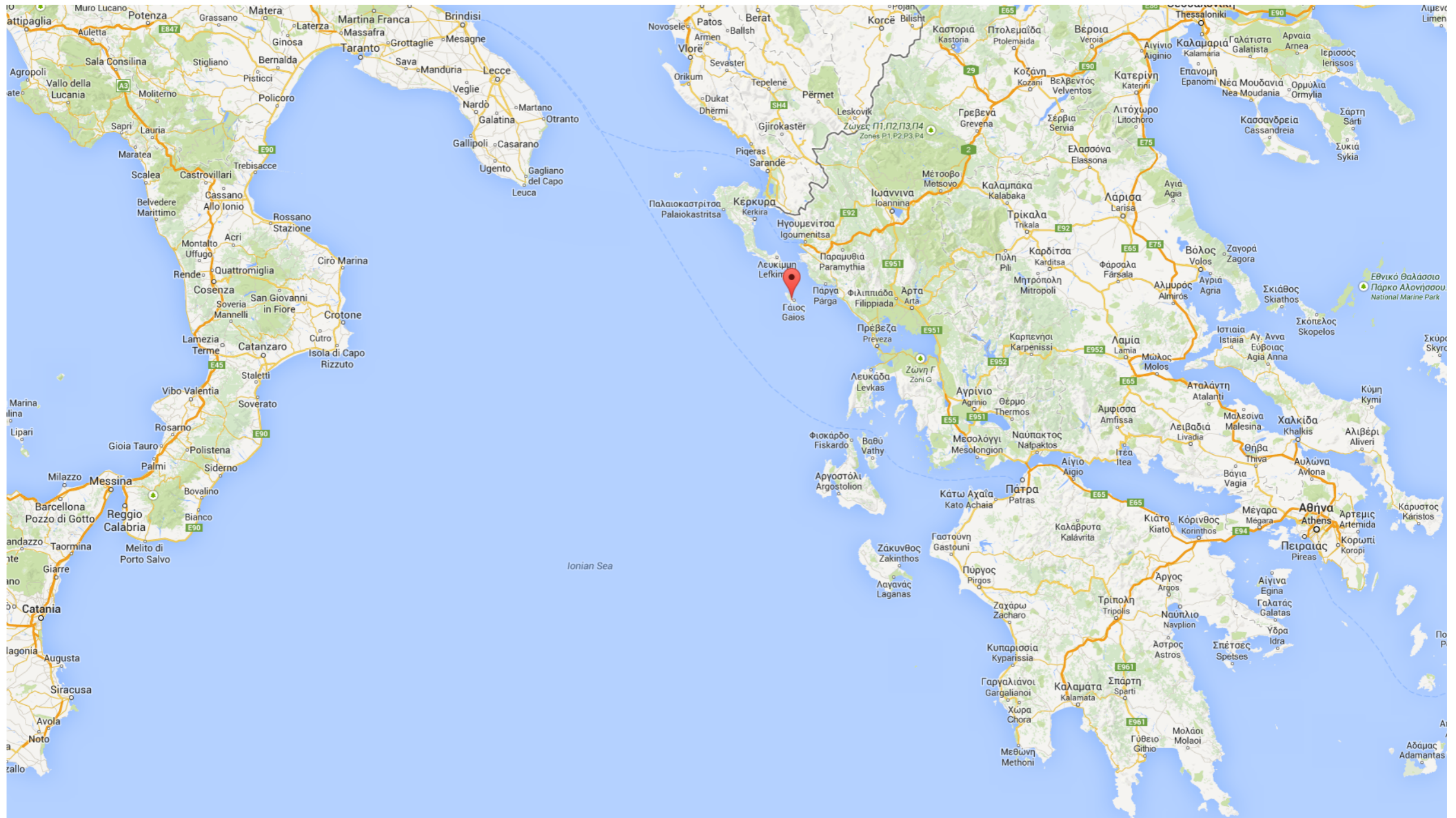


Ready?

Coordinator

Participants

Yes

Coordinator

Participants

# Three Phase Commit (3PC)



Prepare

Coordinator

Participants

Ack

Coordinator

Participants

# Three Phase Commit (3PC)



Commit

Coordinator

Participants

# PAXOS

# PAXOS

- Leslie Lamport. **The part-time parliament**. ACM Transactions on Computer Systems, 16(2):133–169, May 1998.

- Leslie Lamport described in 1990 the algorithm as the solution to a problem of the parliament on a fictitious Greek island called Paxos (not Italy)

- Many readers were so distracted by the description of the activities of the legislators, they did not understand the meaning and purpose of the algorithm. The paper was rejected.

- Leslie Lamport refused to rewrite the paper. He later wrote that he "was quite annoyed at how humorless everyone working in the field seemed to be"

- After a few years, some people started to understand the importance of the algorithm

- After eight years, Leslie Lamport submitted the paper again, basically unaltered. It got accepted!



Leslie Lamport
ACM Turing Award 2014

# Consensus (again)

- We have a collection of processes

- Each process can propose a value

- A single value among the proposed values is chosen

- If no value is proposed, no value should be chosen

- If a value has been chosen, processes should be able to learn the chosen value

# System Model

- <u>Asynchronous</u> communications

- <u>Non-byzantine</u> failures:

  - Nodes crash

  - Messages can be lost, duplicated arbitrarily late

  - No corrupted messages

- <u>Fail-recover</u>: crashed nodes may recover later
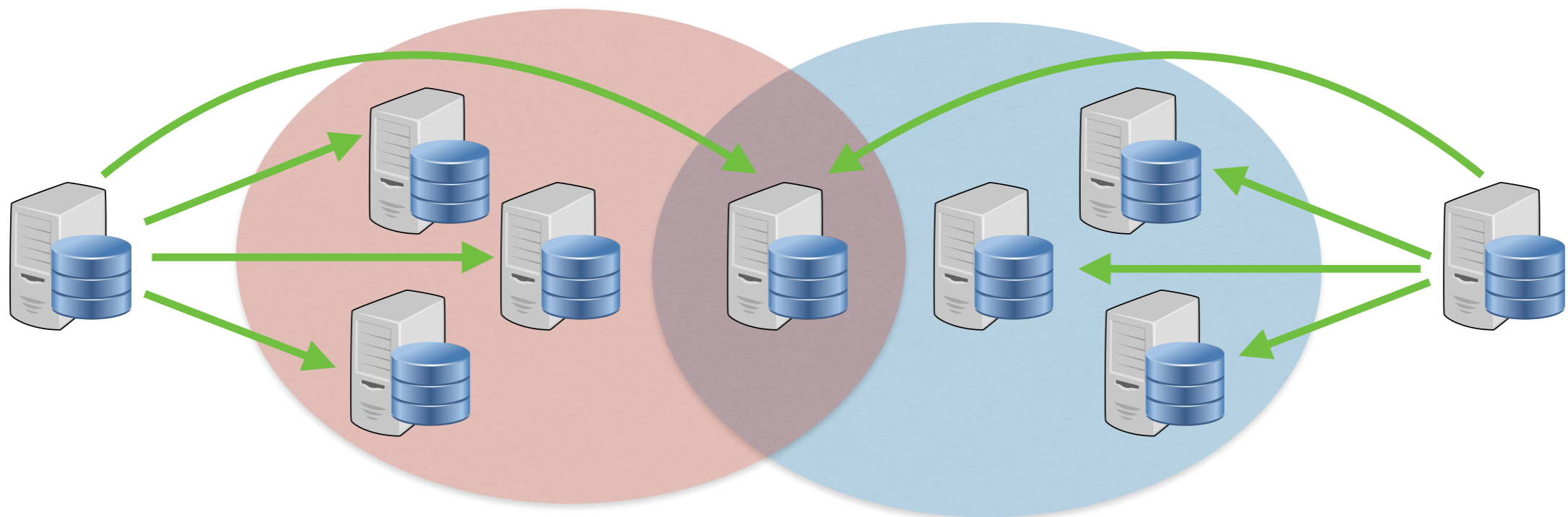
# Roles

- Each nodes has one or more of 3 roles

- <u>Proposer</u>: a node that can propose a certain value for acceptance

- <u>Acceptor</u>: a node that receives proposals from proposers and that can either accept or reject a proposal

- <u>Learner</u>: a node not involved in the decision process that wants to know the final result of the decision process

- We will assume all nodes act both as proposer and as acceptor

# Acceptors

- A single acceptor (a.k.a. coordinator) can fail and block the whole procedure

- There is not a coordinator, but multiple acceptors.

- An acceptor may **accept** a single value (e.g., express a single vote)

- How many acceptors do we need?

# Majority

- To ensure that a single value is chosen, any majority (50% +1) of acceptors is enough

- The intersection of two majorities is not empty

- An acceptor may **accept** a single value

- A value is **chosen** when a majority of acceptors has accepted it

- If a majority choses a value, no other majority can chose a different value

- If an acceptor crashes, chosen value still available

# Acceptors

Since there can be a single proposer, the algorithm must guarantee that
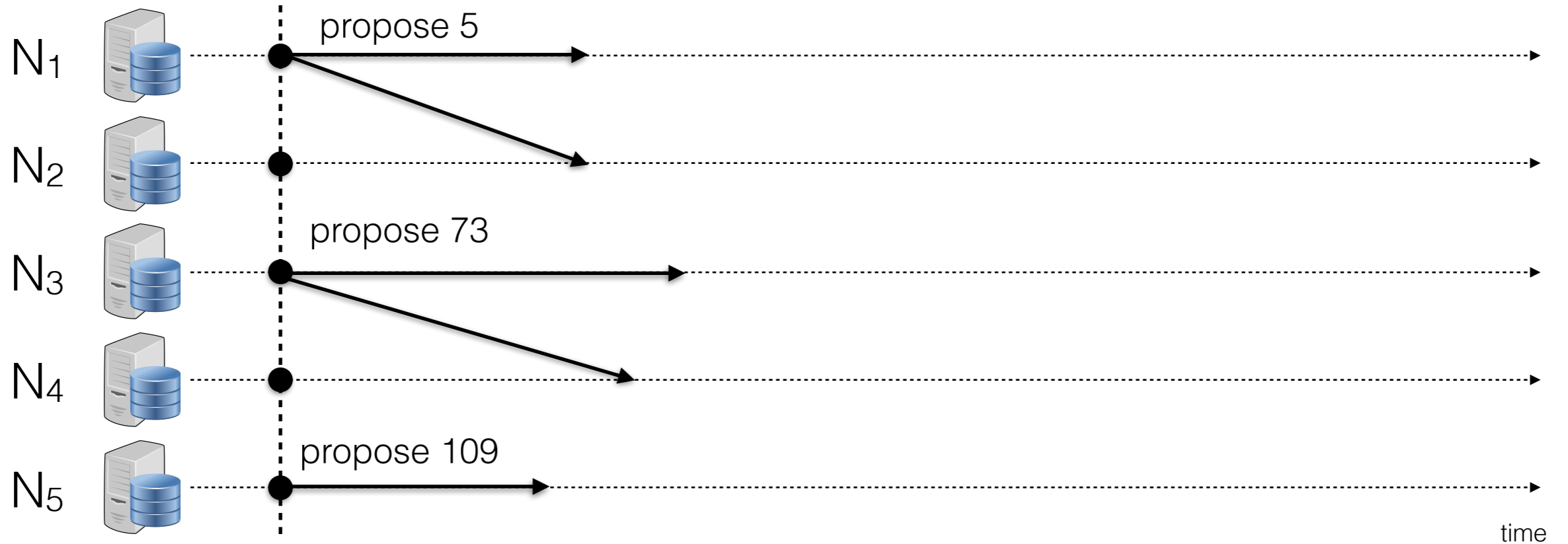
## P1: An acceptor must accept the first proposed value it receives

