

# EECS 591

# DISTRIBUTED SYSTEMS

Manos Kapritsos  
Fall 2021

## **Abstract**

The Paxos algorithm, when presented in plain English, is very simple.

# THREE TYPES OF PROCESSES

**Proposers** A proposer is a process that has a value to propose

**Acceptors  
( $2f+1$ )** Acceptors are the processes that ultimately choose which proposed value will be decided

**Learners** A learner only cares about learning which value was decided

# HOW THE GAME IS PLAYED

- **Election:** Proposers first try to get a majority of acceptors to follow them.
- **Legislating:** After acquiring a majority, a proposer can *try* to enforce her value, by getting acceptors to accept it, **but...**
- **Playing nice:** If an elected proposer finds that some previous value has been proposed, she proposes that value instead.
- **Winning the game:** once a majority of acceptors have accepted a value, the value is **chosen/decided**

# HOW IT IS SUPPOSED TO WORK

Proposer

Greetings, peasants! I am your fearless leader! Grant me your blessing!

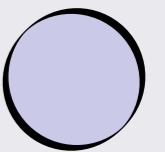
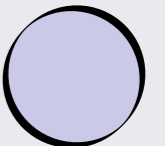
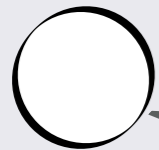
My first decree is:

The value should be 12

Acceptors

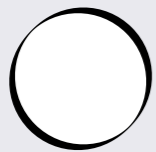
We are with you, oh wise leader!

Sounds good to me!

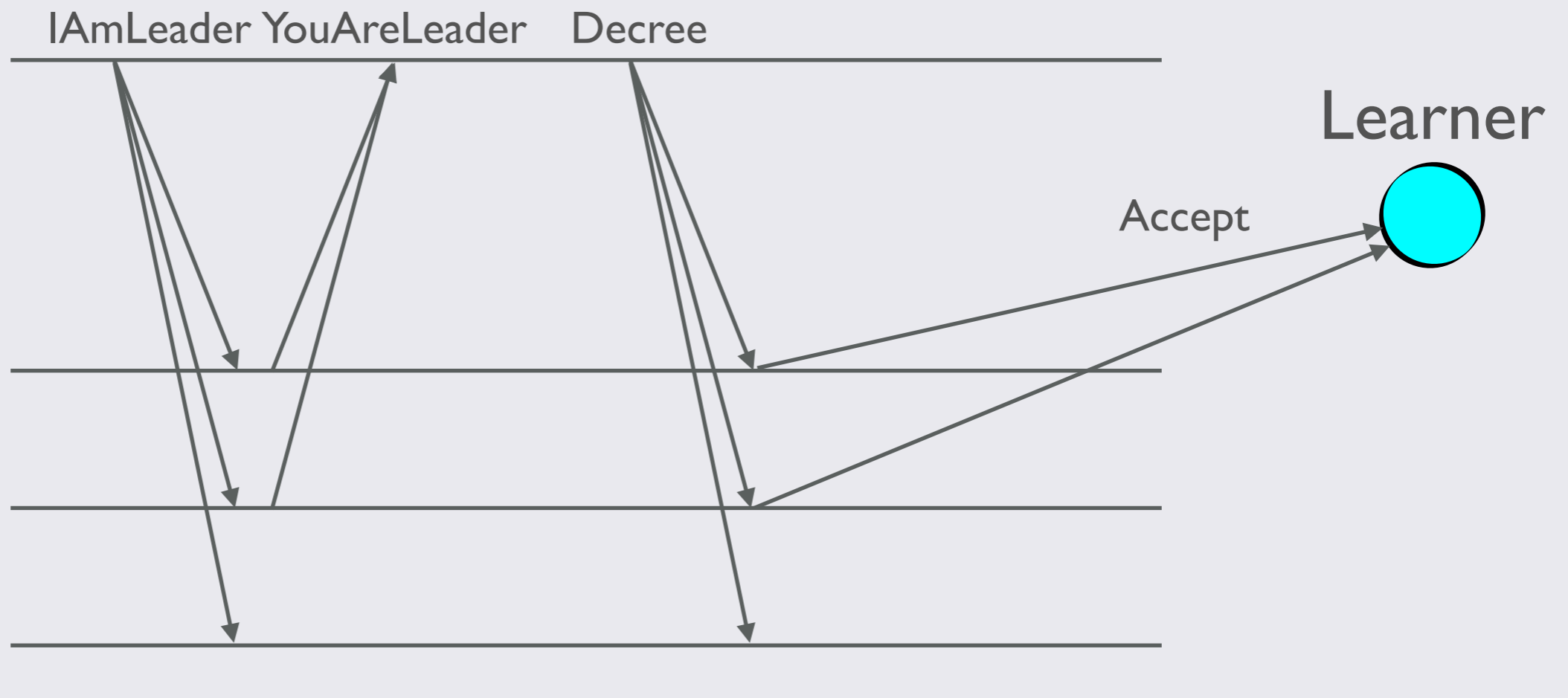
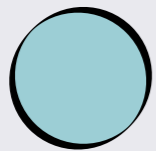
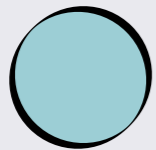
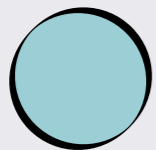


