

CAP

- Presented as Brewer's Conjecture in 2000
- Formalized and proved in **Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services**, by Lynch and Gilbert (2002)
- **Consistency, availability** and **partition-tolerance** cannot be achieved all at the same time in a distributed system.
- Simply, in an asynchronous network that performs as expected, where messages may be lost (partition-tolerance), it is impossible to implement a service providing correct data (consistency) and eventually responding to every request (availability) under every pattern of message loss.
- Slides from http://cs-wwwarchiv.cs.unibas.ch/lehre/hs10/cs341/_Downloads/Workshop/Talks/2010-HS-DIS-I_Giangreco-CAP_Theorem-Talk.pdf

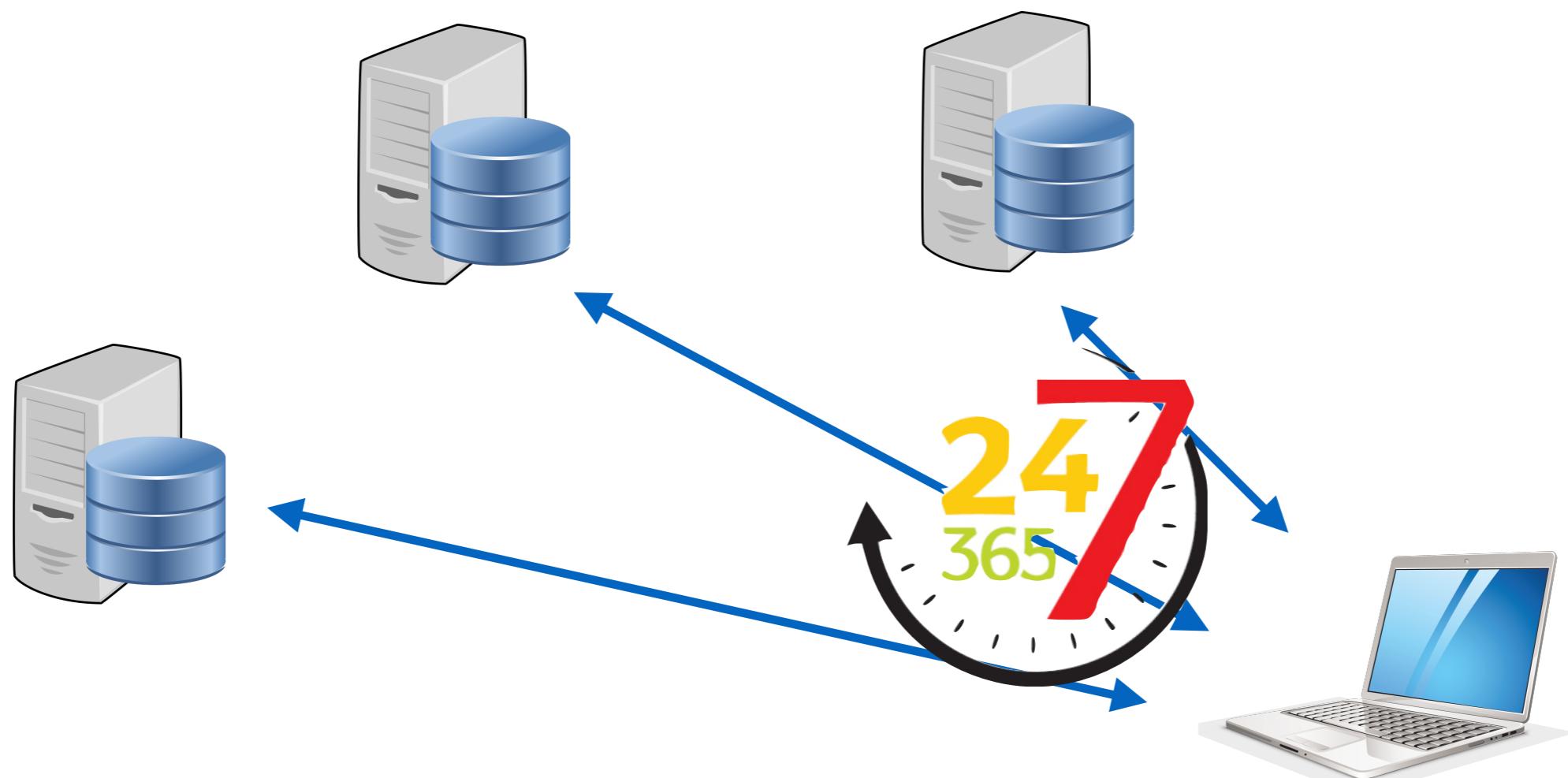
Consistency

- All the nodes in the system see the same state of the data
- Formally, we speak of *atomic* or *linearizable consistency*