Since there can be a single proposer, the algorithm must guarantee that a single value is chosen only when it is accepted by a majority of acceptors
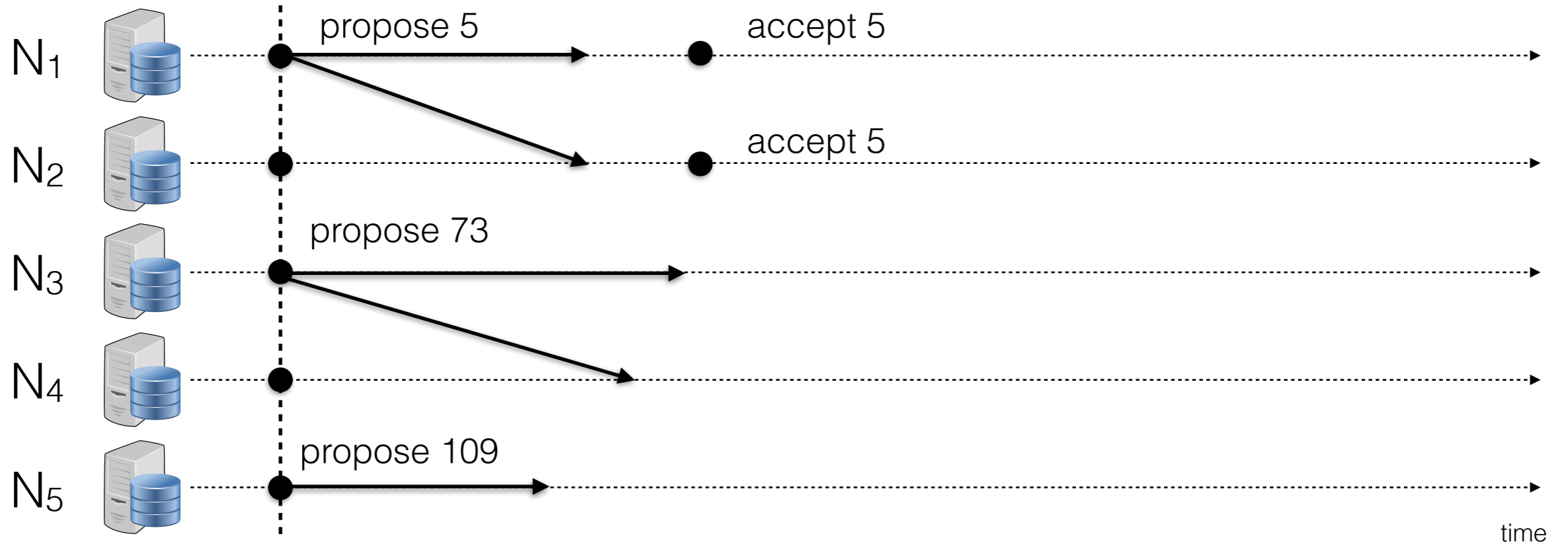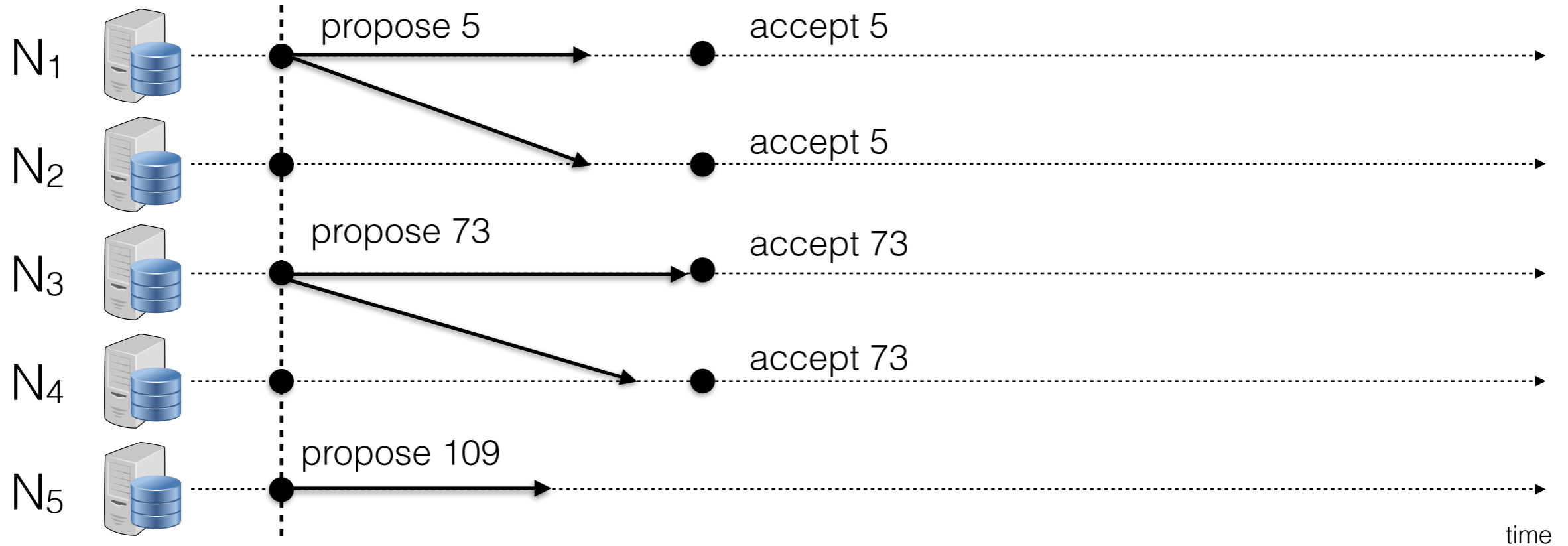
# Split Votes



$N_1$
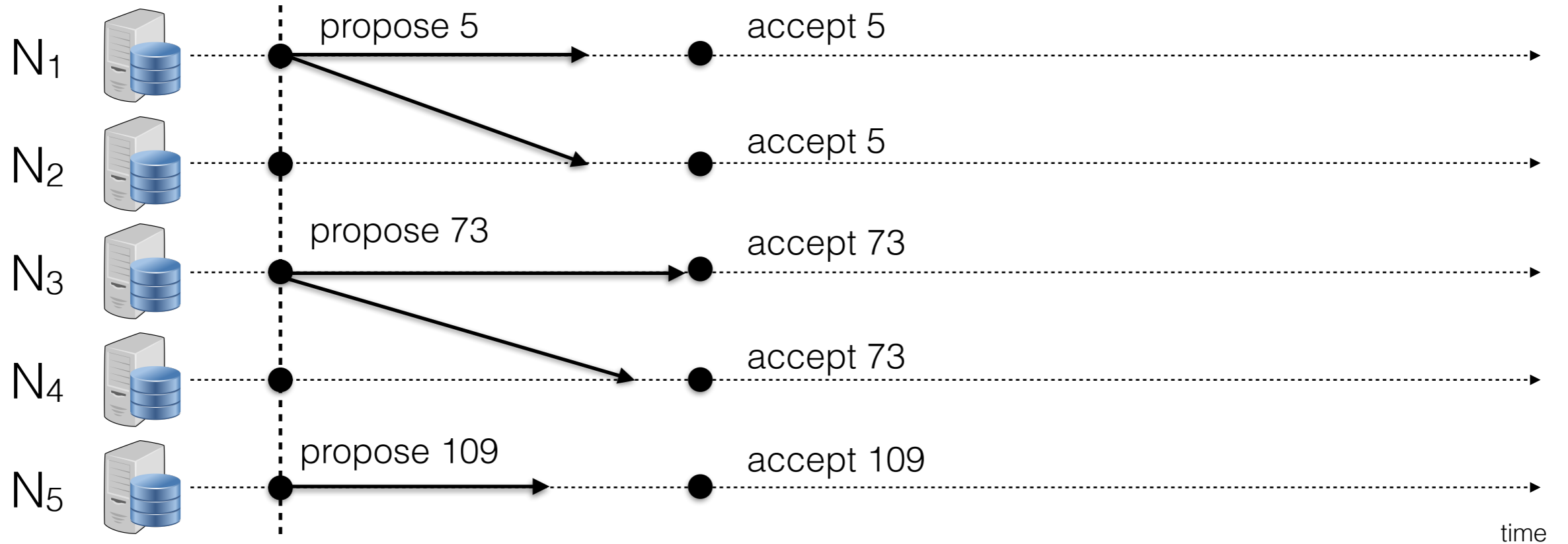
$N_2$

$N_3$

$N_4$

$N_5$

time

# Split Votes



propose 5

propose 73

propose 109

time

# Split Votes



N₁ — propose 5 → accept 5

N₂ — accept 5

N₃ — propose 73

N₄

N₅ — propose 109

time

# Split Votes

# Split Votes

# Split Votes



N₁ propose 5 → accept 5
N₂ accept 5
N₃ propose 73 → accept 73
N₄ accept 73
N₅ propose 109 → accept 109

No majority

time

# Split Votes



N1  propose 5 → accept 5
N2            accept 5
N3  propose 73 → accept 73
N4            accept 73
N5  propose 109 → accept 109