# HOW IT IS SUPPOSED TO WORK

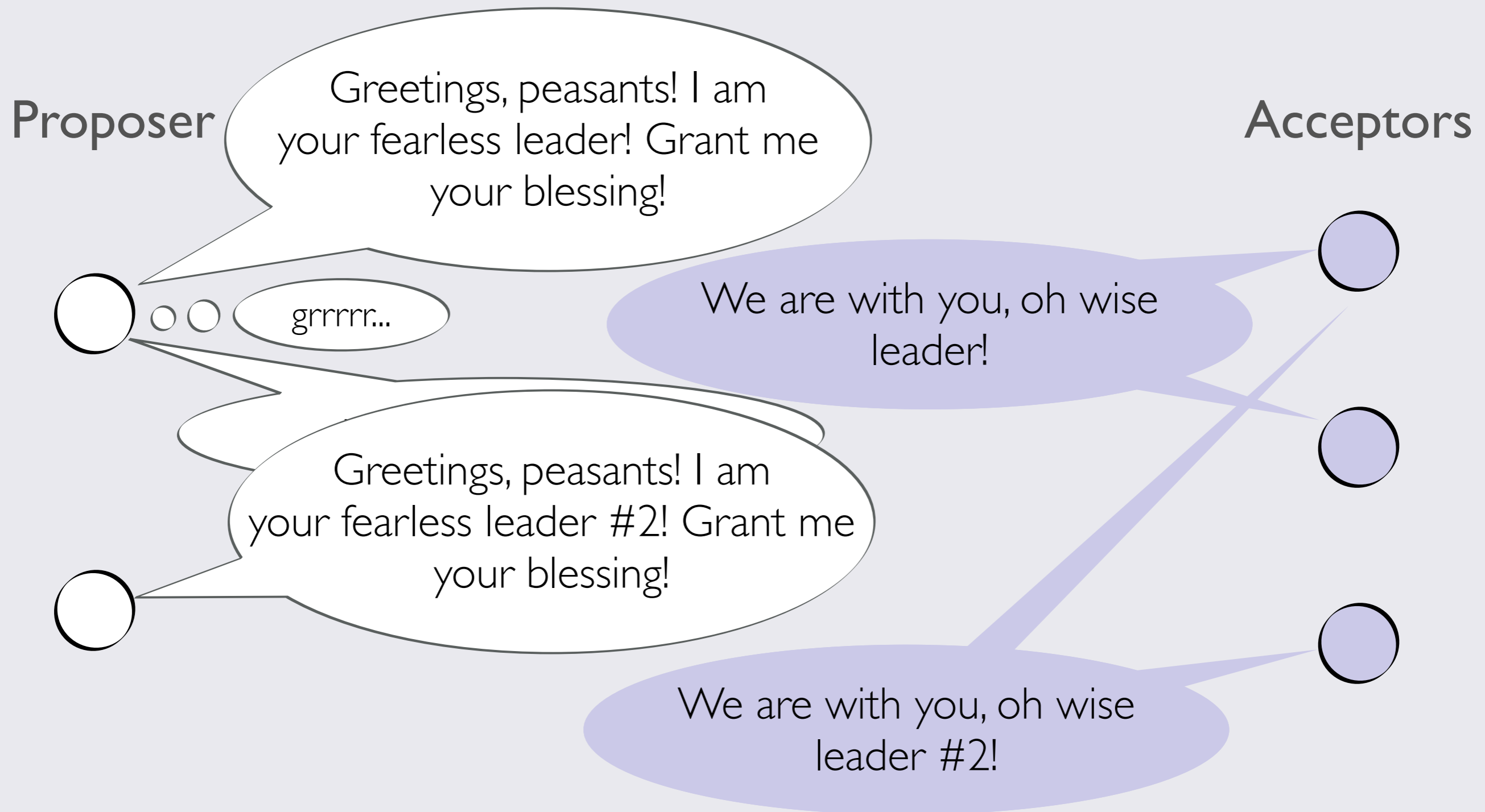
Proposer



Acceptors

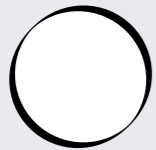


# DEALING WITH MULTIPLE PROPOSERS

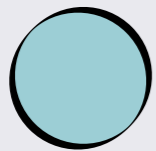
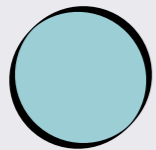