- There exists a sequential order on all operations which is consistent with the order of invocations and responses, such that each operation looks as if it were completed at a single instant.



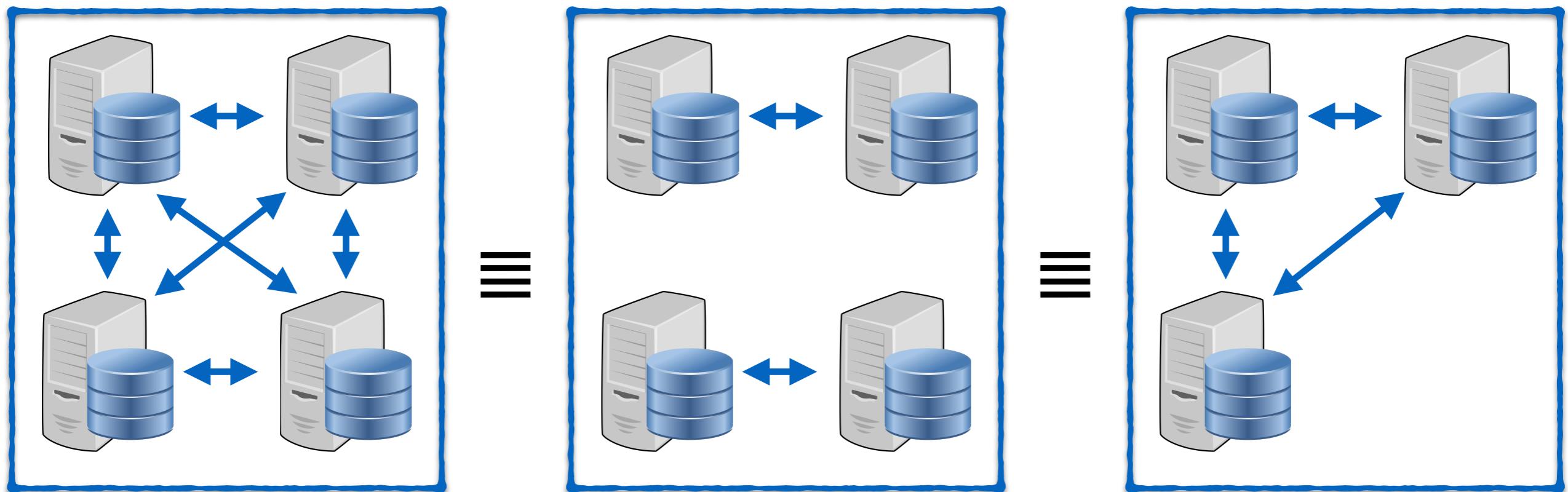
Availability

- Every request received by a non-failing node should be processed and must result in a response



Partition Tolerance

- If some nodes crash and/or some communications fail, system still performs as expected



CAP Theorem 1

It is impossible in the **asynchronous network model** to implement a read/write data object that guarantees the following properties:

- Availability
- Atomic consistency

in all fair executions (including those in which messages are lost).

Asynchronous, i. e. there is no clock, nodes make decisions based only on the messages received and local computation.