No majority

time

**An acceptor must accept the first proposed value it receives**

**(and can accept more than one proposed value)**

# Split Votes



N$_1$ — propose 5 → accept 5

N$_2$ — accept 5

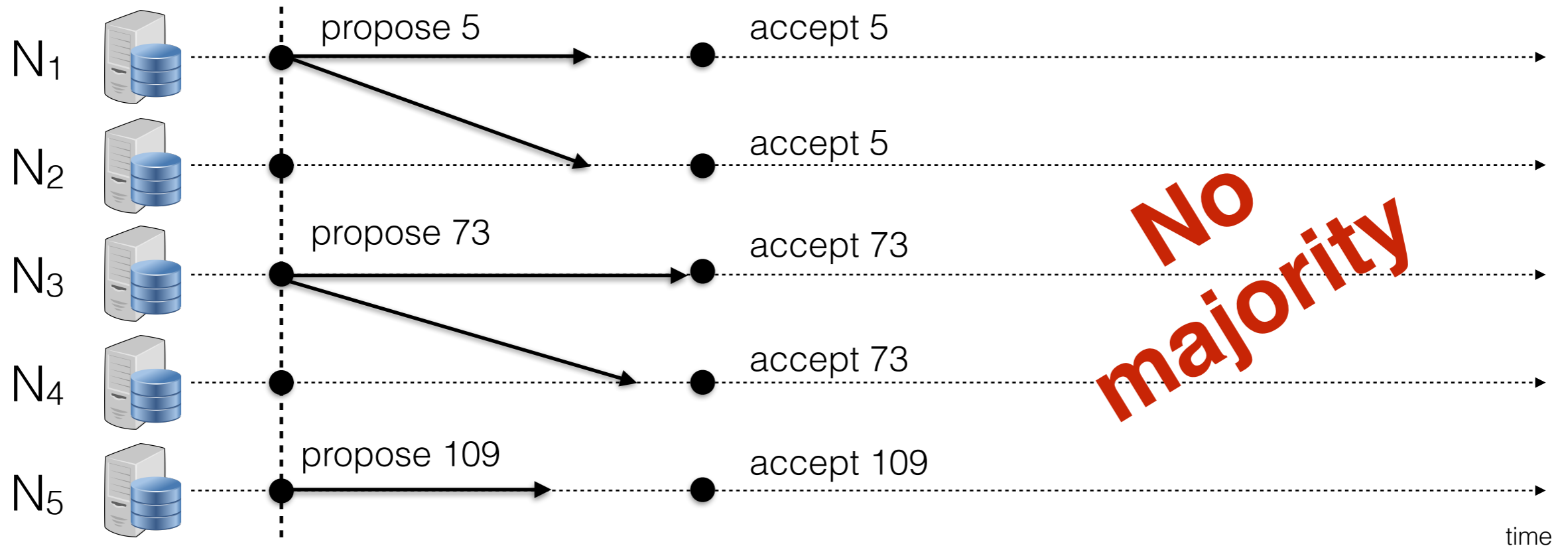N$_3$ — propose 73 → accept 73

N$_4$ — accept 73

N$_5$ — propose 109 → accept 109

No majority

time

**An acceptor must accept the first proposed value it receives**

**(and can accept more than one proposed value)**

An acceptor can **accept** any number of proposed values, but an accepted

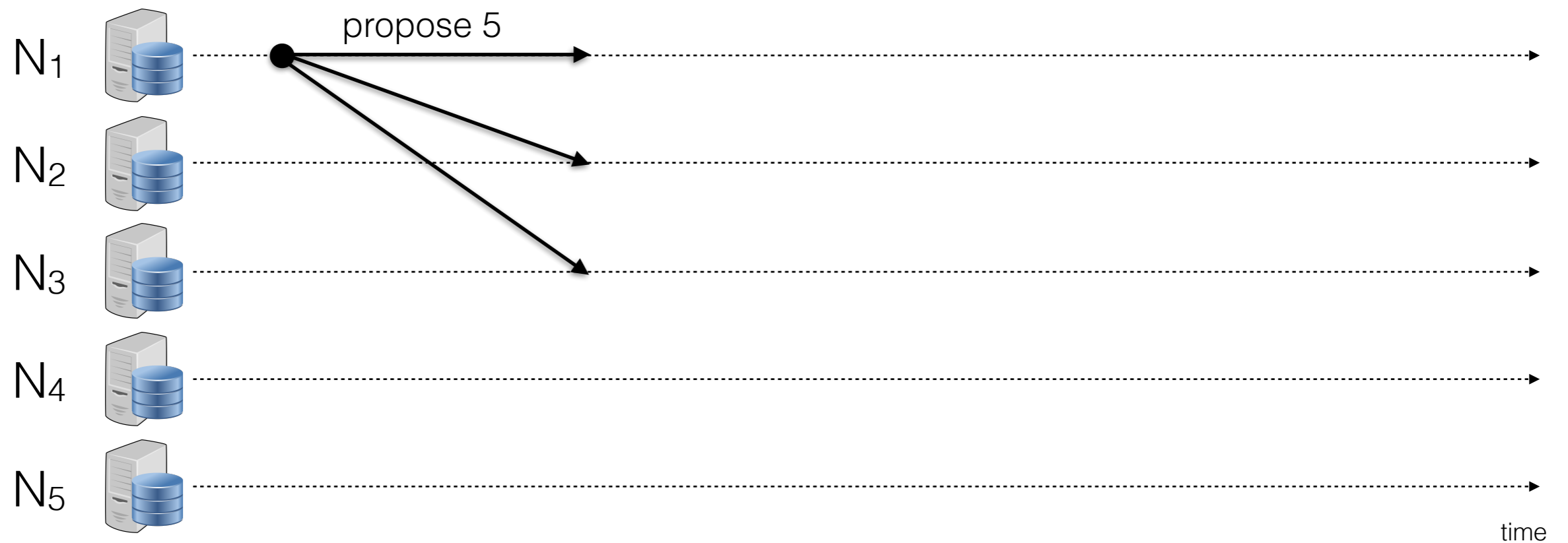proposed value may not necessarily be **chosen**.

# Conflicting choices

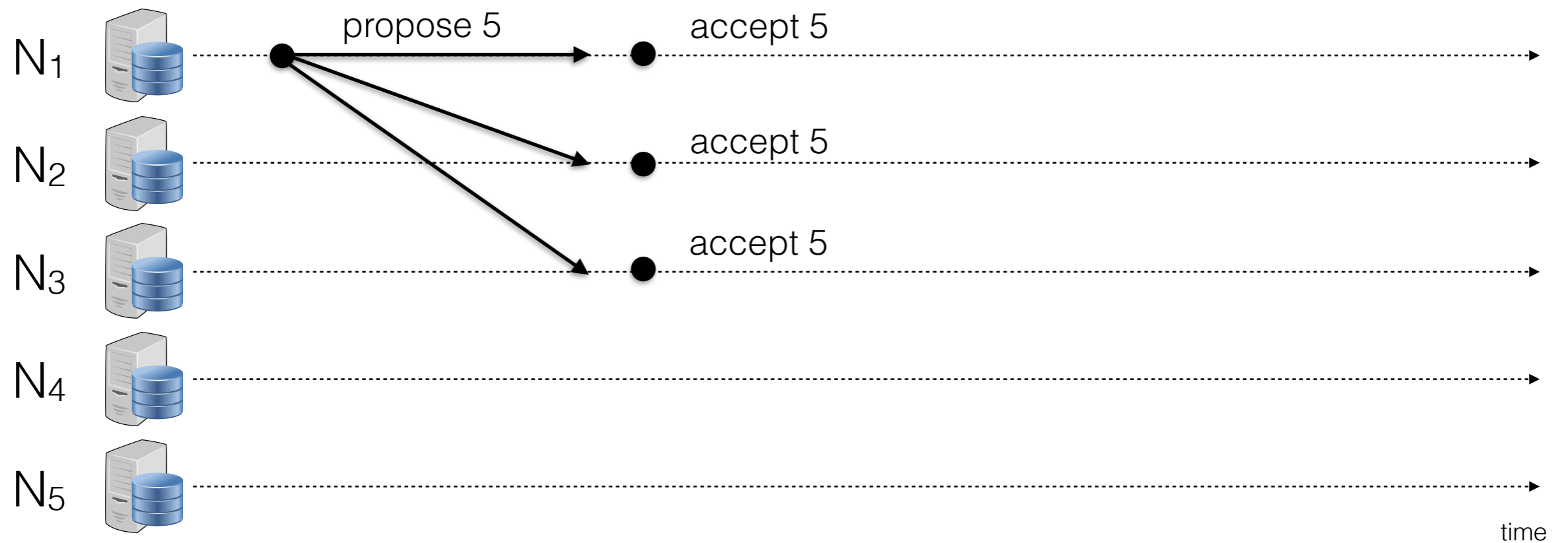Should acceptors accept _every_ proposed value? No
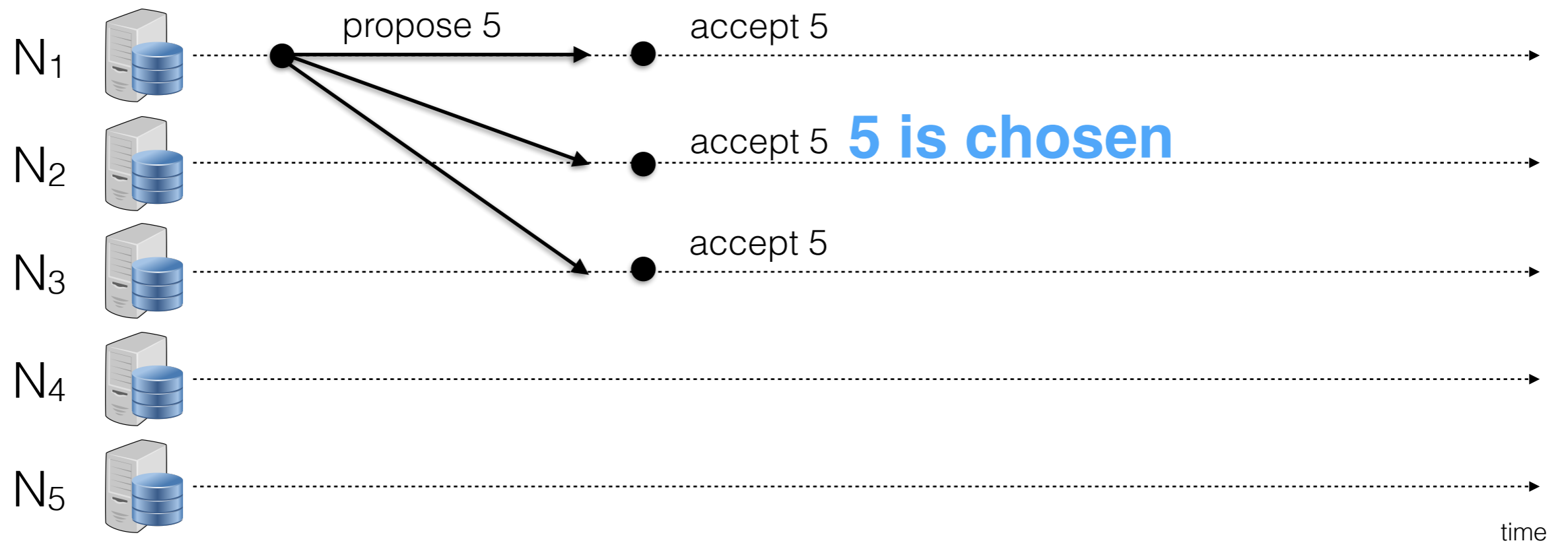


$N_1$

$N_2$

$N_3$

$N_4$

$N_5$

time

# Conflicting choices

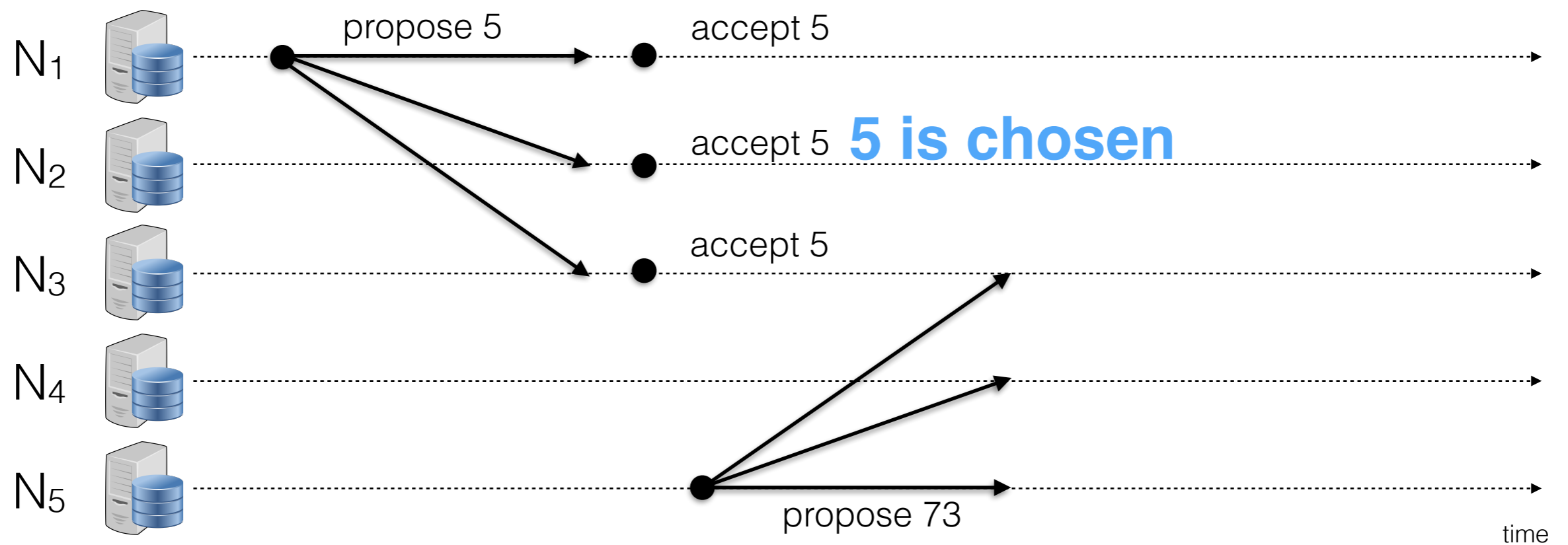Should acceptors accept __every__ proposed value? No