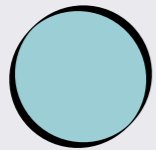


# DEALING WITH MULTIPLE PROPOSERS

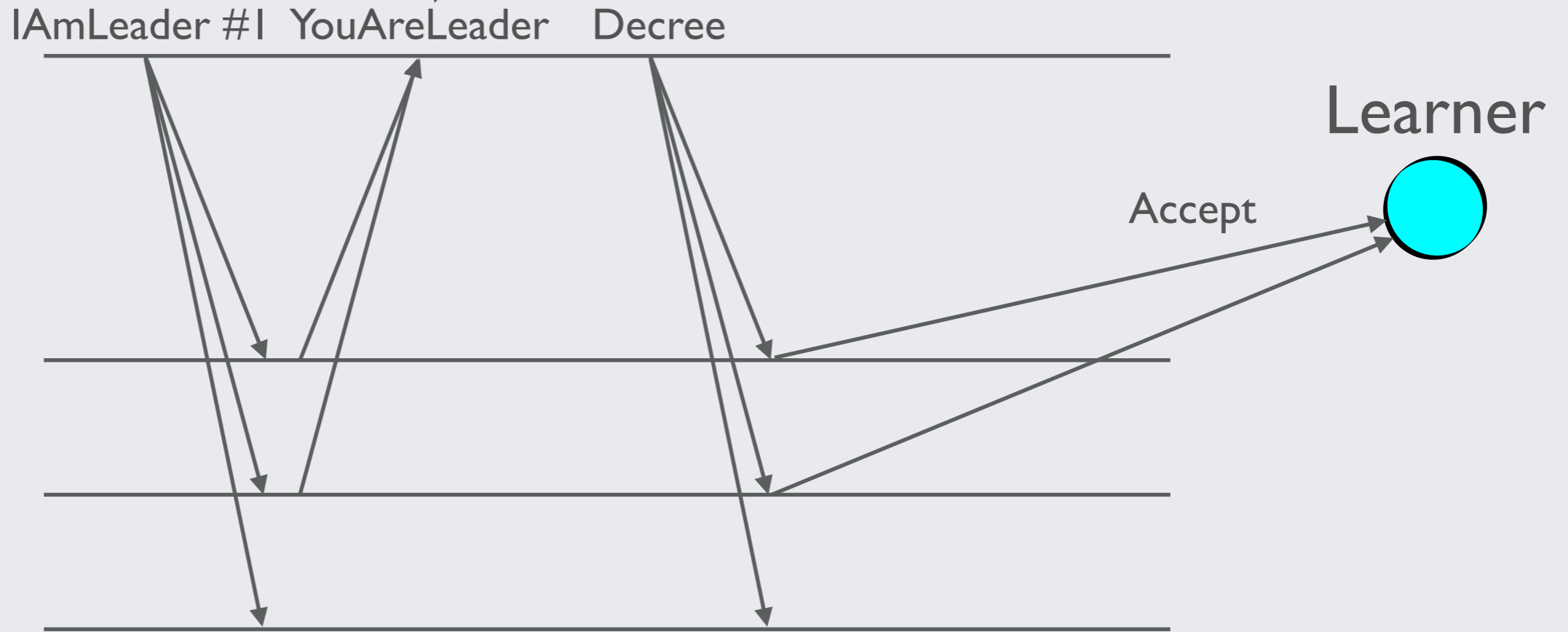
Proposer



Acceptors



- I swear I won't follow an earlier leader!
- And, btw, here is my current accepted value (if any) by leader x.