CAP Theorem 2

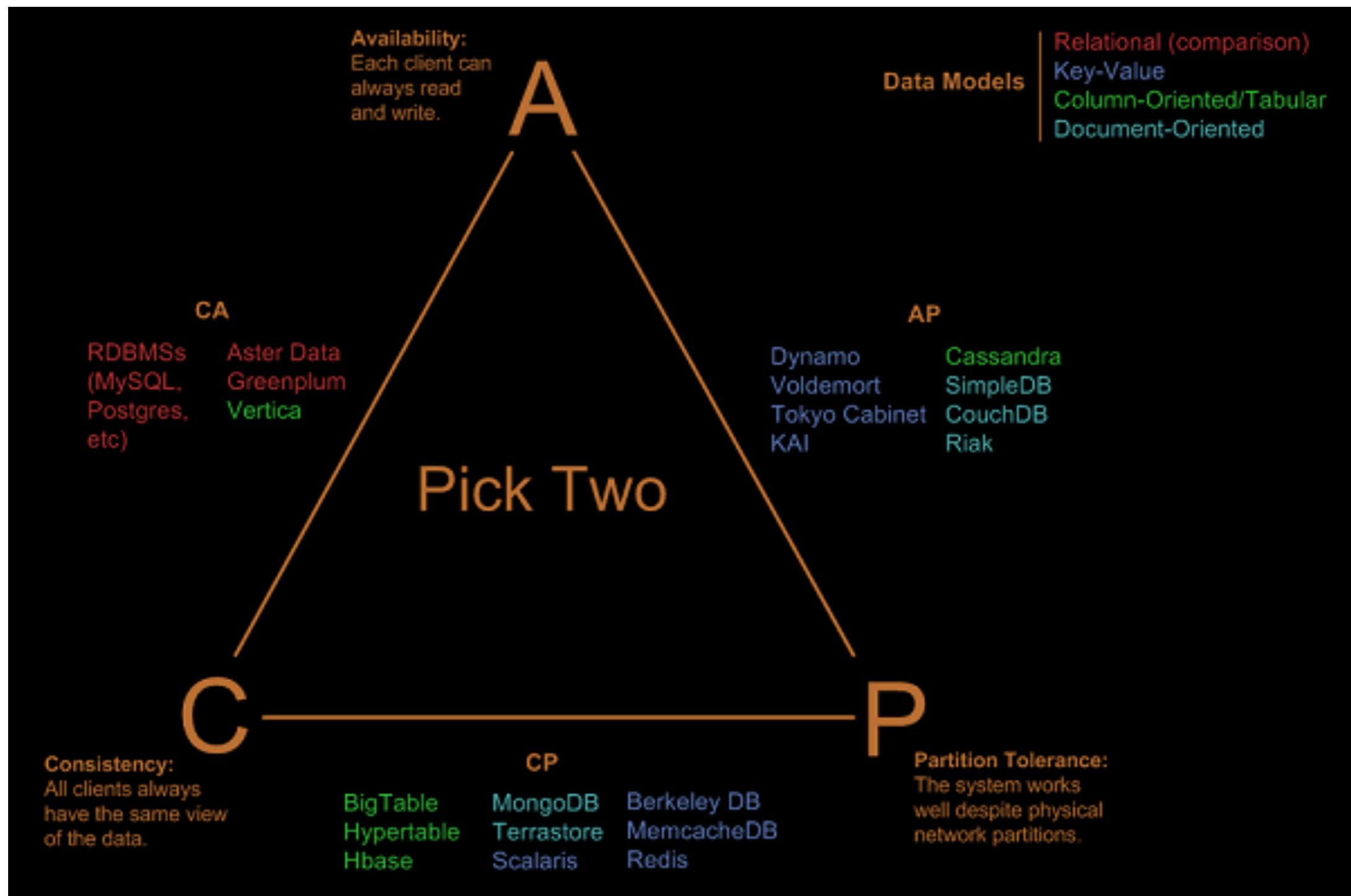
It is impossible in the **partially synchronous network model** to implement a read/write data object that guarantees the following properties:

- Availability
- Atomic consistency

in all fair executions (including those in which messages are lost).

Partially synchronous, i. e. every node has a clock, and all clocks increase at the same rate. However, they are not synchronized.