# Conflicting choices

Should acceptors accept <u>every</u> proposed value? No

$N_1$

propose 5 — accept 5

$N_2$

accept 5

$N_3$

accept 5

$N_4$

$N_5$

time

# Conflicting choices

Should acceptors accept <u>every</u> proposed value? No

$N_1$ — propose 5 → accept 5

$N_2$ — accept 5 **5 is chosen**

$N_3$ — accept 5

$N_4$

$N_5$

time

# Conflicting choices

Should acceptors accept <u>every</u> proposed value? No



$N_1$ — propose 5 → accept 5

$N_2$ — accept 5 **5 is chosen**

$N_3$ — accept 5

$N_4$

$N_5$ — propose 73

time

# Conflicting choices

Should acceptors accept <u>every</u> proposed value? No

$N_1$ propose 5 accept 5

$N_2$ accept 5 **5 is chosen**

$N_3$ accept 5 accept 73

$N_4$ accept 73

$N_5$ propose 73 accept 73

time

# Conflicting choices

Should acceptors accept <u>every</u> proposed value? No



$N_1$  propose 5  accept 5

$N_2$  accept 5  **5 is chosen**

$N_3$  accept 5  accept 73

$N_4$  accept 73  **73 is chosen**

$N_5$  accept 73

propose 73

time

# Conflicting choices

Should acceptors accept <u>every</u> proposed value? No



Once a value has been **chosen**, proposers must **propose** that same value

# Conflicting choices

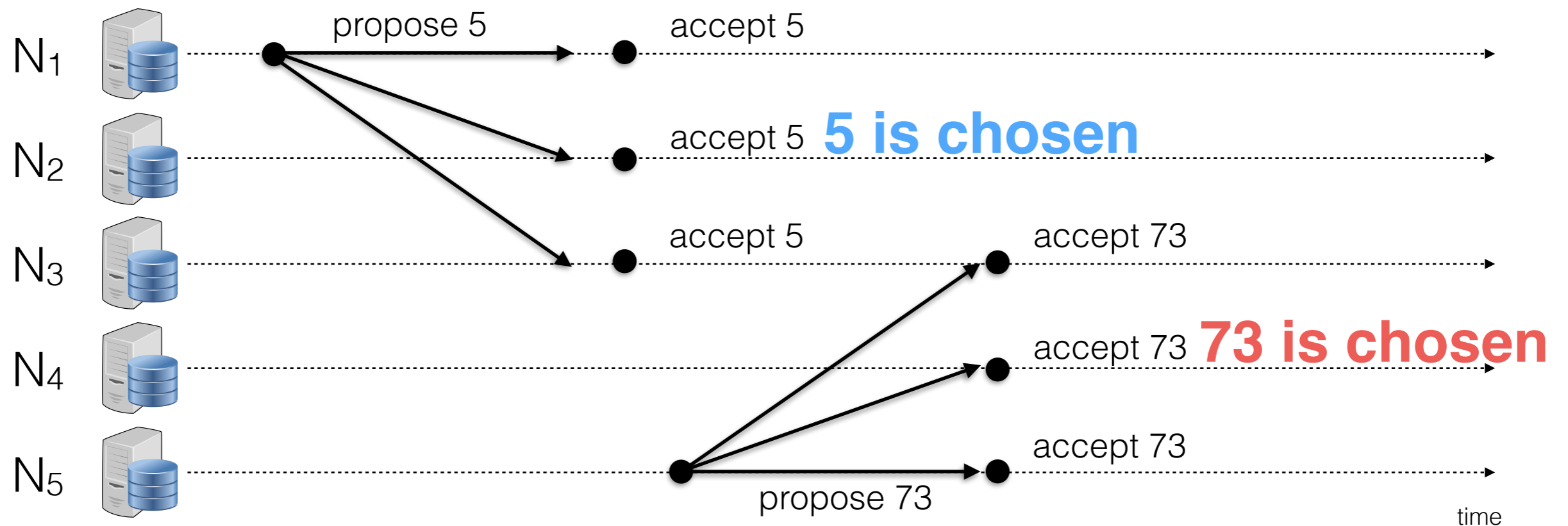Should acceptors accept <u>every</u> proposed value? No



Once a value has been **chosen**, proposers must **propose** that same value

**We need a two-phase protocol**

# Conflicting choices



$N_1$

$N_2$

$N_3$

$N_4$

$N_5$

time

Once a value has been **chosen**, older proposed values must be ignored

**We need a to order proposed values and reject old ones**

# Conflicting choices



Once a value has been **chosen**, older proposed values must be ignored

**We need a to order proposed values and reject old ones**

# Conflicting choices



Once a value has been **chosen**, older proposed values must be ignored

**We need a to order proposed values and reject old ones**

# Conflicting choices



Once a value has been **chosen**, older proposed values must be ignored
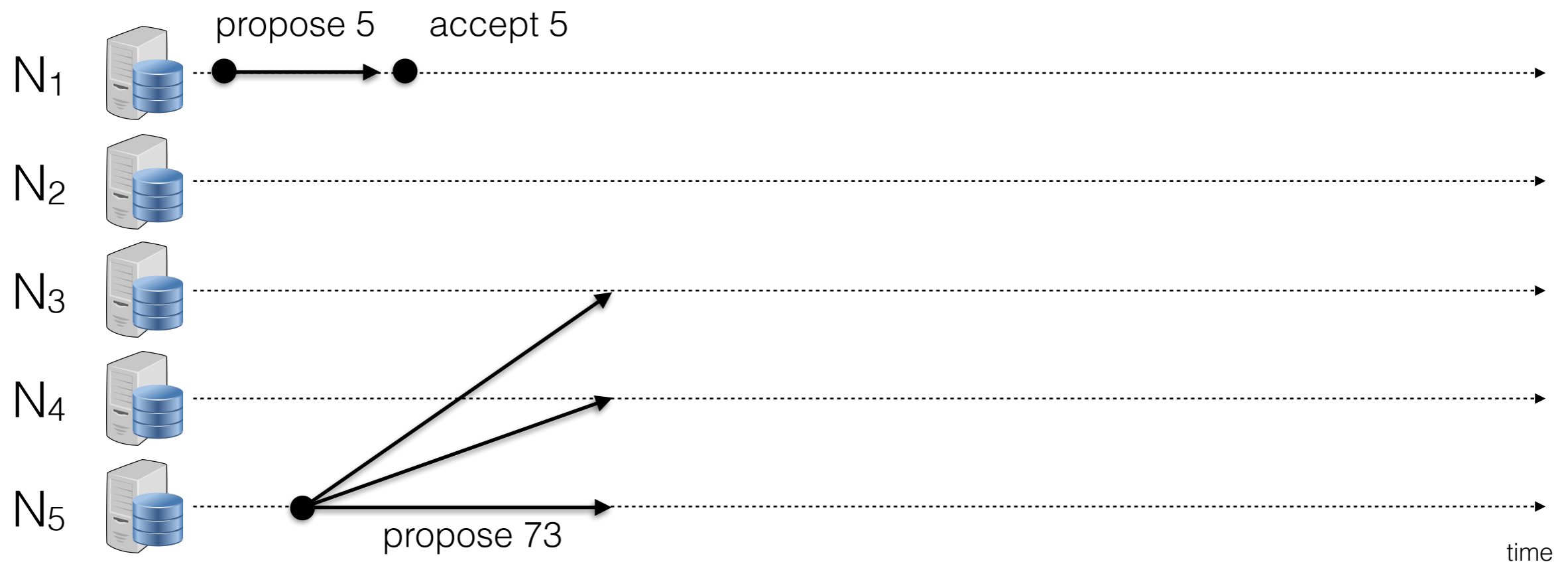
**We need a to order proposed values and reject old ones**
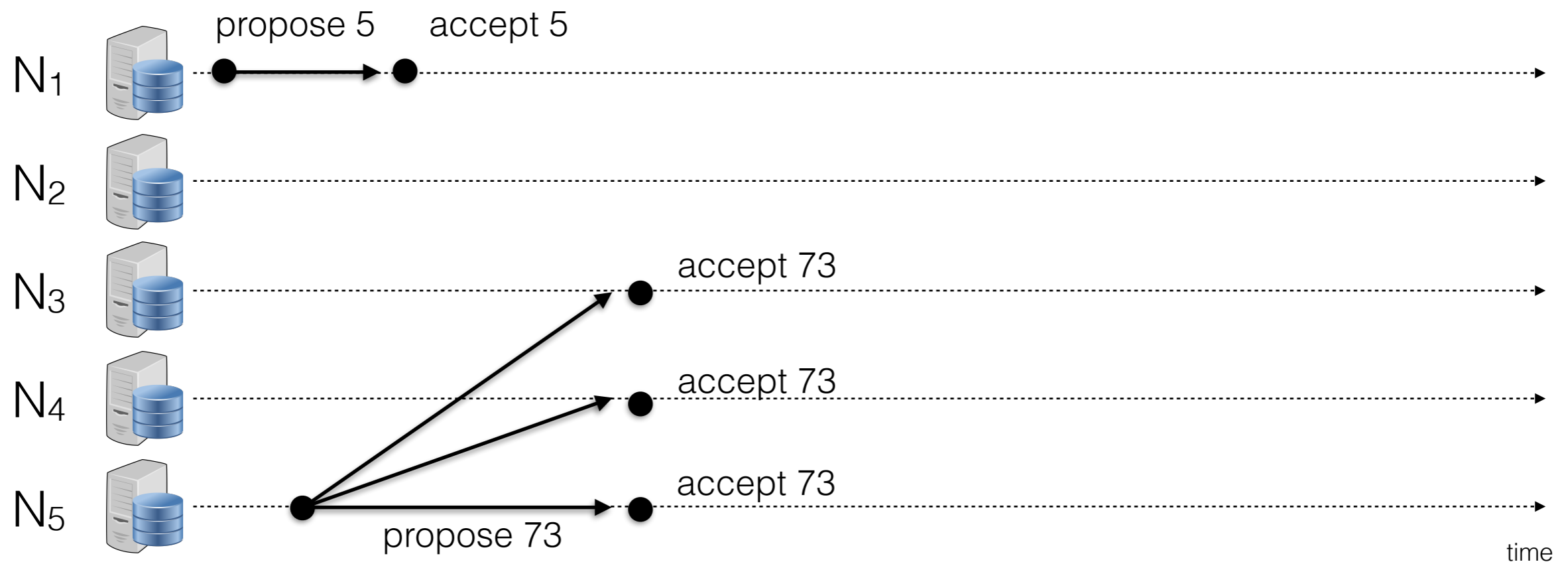
# Conflicting choices



