# Distributed
# Coordination
# (Mutual Exclusion, Consensus)

# SURVEY FEEDBACK

- Breadth vs Depth
- Example Use Cases
- Project Difficulty
- Using cloud trial version – hybrid + on premise VMs
- Programming Language - Go

Prof. Tim Wood & Prof. Roozbeh Haghnazar

# SCHEDULE

- Remaining Topics
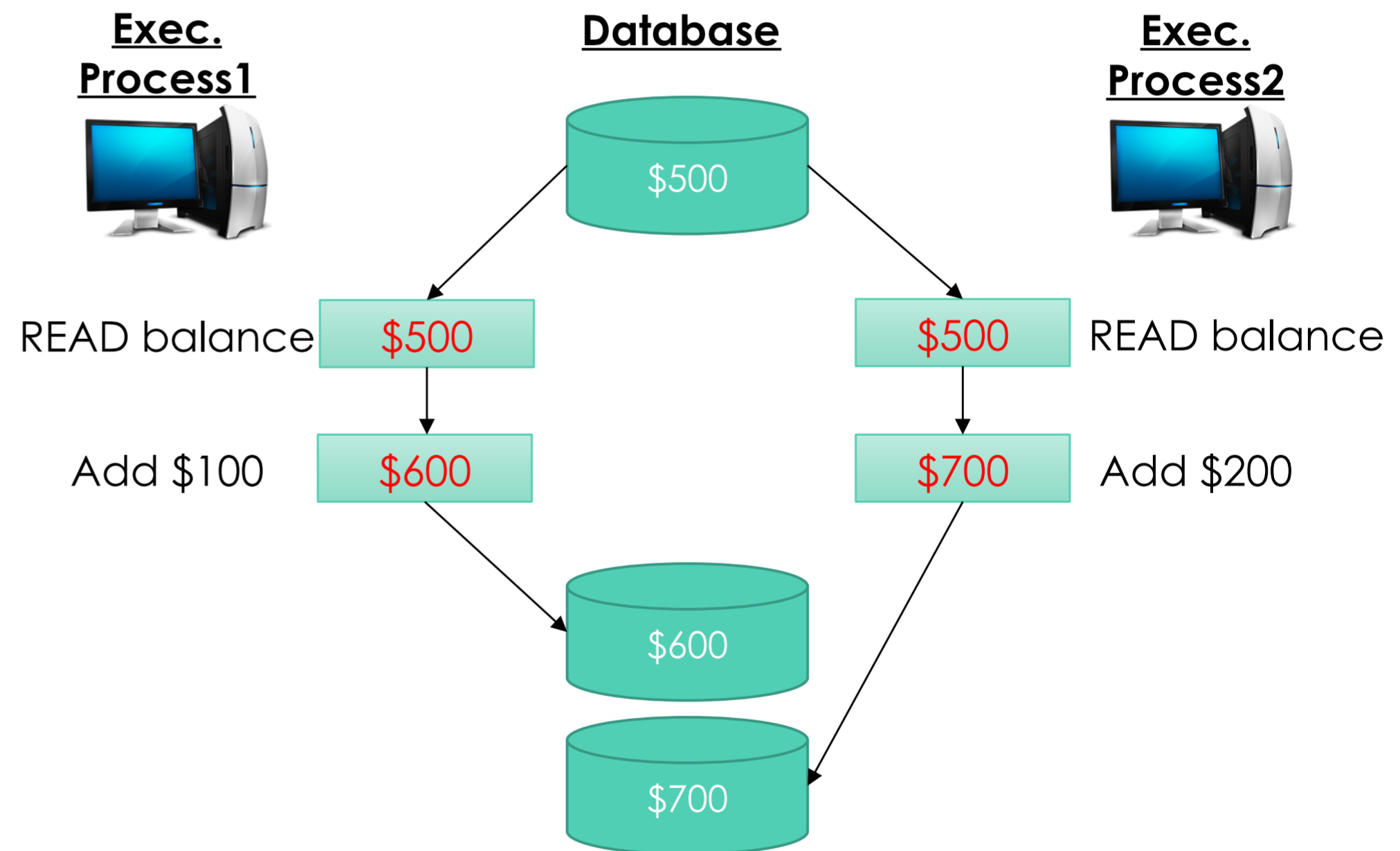- Midterm
- Final Project