Consequences of CAP



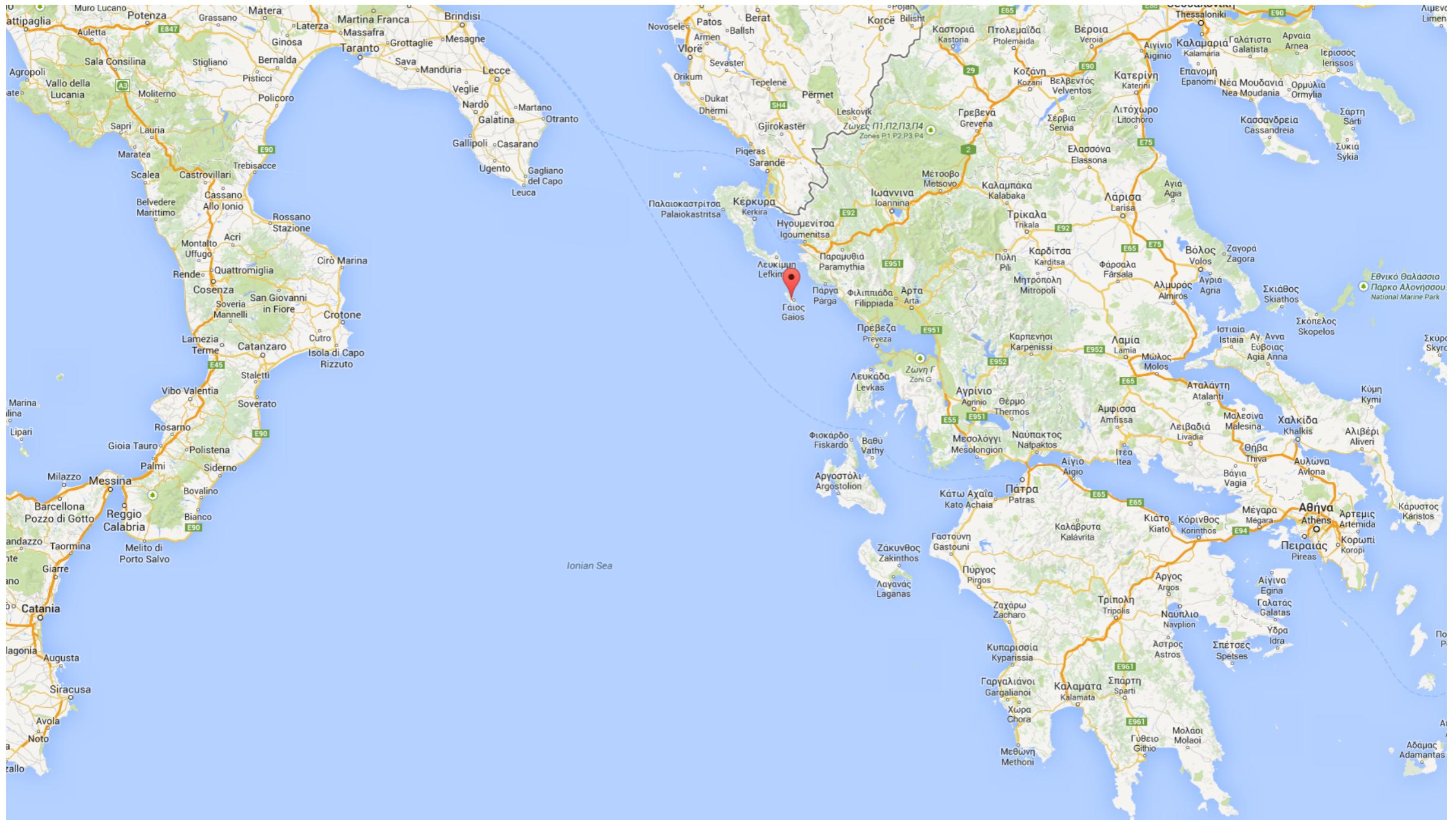
Consequences of CAP

- When partitions are rare (e.g., parallel systems), CAP should allow perfect C and A most of the time
- In distributed systems it is not possible to avoid network partitions.
- There is not a need to choose between either C or A, instead, it is more an act of balancing between the two properties.