propose 5    accept 5

$N_1$

$N_2$

$N_3$    accept 73

**73 is chosen**

$N_4$    accept 73

$N_5$    accept 73

propose 73

time

Once a value has been **chosen**, older proposed values must be ignored

**We need a to order proposed values and reject old ones**

# Conflicting choices



Once a value has been **chosen**, older proposed values must be ignored
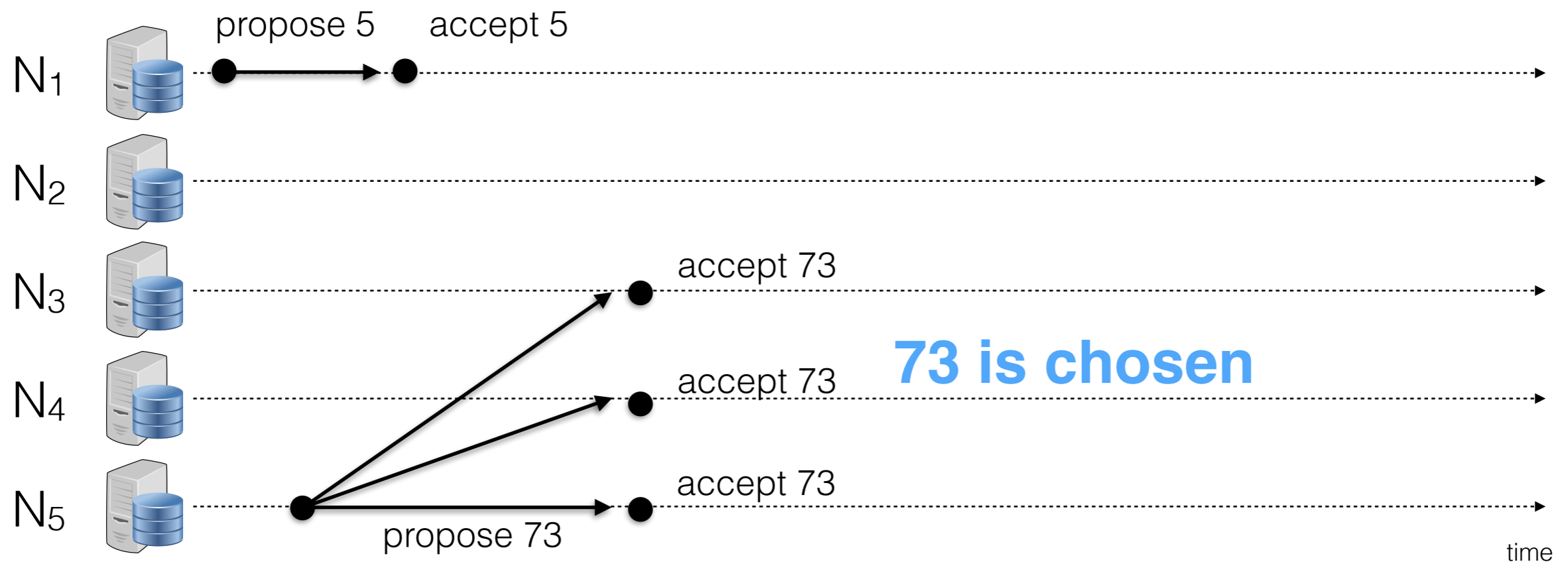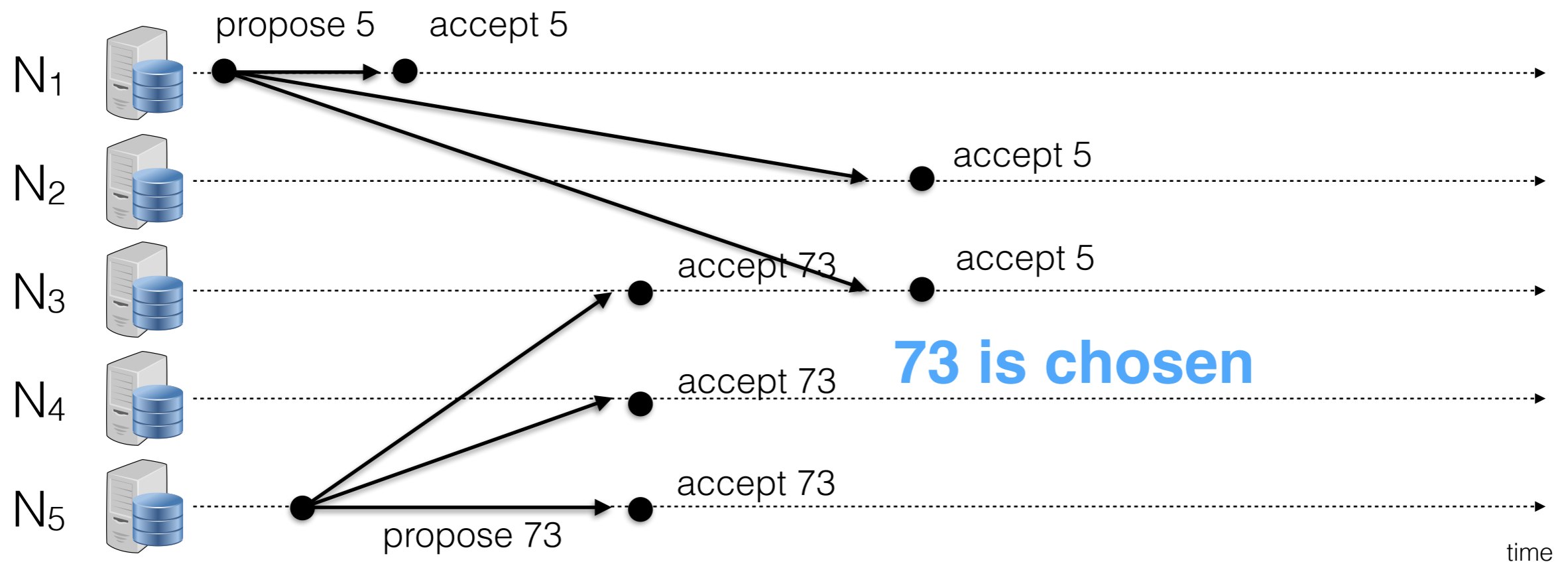
**We need a to order proposed values and reject old ones**

# Conflicting choices



Once a value has been **chosen**, older proposed values must be ignored

**We need a to order proposed values and reject old ones**
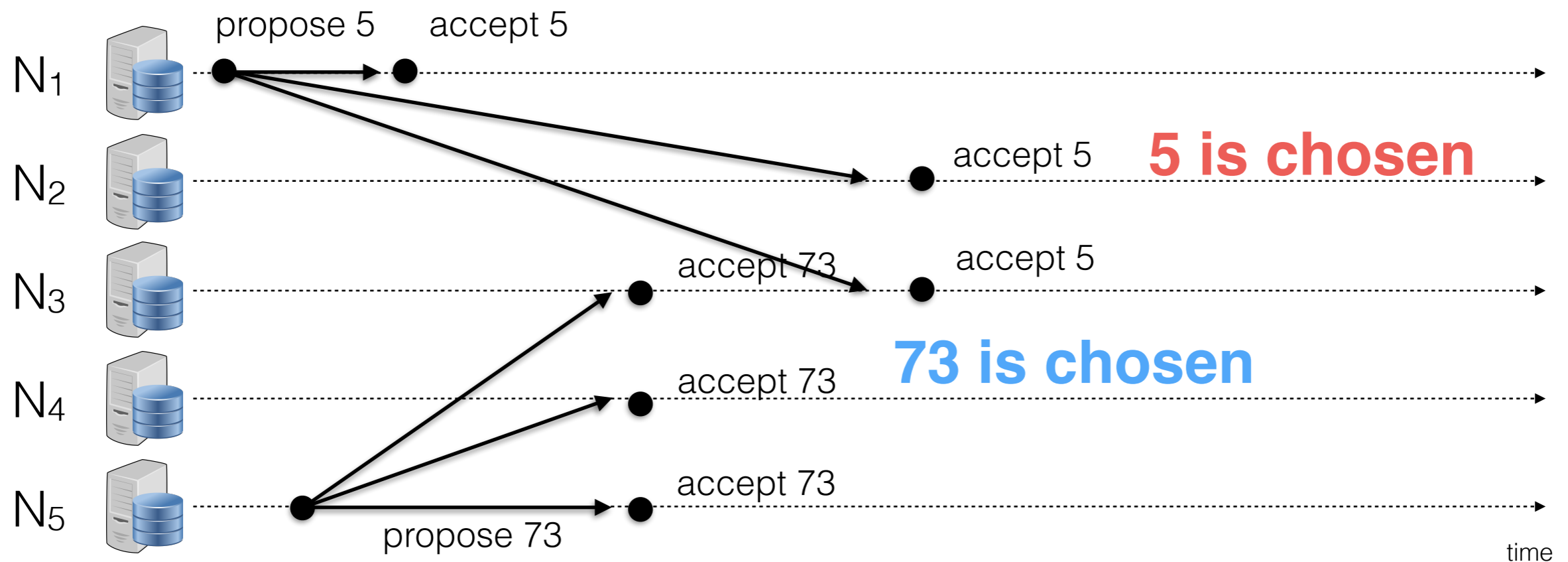
# Proposals

- A **proposal** *(v, n)* consists in the **proposed value** *v* and a **proposal number** *n*

- Whenever a proposer issues a new proposal, it chooses a <u>strictly-increasing</u> proposal number

- For a given proposer, a new proposal number must be greater than anything it has seen/used before

- Simple implementation

| ROUNDNUMBER | NODEID |
|:---:|:---:|

- Each server stores MAXROUND, i.e., the largest round number it has seen so far

- To generate a new proposal number: increment MAXROUND and concatenate with NODEID

- Proposers must persist MAXROUND on disk: must not reuse proposal numbers after crash/restart

# Phase 1: PREPARE

- A proposer broadcasts a *prepare proposal (v, n)* with its own proposal value *v* and proposal number *n*

- An acceptor receives a *prepare proposal (v, n)*:

  - if it has never received a *prepare proposal*:

    - it *promises* to never accept *proposal numbers* lesser than $n$ ($n_{min} = n$)

    - it returns ($\varnothing$, *0*);

  - otherwise it checks its *promise*:

    - if $n > n_{min}$ it sends its last *accepted proposal ($v_{last}$, $n_{last}$)*; note that $n_{last} < n$.

    - otherwise it does nothing

# Phase 2: PROPOSE

- When a proposer receives a majority of responses, it broadcast to them a *propose proposal (v', n)*

  - If all acceptors returned (∅, 0), *v'* is its own proposal value *v*

  - Otherwise *v'* is the $v_{last}$ proposal value in the returned proposals with the greatest proposal number $n_{last}$

- If a majority of acceptors replies with ACK, the proposal is chosen

- Note that after a timeout, the proposer gives up and may send a new proposal

# PAXOS Algorithm

| Proposers | Acceptors |
|---|---|

0) $n_p = 0$  highest prepare number seen
$n_a = 0$  highest accepted proposal number
$v_a = \varnothing$  highest accepted proposal value

---

1) Choose new proposal number $n > n_p$

2) Broadcast PREPARE($v,n$) to all nodes

3) Handle PREPARE($v,n$):
   If ($n > n_p$) then
   $n_p = n$
   REPLY($v_a$, $n_a$)

4) If REPLY($v_a$, $n_a$) from majority:
   $v' = v_a$ with greatest $n_a$
   If $v' = \varnothing$ then
   $v' = v$

5) Broadcast PROPOSE($v',n$) to all nodes

6) Handle PROPOSE($v',n$):
   If ($n \geq n_p$) then
   $n_p = n$
   ($v_a$, $n_a$) = ($v',n$)
   ACCEPT($v'$, $n$)

7) If ACCEPT($v'$, $n$) from majority:
   Broadcast DECIDED($v'$) to all nodes

# PAXOS Algorithm

**Proposers**  |  **Acceptors**

0) $n_p = 0$  highest prepare number seen
$n_a = 0$  highest accepted proposal number
$v_a = \varnothing$  highest accepted proposal value

**These value must stably persist on disk**

1) Choose new proposal number $n > n_p$

2) Broadcast PREPARE($v,n$) to all nodes

3) Handle PREPARE($v,n$):
    If ($n > n_p$) then
        $n_p = n$
        REPLY($v_a$, $n_a$)

4) If REPLY($v_a$, $n_a$) from majority:
    $v' = v_a$ with greatest $n_a$
    If $v' = \varnothing$ then
        $v' = v$

5) Broadcast PROPOSE($v',n$) to all nodes

6) Handle PROPOSE($v',n$):
    If ($n \geq n_p$) then
        $n_p = n$
        ($v_a$, $n_a$) = ($v',n$)
        ACCEPT($v'$, $n$)

7) If ACCEPT($v'$, $n$) from majority:
    Broadcast DECIDED($v'$) to all nodes

# Learning a decision

- After a proposal is chosen, only the proposer knows about it!

- How do the other nodes get informed?

    1. The proposer could inform all nodes directly

        - If the proposer fails, the others are not informed (directly)...

    2. The acceptors could broadcast every time they accept a proposal

        - Much more fault-tolerant

        - Many accepted proposals may not be chosen...

        - Moreover, choosing a value costs $O(n^2)$ messages without failures!

    3. The proposer could inform some nodes directly

        - They will broadcast the decision to other nodes

# PAXOS is safe!

If a proposal *(v,n)* is **chosen**, then for every proposed proposal *(u,m)* for which $m > n$ it holds that $v = u$

- Assume that there is a proposed proposal *(u,m)* for which $m > n$ and $u \neq v$. Consider such a proposal with the smallest *m*.

- Consider the non-empty intersection *S* of the two majority sets of nodes that are acceptors for *(v,n)* and *(u,m)* proposals.

- Since proposal *(v,n)* has been accepted and $m > n$, nodes in *S* must have received PREPARE(*u,m*) after *(v,n)* has been accepted, thus returning REPLY(*v,n'*), with $n \leq n' < m$.

- As a consequence, the proposer of *(u,m)* should propose *(v, m)*, hence $u = v$, that is a contradiction.

# PAXOS is correct!

If a value is chosen, all acceptors choose this value

- Once a proposal *(v,n)* is chosen, each following proposal *(u,m)* has the same proposal value, i.e., *u = v*, according to the previous theorem.

- Since every following proposal has the same value *v*, every proposal that is accepted after *(v,n)* is chosen will have the same proposal value *v*.

- Since no other value than *v* can be accepted, no other values can be chosen.
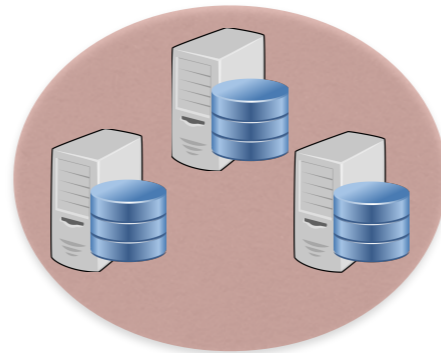
# PAXOS is great?

- PAXOS is a deterministic algorithm working for asynchronous systems and tolerating $f < n/2$ failures.

- Many optimizations exists (Multi-PAXOS, Disk-PAXOS, RAFT)

- **FLP Theorem** [1985]. No totally correct consensus algorithm exists (for the given system model).