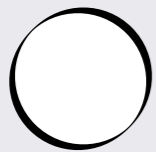




# DEALING WITH MULTIPLE PROPOSERS

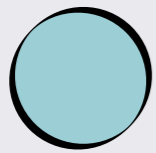
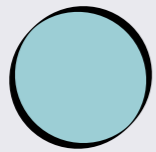
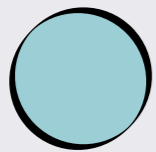
- I swear I won't follow an earlier leader!
- And, btw, here is my current accepted value (if any) by leader x.

Proposer #1



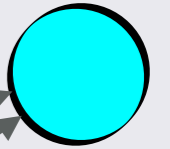
IAmLeader #1 YouAreLeader Decree

Acceptors

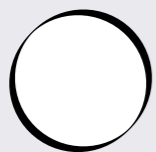


Learner

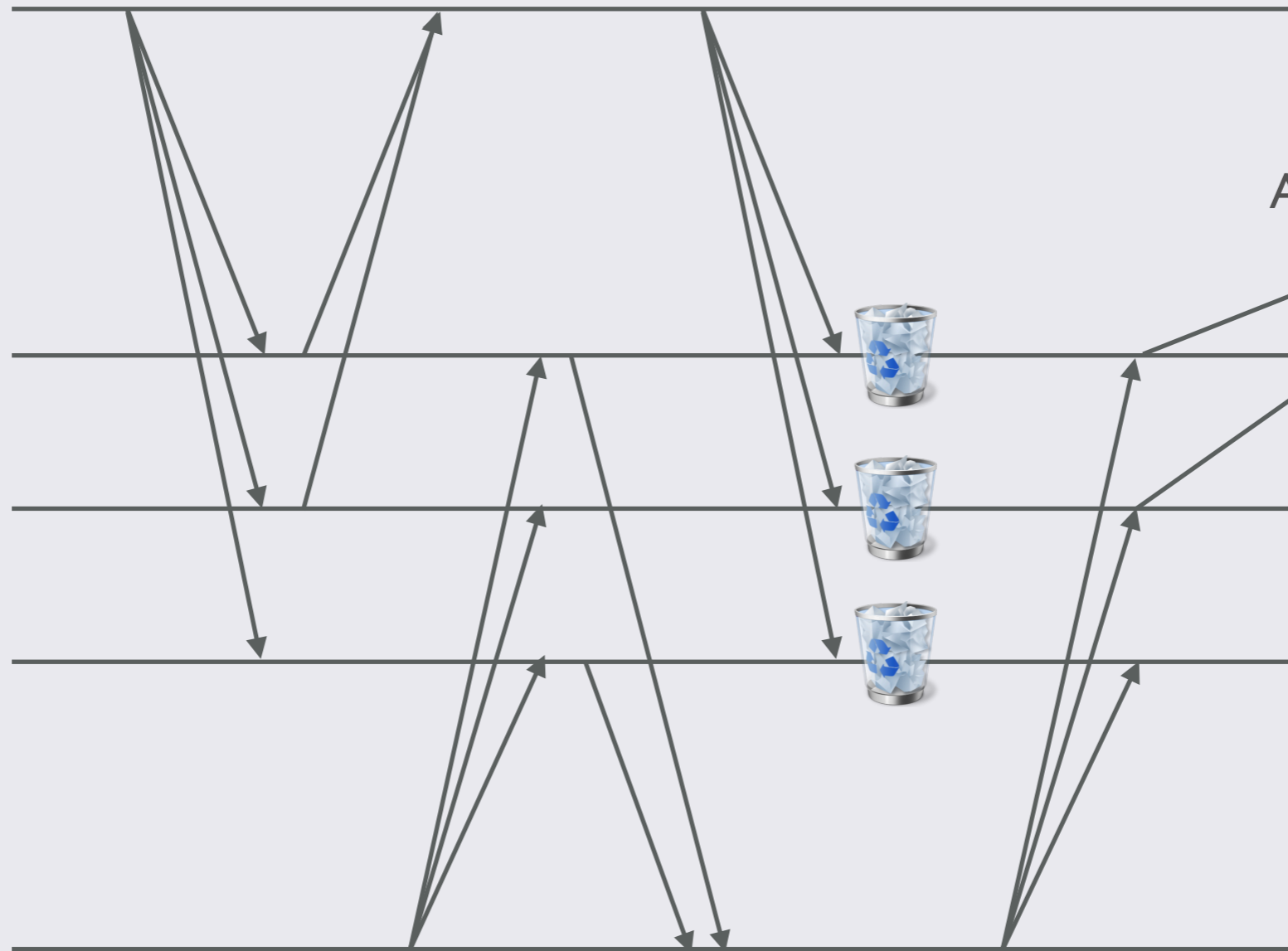
Accept



Proposer #2



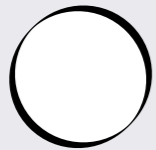
IAmLeader #2 YouAreLeader Decree



# DEALING WITH MULTIPLE PROPOSERS

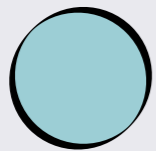
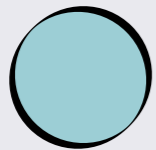
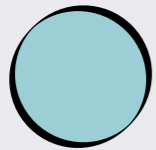
- I swear I won't follow an earlier leader!
- And, btw, here is my current accepted value (if any) by leader x.