Practical Consequences of CAP

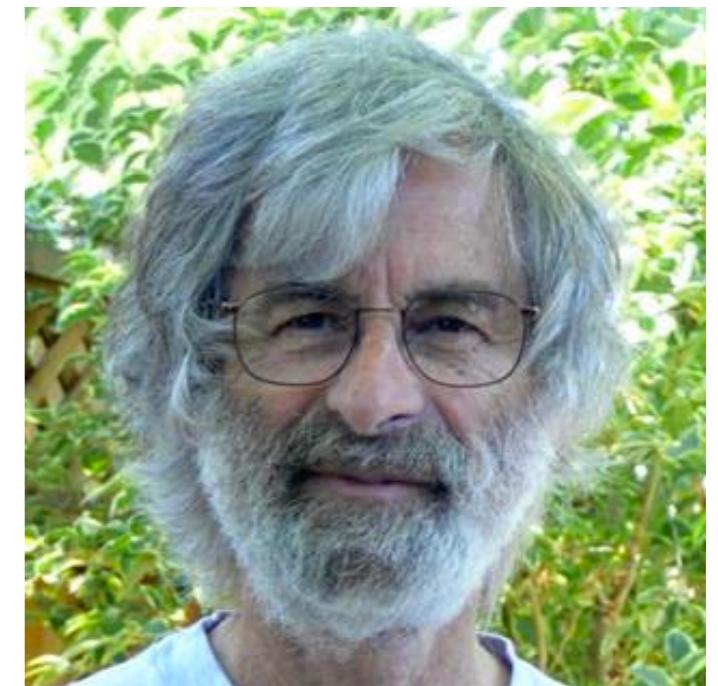
- Many system designs used in early distributed relational database systems did not take into account partition tolerance (e.g. they were CA designs).
- There is a tension between strong consistency and high availability during network partitions. A distributed system consisting of independent nodes connected by an unpredictable network cannot behave in a way that is indistinguishable from a non-distributed system.
- There is a tension between strong consistency and performance in normal operation. Strong consistency requires that nodes communicate and agree on every operation. This results in high latency during normal operation.
- If we do not want to give up availability during a network partition, then we need to explore whether consistency models other than strong consistency are workable for our purposes.

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

- Asynchronous communications
- Non-byzantine failures:
 - Nodes crash
 - Messages can be lost, duplicated arbitrarily late
 - No corrupted messages
- Fail-recover: crashed nodes may recover later

Roles

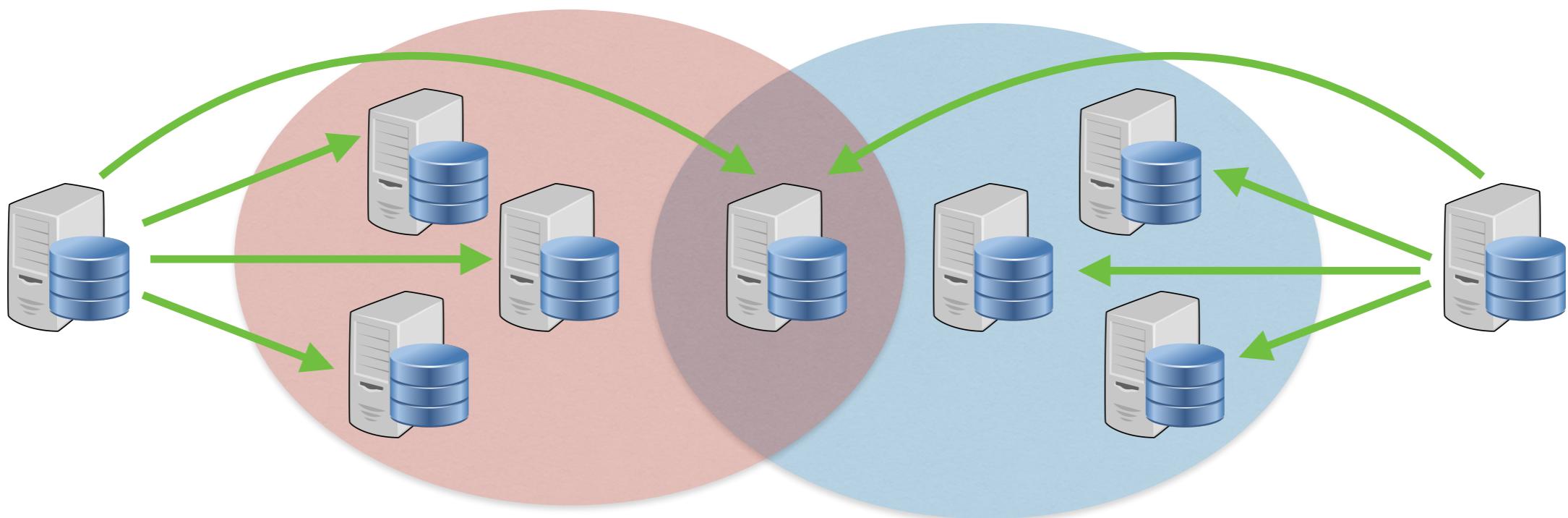
- Each nodes has one or more of 3 roles
- Proposer: a node that can propose a certain value for acceptance
- Acceptor: a node that receives proposals from proposers and that can either accept or reject a proposal
- Learner: a node not involved in the decision process that wants to know the final result 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 chooses a value, no other majority can choose 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