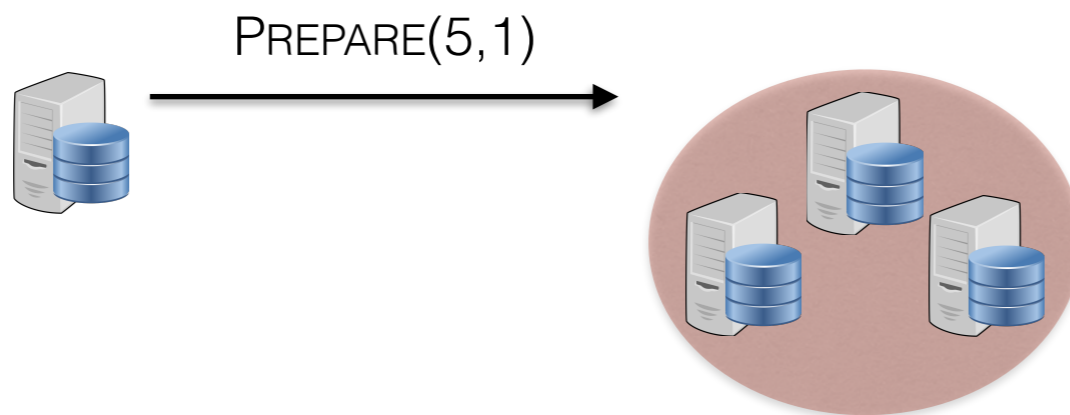- PAXOS disproves the FLP Theorem?

# No Liveness Guarantees

- PAXOS only guarantees that if a value is chosen, the other nodes can only choose the same value

- PAXOS does not guarantee that a value is chosen!
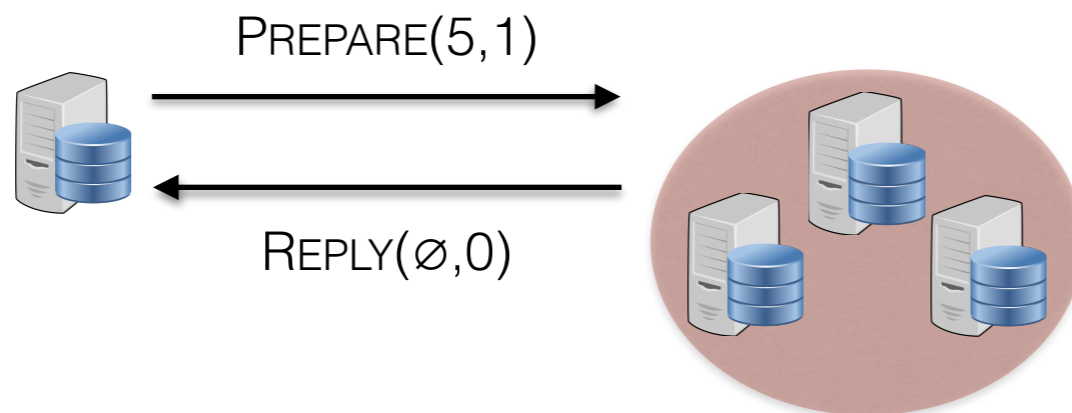
# No Liveness Guarantees

- PAXOS only guarantees that if a value is chosen, the other nodes can only choose the same value

- PAXOS does not guarantee that a value is chosen!

time

# No Liveness Guarantees
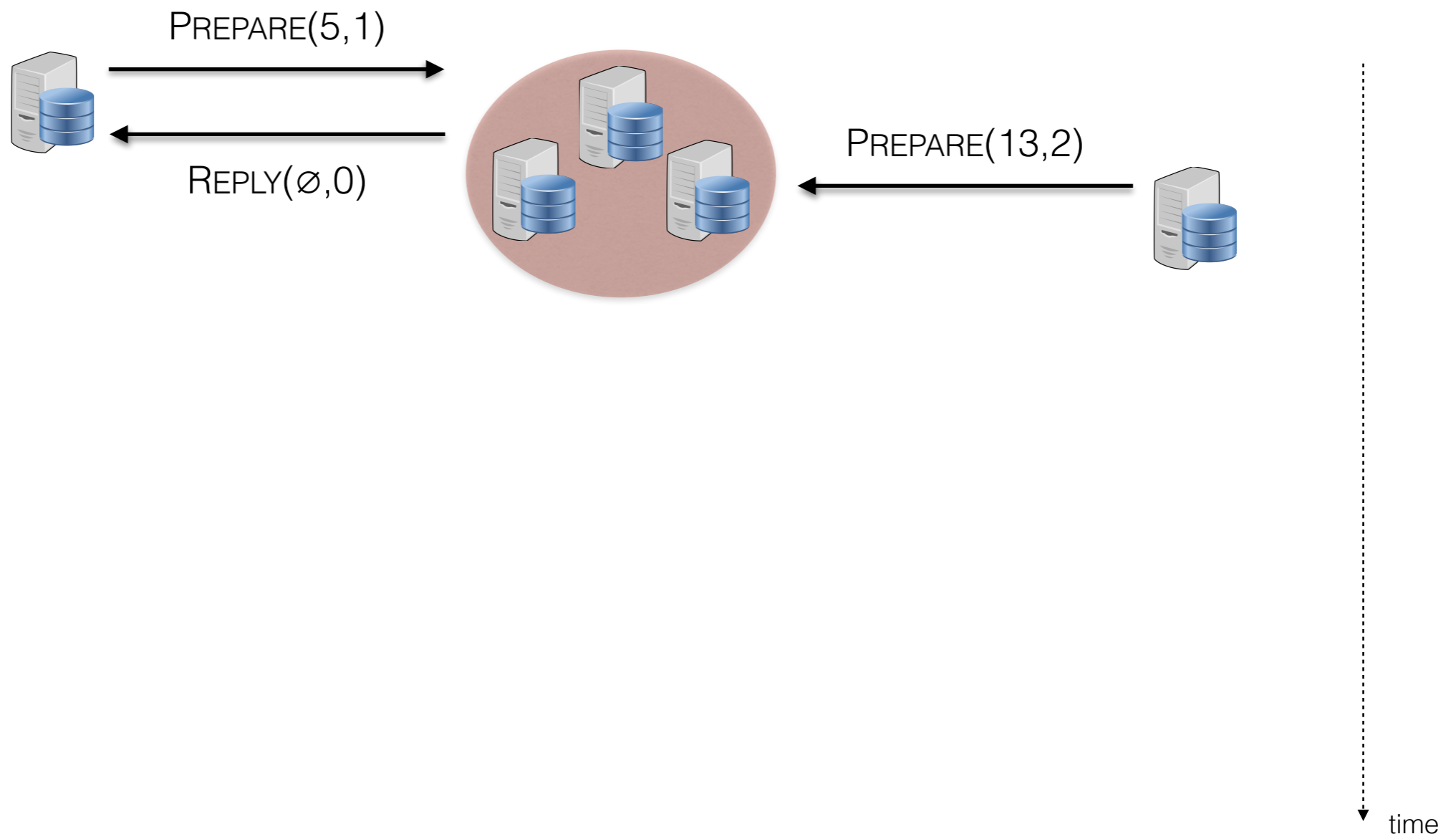
- PAXOS only guarantees that if a value is chosen, the other nodes can only choose the same value

- PAXOS does not guarantee that a value is chosen!
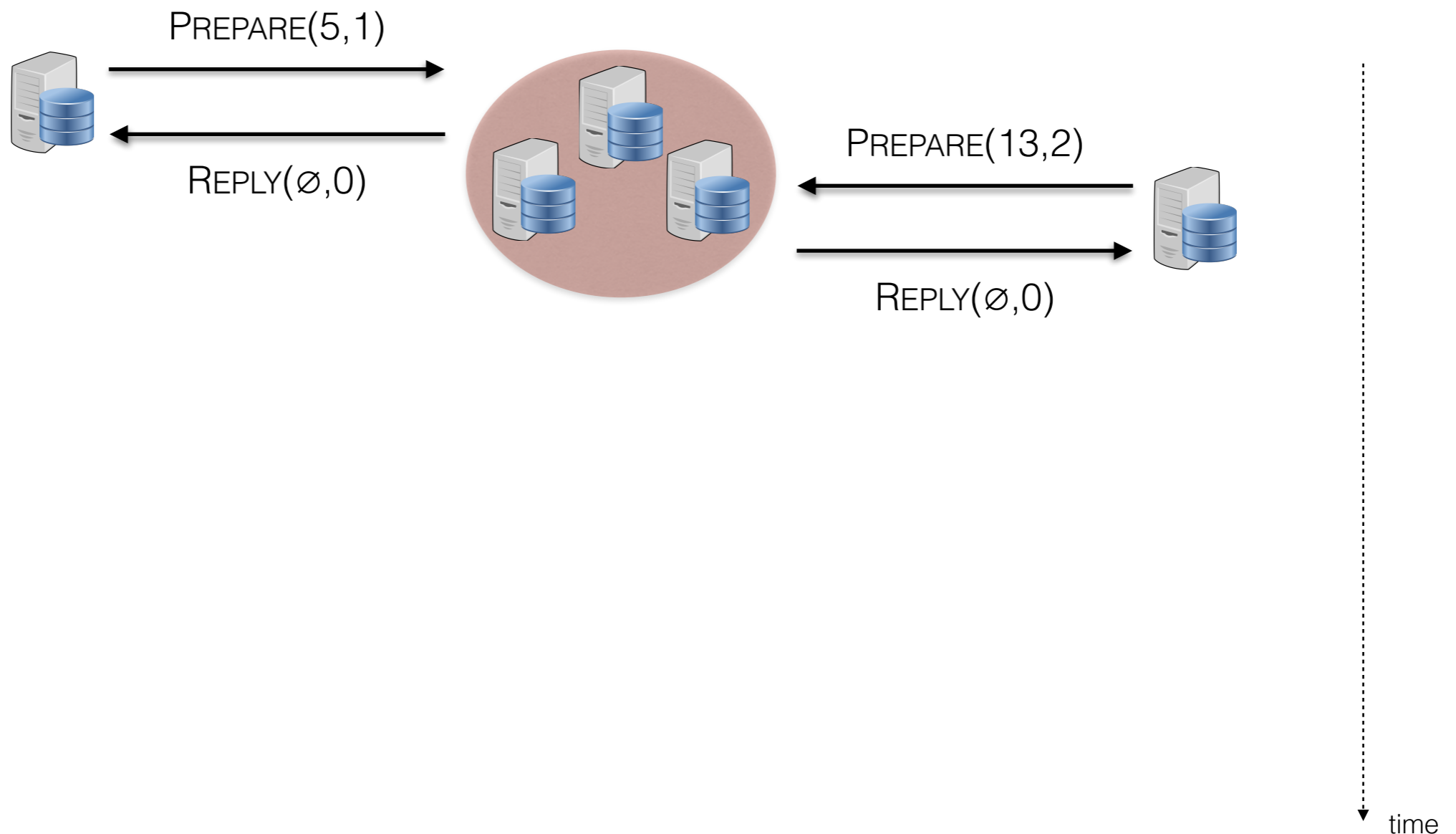
PREPARE(5,1)

time

# No Liveness Guarantees

- PAXOS only guarantees that if a value is chosen, the other nodes can only choose the same value

- PAXOS does not guarantee that a value is chosen!