Proposer #1

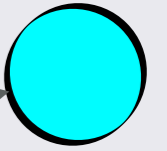


IAmLeader #1 YouAreLeader Decree

Acceptors

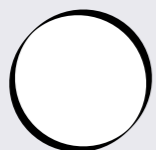


Learner

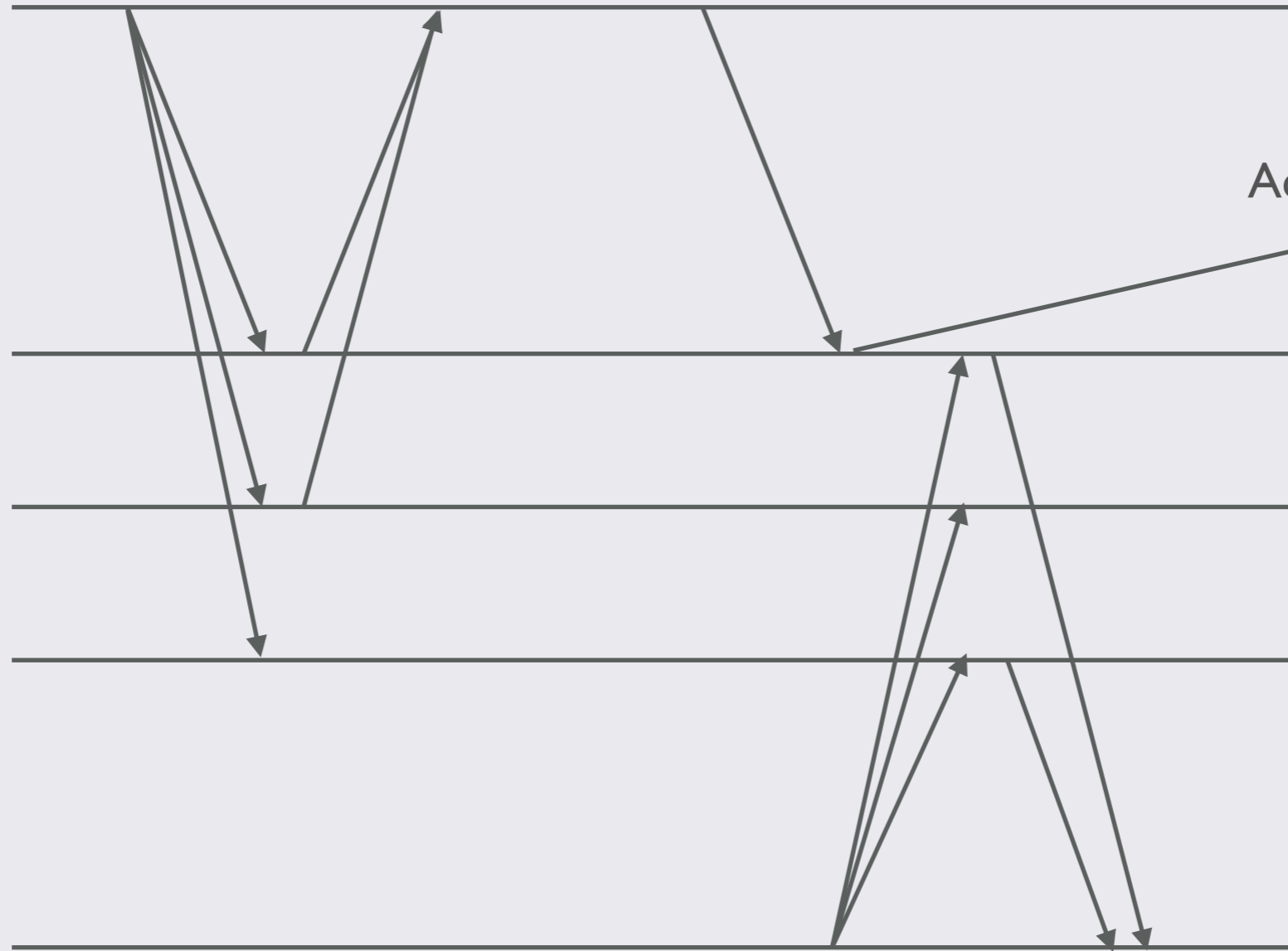


Accept

Proposer #2



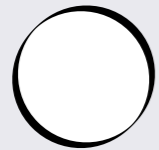
IAmLeader #2 YouAreLeader