Prof. Tim Wood & Prof. Roozbeh Haghnazar

# THIS WEEK: DISTRIBUTED COORDINATION

- Distributed Locking

- Consensus

- Elections

- State Machine Replication

- Blockchain

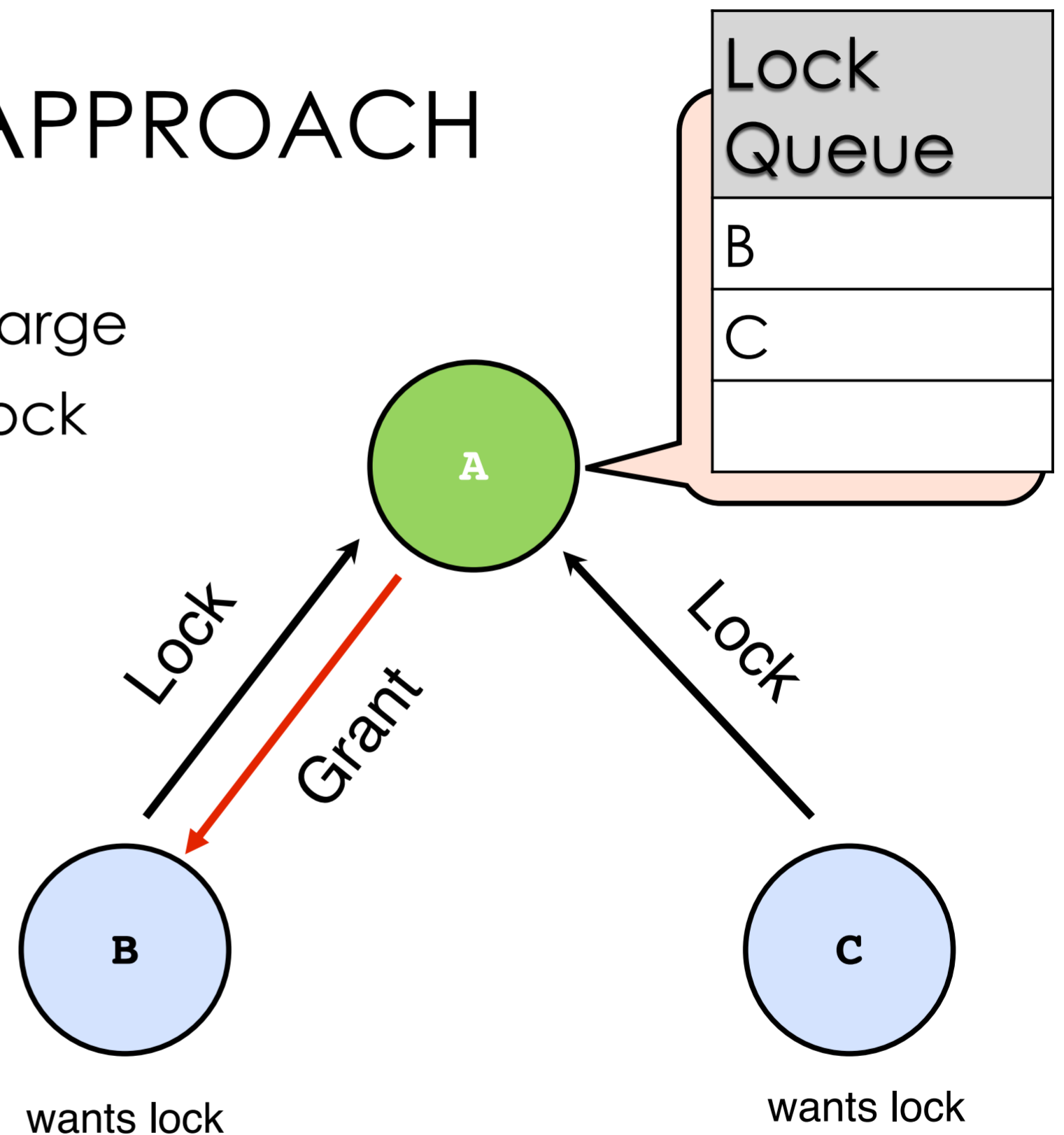Prof. Tim Wood & Prof. Roozbeh Haghnazar

# WHY LOCK?

- Locks let us protect a **shared** resource
  - A database, values in shared memory, files on a shared file system, throttle control on a drone, etc

- How to manage a lock in a distributed environment?

- How do locks limit scalability?