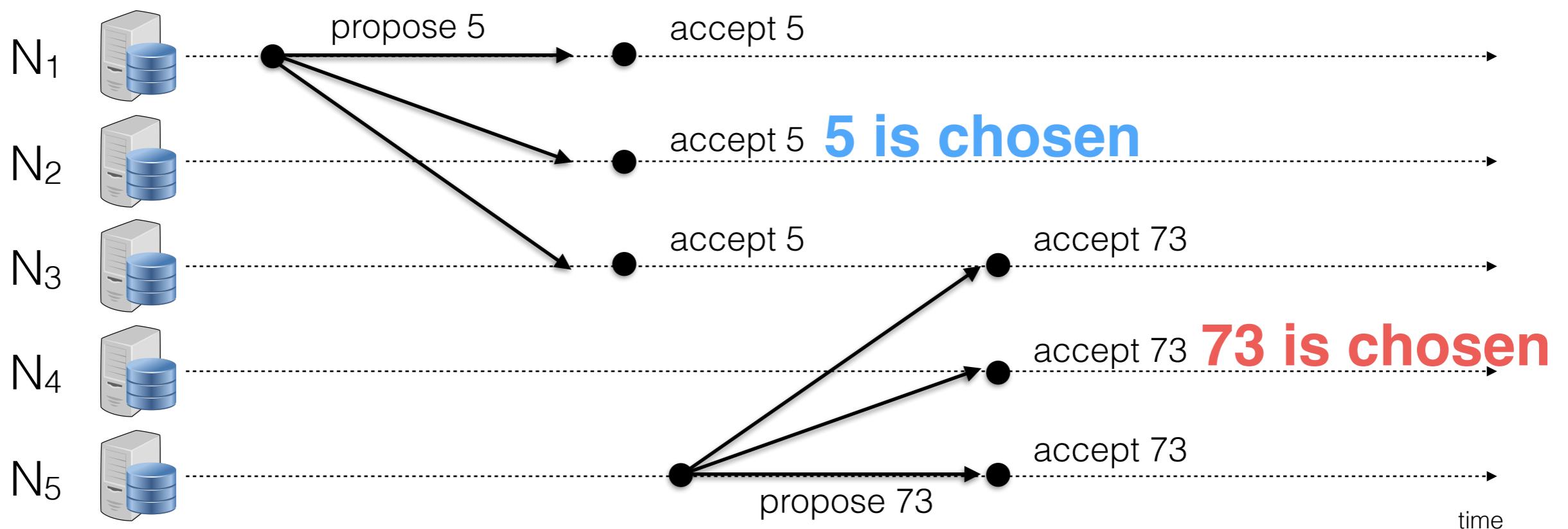


**An acceptor must accept the first proposed value it receives
(and can accept more than one proposed value)**

An acceptor can **accept** any number of proposed values, but an accepted proposed value may not necessarily be **chosen**.