# THE CRUCIAL *YouAreLeader* MESSAGE

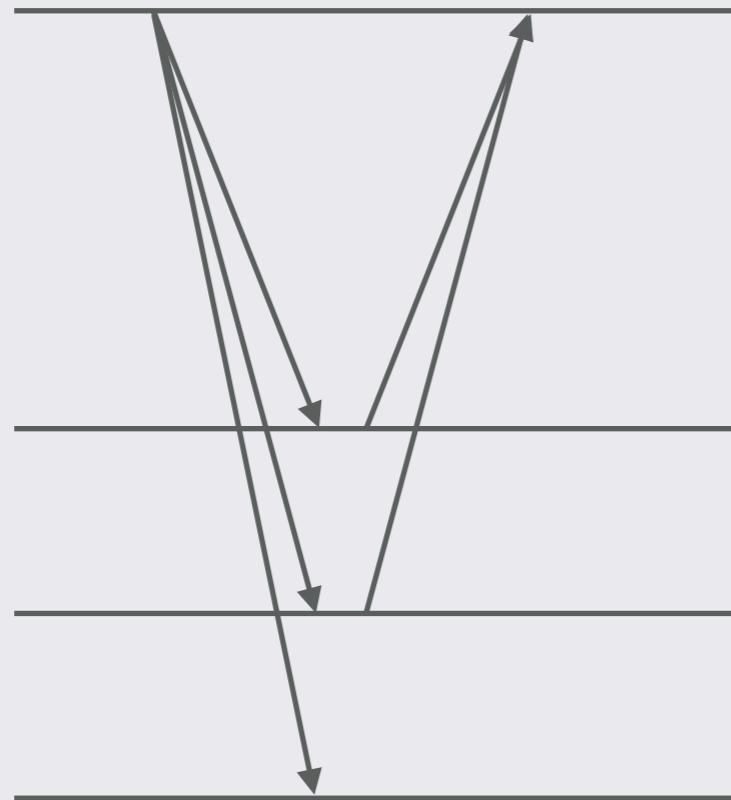
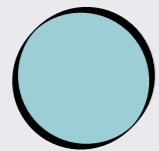
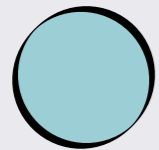
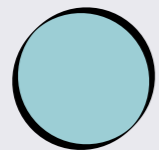
- I swear I won't follow an earlier leader!
- And, btw, here is my current accepted value (if any) by leader x.