PREPARE(5,1)

REPLY(∅,0)

time

# No Liveness Guarantees

- PAXOS only guarantees that if a value is chosen, the other nodes can only choose the same value

- PAXOS does not guarantee that a value is chosen!
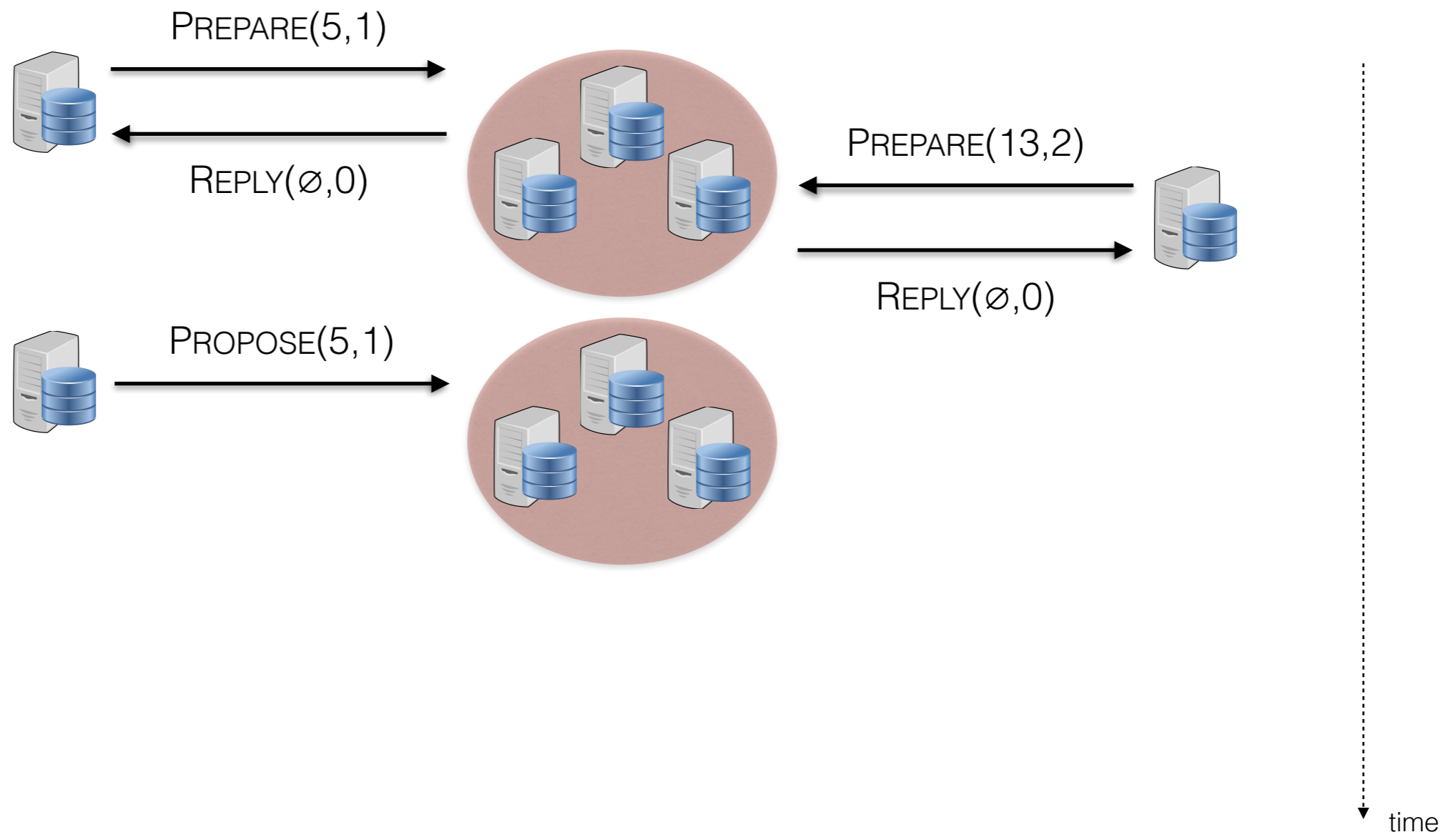
PREPARE(5,1)

REPLY(∅,0)

PREPARE(13,2)

time

# No Liveness Guarantees

- PAXOS only guarantees that if a value is chosen, the other nodes can only choose the same value

- PAXOS does not guarantee that a value is chosen!

PREPARE(5,1)

REPLY(∅,0)
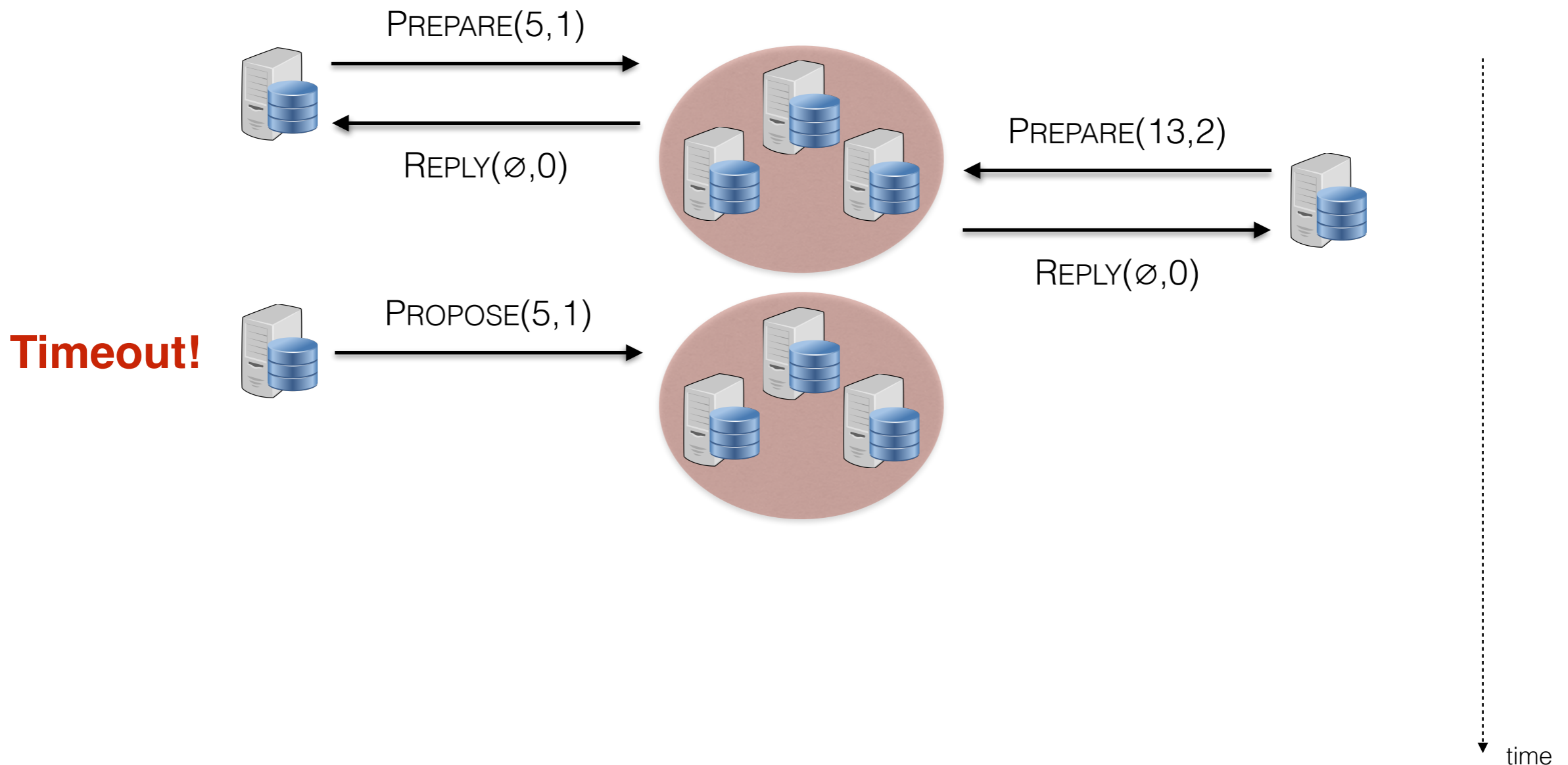
PREPARE(13,2)

REPLY(∅,0)

time

# No Liveness Guarantees

- PAXOS only guarantees that if a value is chosen, the other nodes can only choose the same value

- PAXOS does not guarantee that a value is chosen!

PREPARE(5,1)

REPLY(∅,0)
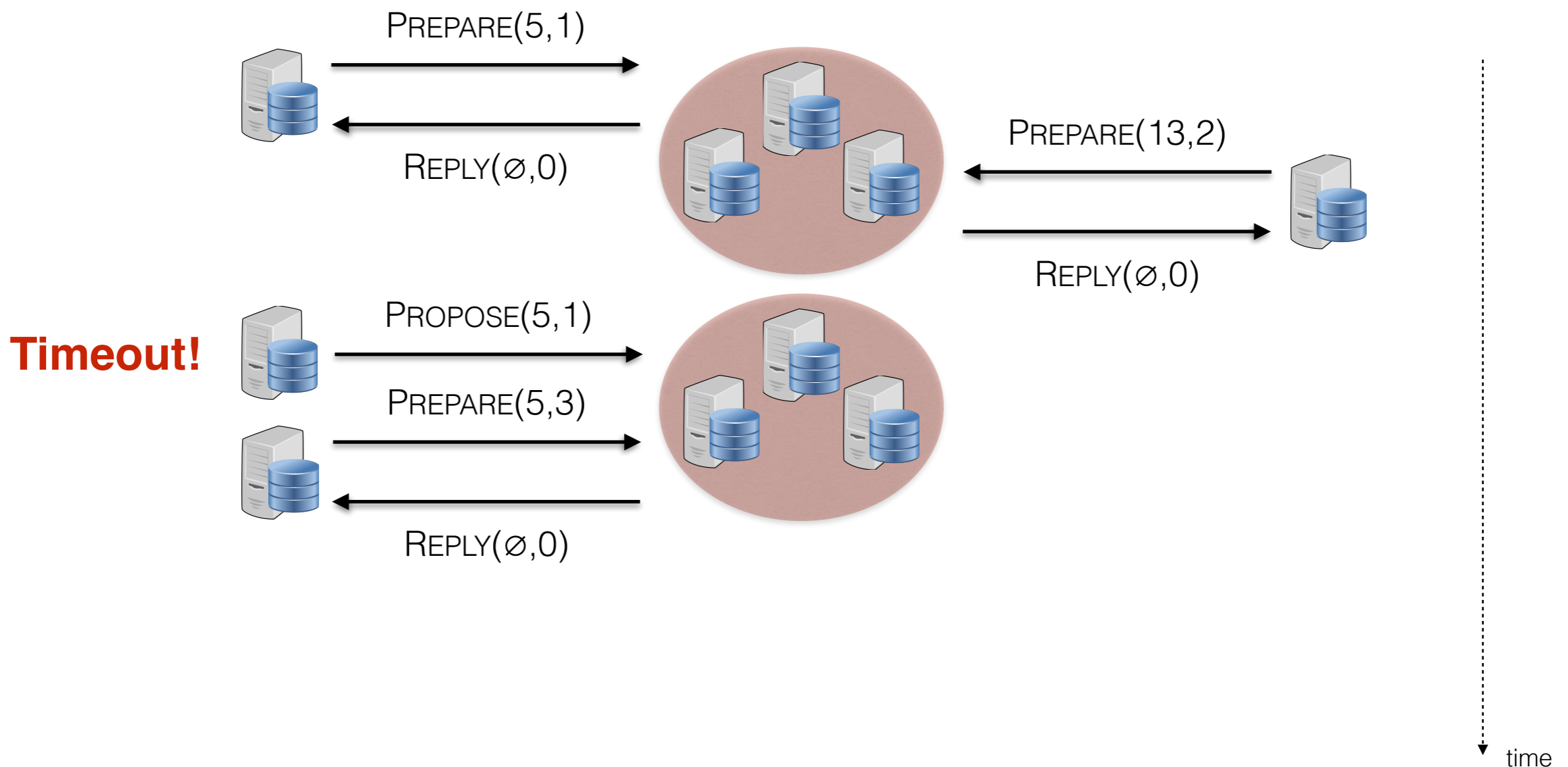
PREPARE(13,2)

REPLY(∅,0)

PROPOSE(5,1)

time

# No Liveness Guarantees

- PAXOS only guarantees that if a value is chosen, the other nodes can only choose the same value

- PAXOS does not guarantee that a value is chosen!