Conflicting choices

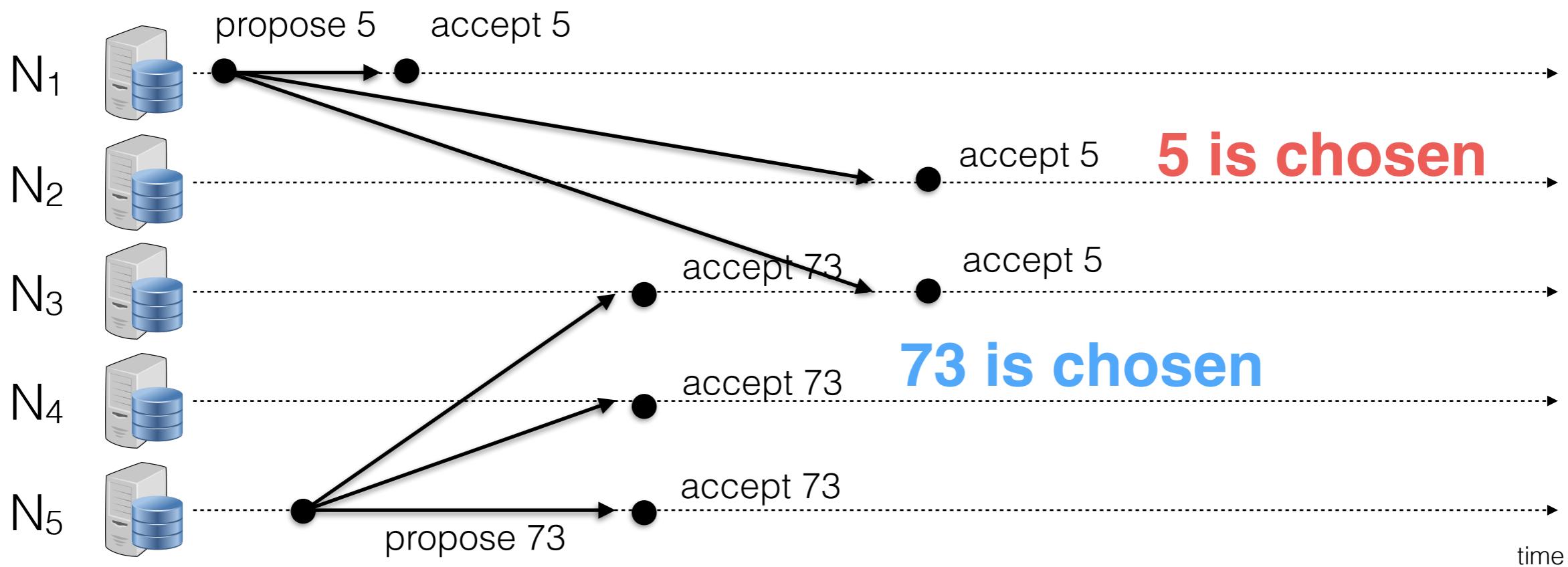
Should acceptors accept every proposed value? No



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

We need a two-phase protocol

Conflicting choices

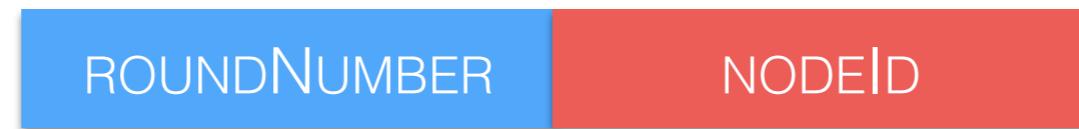


Once a value has been **chosen**, older proposed values must be ignored

We need 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 strictly-increasing proposal number
- For a given proposer, a new proposal number must be greater than anything it has seen/ used before
- Simple implementation



- 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 $(\emptyset, 0)$;
 - 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 $(\emptyset, 0)$, 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
 $v_a = \emptyset$ 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

4) If REPLY(v_a, n_a) from majority:

$v' = v_a$ with greatest n_a

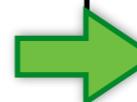
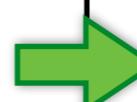
If $v' = \emptyset$ then

$v' = v$

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

7) If ACCEPT(v', n) from majority:

Broadcast DECIDED(v') to all nodes



3) Handle PREPARE(v, n):

If ($n > n_p$) then

$n_p = n$

REPLY(v_a, n_a)

6) Handle PROPOSE(v', n):

If ($n \geq n_p$) then

$n_p = n$

$(v_a, n_a) = (v', n)$

ACCEPT(v', n)

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 $\text{PREPARE}(u, m)$ after (v, n) has been accepted, thus returning $\text{REPLY}(v, n')$, 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