Proposer #1



IAmLeader #1    *YouAreLeader*

Acceptors



1. Wait for a majority of *YouAreLeader* messages before proceeding.
2. If none of them contain a previously accepted value, propose your own  
Otherwise, propose the value of the **most recent** leader.

# THE CRUCIAL *YouAreLeader* MESSAGE

1. Wait for a majority of *YouAreLeader* messages before proceeding.
2. If none of them contain a previously accepted value, propose your own  
Otherwise, propose the value of the **most recent** leader.






## Important

By consulting a majority, the new leader makes sure she cannot have missed a chosen value

(a value must be accepted by a majority to be chosen, and any two majorities overlap!)






# EXAMPLES OF ACCEPTOR STATES

(as leader #50 comes to power)

Acceptors	Value	By leader
	x	37
	-	-
	-	-
	-	-
	-	-

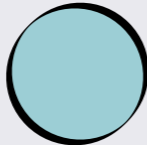

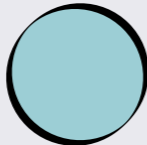


# EXAMPLES OF ACCEPTOR STATES

(as leader #50 comes to power)

Acceptors	Value	By leader
	x	37
	y	42
	-	-
	x	37
	y	41

# EXAMPLES OF ACCEPTOR STATES

(as leader #50 comes to power)

Acceptors	Value	By leader
	x	37
	y	42
	-	-
	x	37
	x	41

# OVERVIEW OF PAXOS

## Proposer

Send *IAmLeader(n)* to all  
Wait for a majority of responses

Once majority is received, send  
*Propose(n, V)* where  $V$  is the highest-leader  
proposal among the responses (or my own  
value, if none of the responses had a value)

## Acceptor

If  $n$  is the highest leader # I have seen:  
respond with  
*YouAreLeader(Value, LeaderWhoProposedValue)*

If  $n$  is the highest leader # I have seen, send  
*Accept(n, V)* to the learner



# TOLERATING $f$ FAILURES

## Safety

- There are  $2f + 1$  acceptors
- A value is only chosen if accepted by a majority ( $f + 1$ )
- So, even if  $f$  of those acceptors fail, one will remain and **will be part** of any future majority

## Liveness

- The leader always waits for  $f + 1$  responses. So, even if  $f$  replicas fail, it will not block

# THE THREAT TO LIVENESS: DUELING PROPOSERS

Greetings, peasants! I am  
your fearless leader #1! Grant me  
your blessing!

Greetings, peasants! I am  
your fearless leader #3! Grant me  
your blessing!

Greetings, peasants! I am  
your fearless leader #5! Grant me  
your blessing!

Greetings, peasants! I am  
your fearless leader #7! Grant me  
your blessing!

⋮

Greetings, peasants! I am  
your fearless leader #2! Grant me  
your blessing!

Greetings, peasants! I am  
your fearless leader #4! Grant me  
your blessing!

Greetings, peasants! I am  
your fearless leader #6! Grant me  
your blessing!

Greetings, peasants! I am  
your fearless leader #8! Grant me  
your blessing!

⋮

# THE THREAT TO LIVENESS: DUELING PROPOSERS

This problem can be avoided during synchrony  
(proposer faults can be detected accurately using timeouts)

It's **impossible** to avoid during asynchrony!

Well, we kind of knew that already...

# THE BEAUTY OF PAXOS

Paxos **cannot** be both safe and live during asynchrony!

(that would violate FLP)

So it's doing the next best thing:  
staying **safe all the time** and achieving liveness  
when the system starts behaving synchronously

# USING (MULTI)PAXOS TO IMPLEMENT STATE MACHINE REPLICATION

The original Paxos algorithm achieves agreement on **one** value

SMR required replicas to agree on the **sequence** of commands that will be executed

*3. Ensure that all replicas go through the same sequence of state transitions*



MultiPaxos: Run an instance of Paxos for each slot in the sequence

**Important:** we don't need to run phase **1** (election) every time!