PREPARE(5,1)

REPLY($\varnothing$,0)

PREPARE(13,2)

REPLY($\varnothing$,0)

**Timeout!**
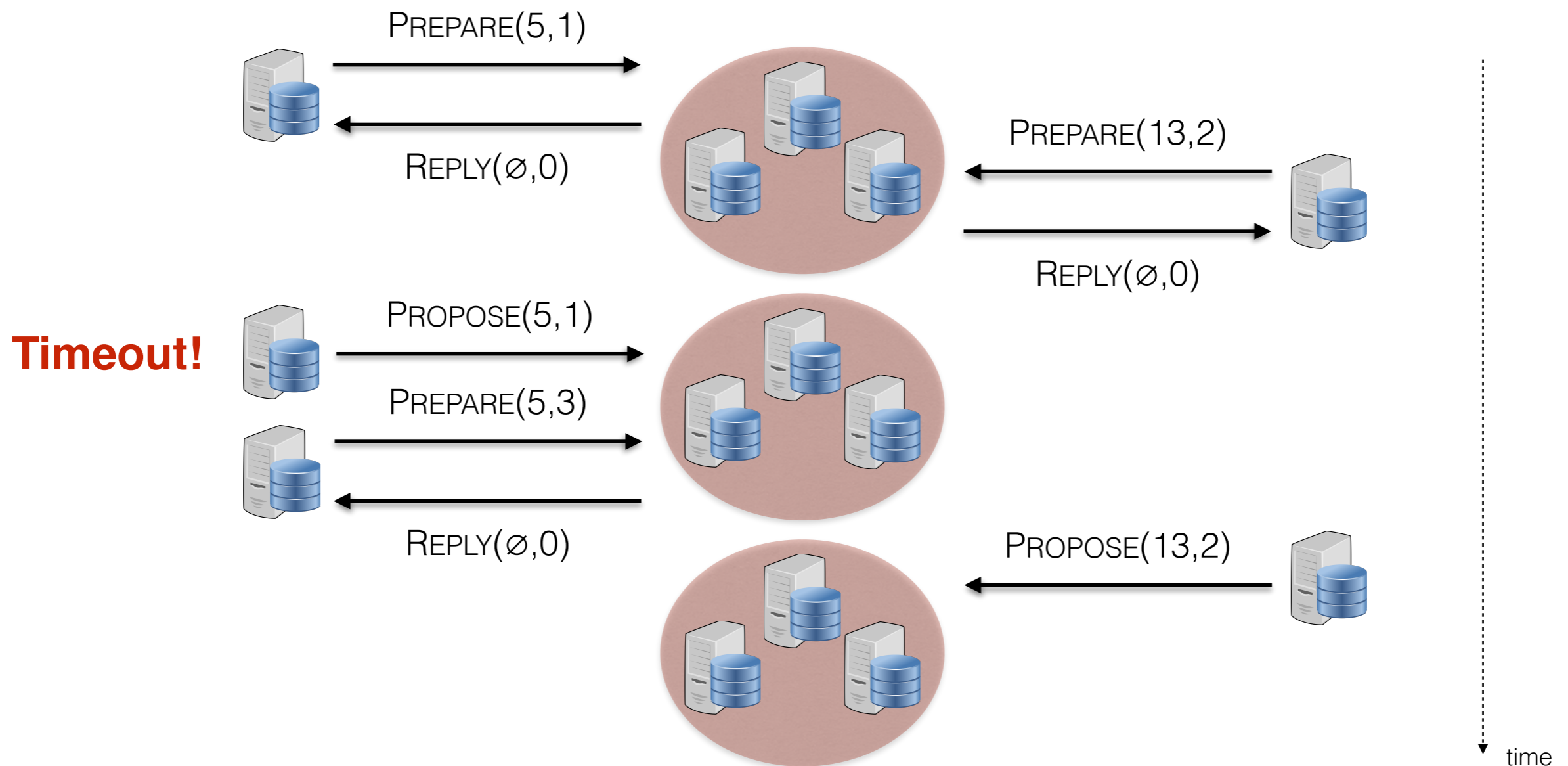
PROPOSE(5,1)

time

# No Liveness Guarantees

- PAXOS only guarantees that if a value is chosen, the other nodes can only choose the same value

- PAXOS does not guarantee that a value is chosen!

PREPARE(5,1)

REPLY(∅,0)

PREPARE(13,2)

REPLY(∅,0)

**Timeout!**

PROPOSE(5,1)

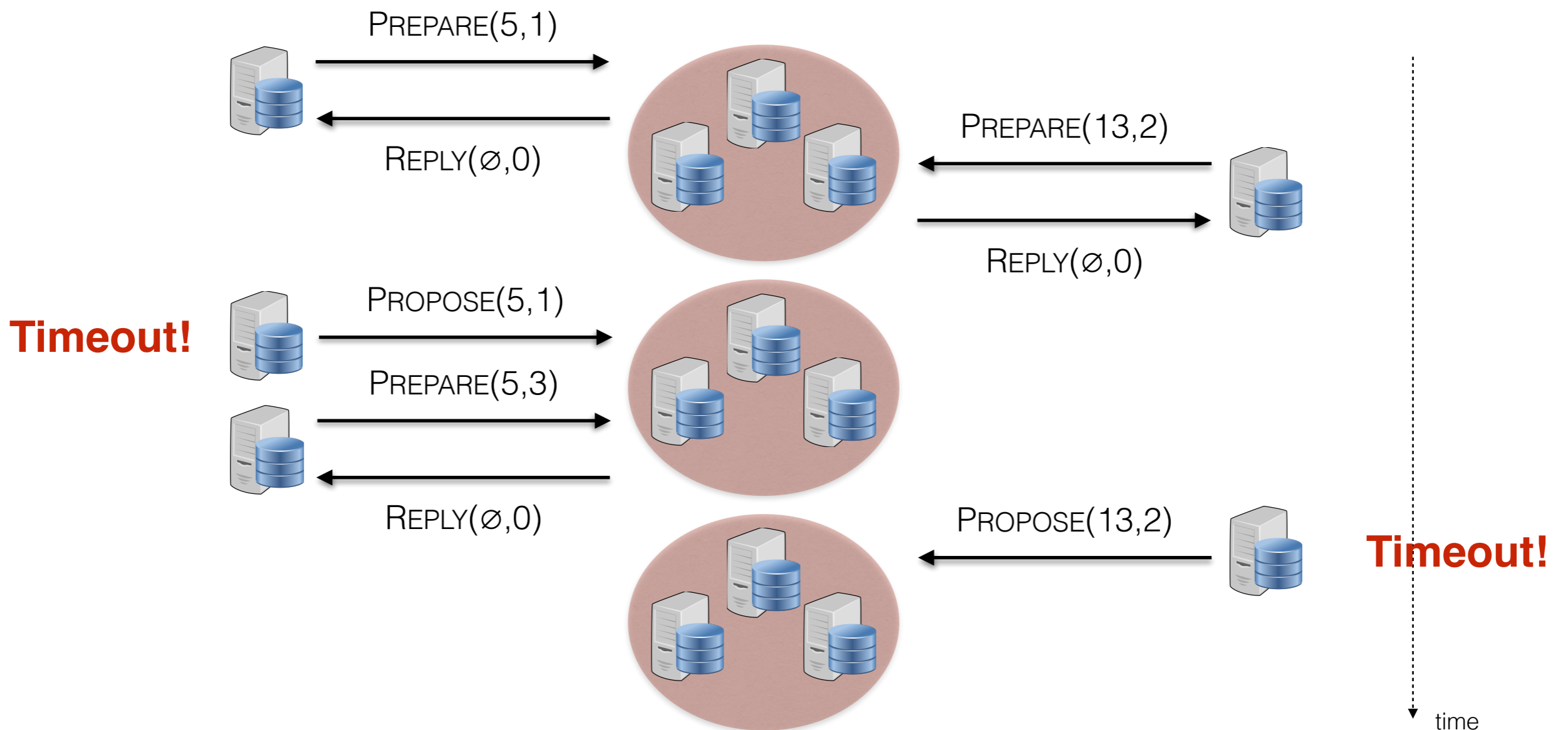PREPARE(5,3)

REPLY(∅,0)

time

# No Liveness Guarantees

- PAXOS only guarantees that if a value is chosen, the other nodes can only choose the same value

- PAXOS does not guarantee that a value is chosen!
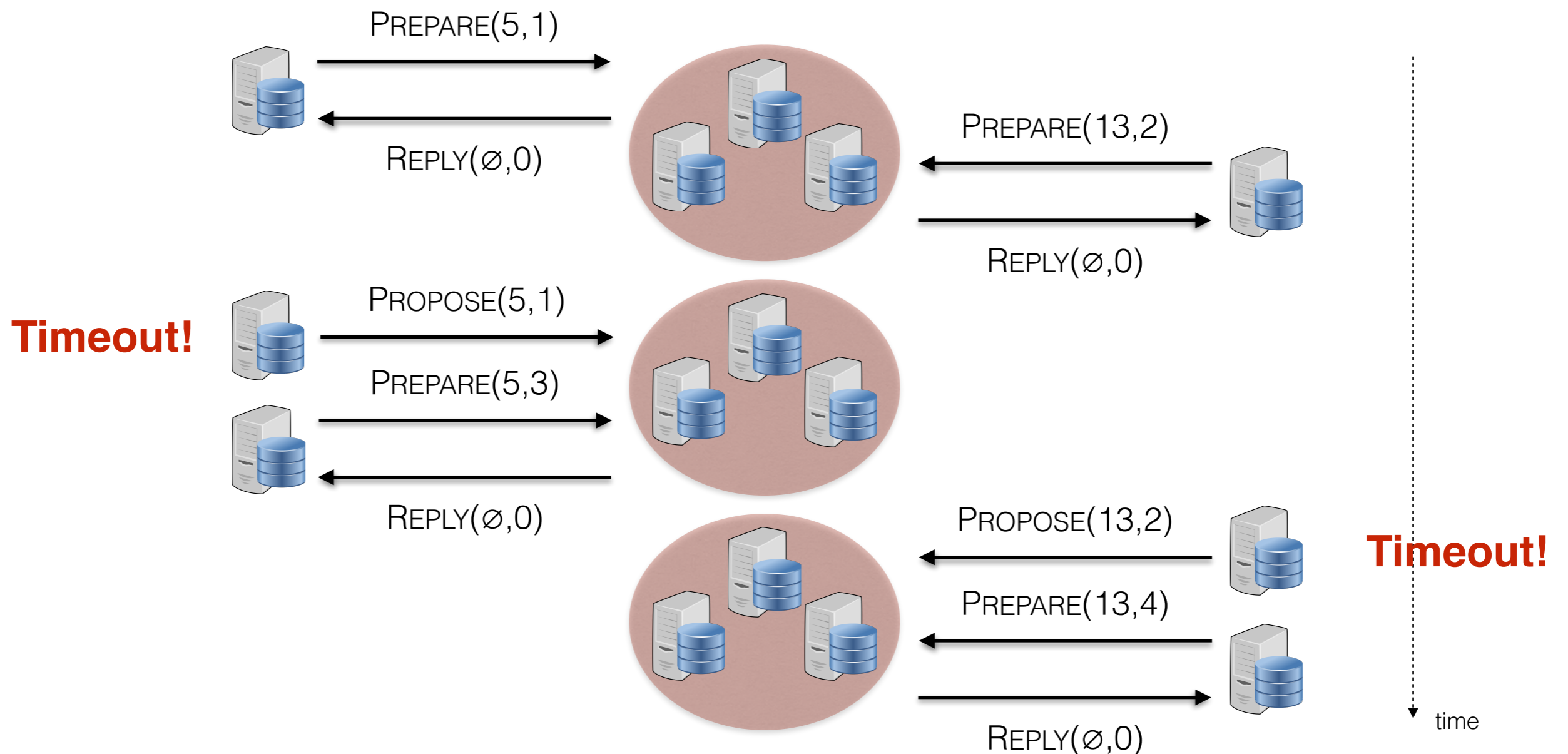
# No Liveness Guarantees

- PAXOS only guarantees that if a value is chosen, the other nodes can only choose the same value

- PAXOS does not guarantee that a value is chosen!

# No Liveness Guarantees

- PAXOS only guarantees that if a value is chosen, the other nodes can only choose the same value

- PAXOS does not guarantee that a value is chosen!

# Correctness vs. Termination

- In asynchronous systems, we cannot guarantee termination and correctness at the same time

- PAXOS is correct, so termination is not guaranteed

- PAXOS cannot guarantee that a consensus is reached in a finite number of steps

- In practice, PAXOS can be optimized to reduce probability of no termination

- For example, the acceptors could send NAK if they do not accept a prepare message or a proposal (this optimization increases the message complexity)

- PAXOS is used in Apache's Zookeeper and Google's Chubby