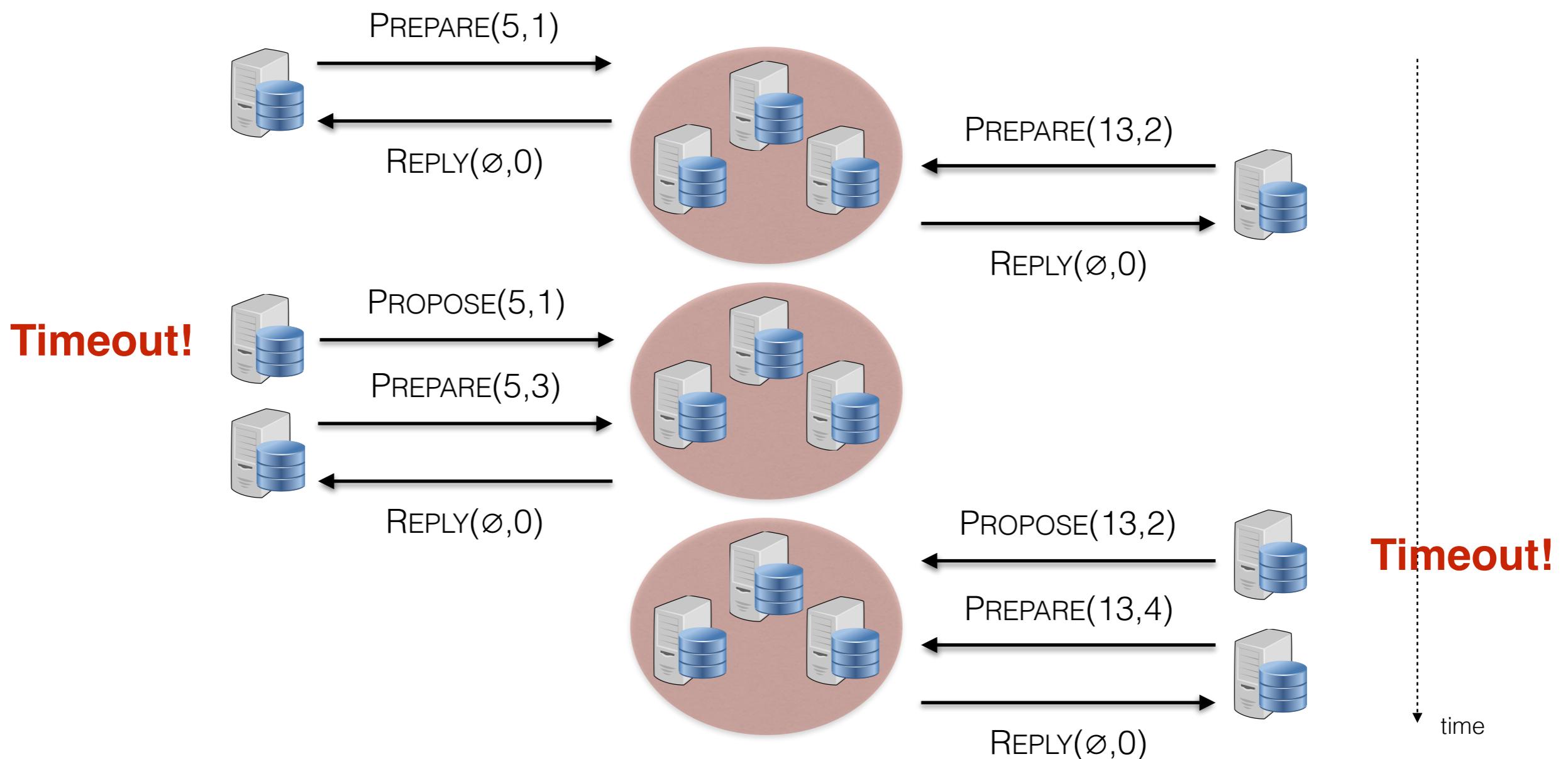
- Once a proposal (v,n) is chosen, each following proposal (u,m) has the same proposal value, i.e., $u = v$, according to the previous theorem.
- Since every following proposal has the same value v , every proposal that is accepted after (v,n) is chosen will have the same proposal value v .
- Since no 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!



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