**Exec. Process1**

**Database**

**Exec. Process2**

$500

READ balance | $500 | | $500 | READ balance

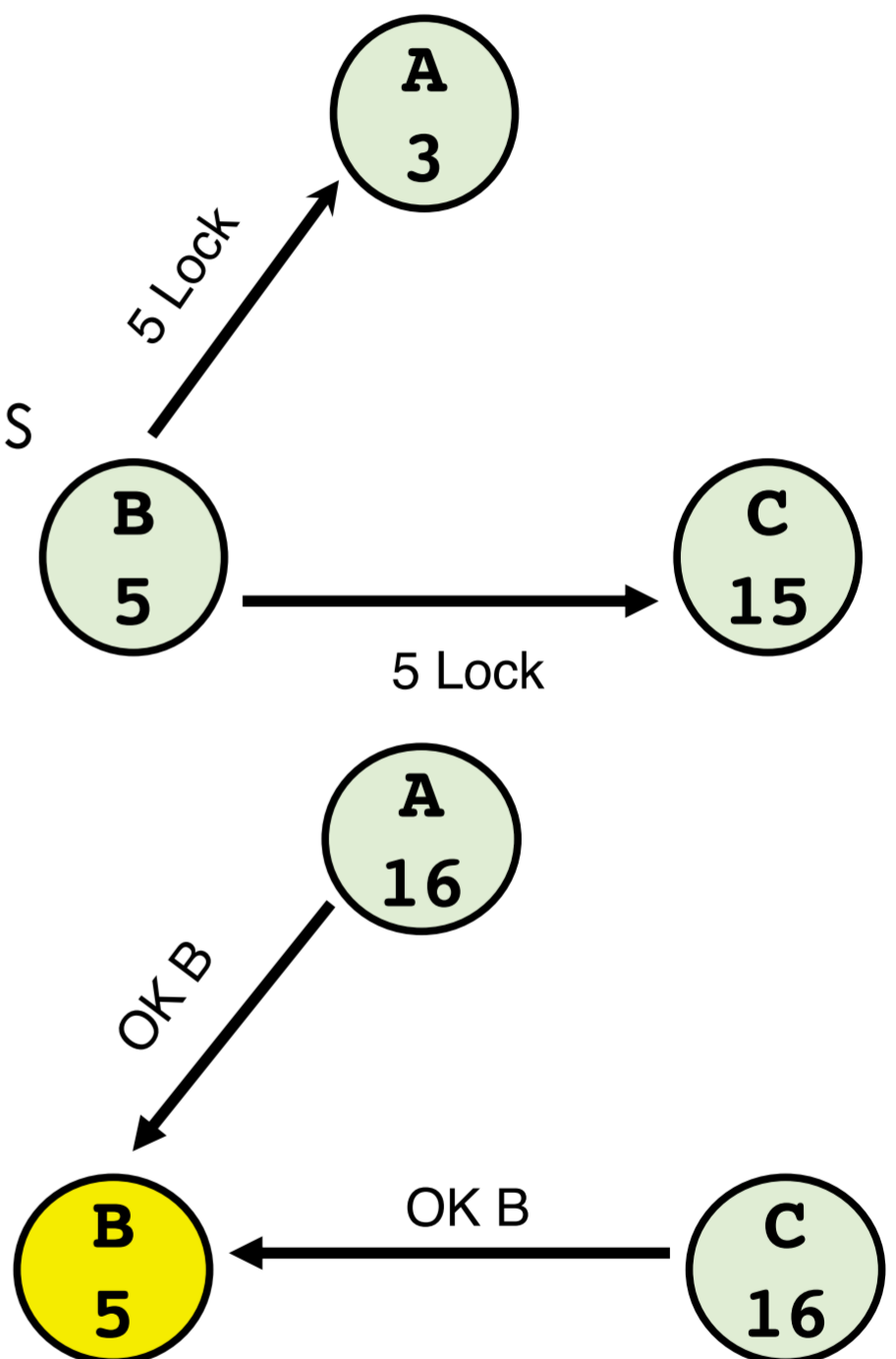Add $100 | $600 | | $700 | Add $200

$600

$700

# CENTRALIZED APPROACH

- Simplest approach: put one node in charge
- Other nodes ask coordinator for each lock
  - Block until they are granted the lock
  - Send release message when done
- Coordinator can decide what order to grant lock
- Do we get:
  - Mutual exclusion?
  - Progress?
  - Resilience to failures?
  - Balanced load?

Lock Queue

| |
|---|
| B |
| C |
| |

A

Lock

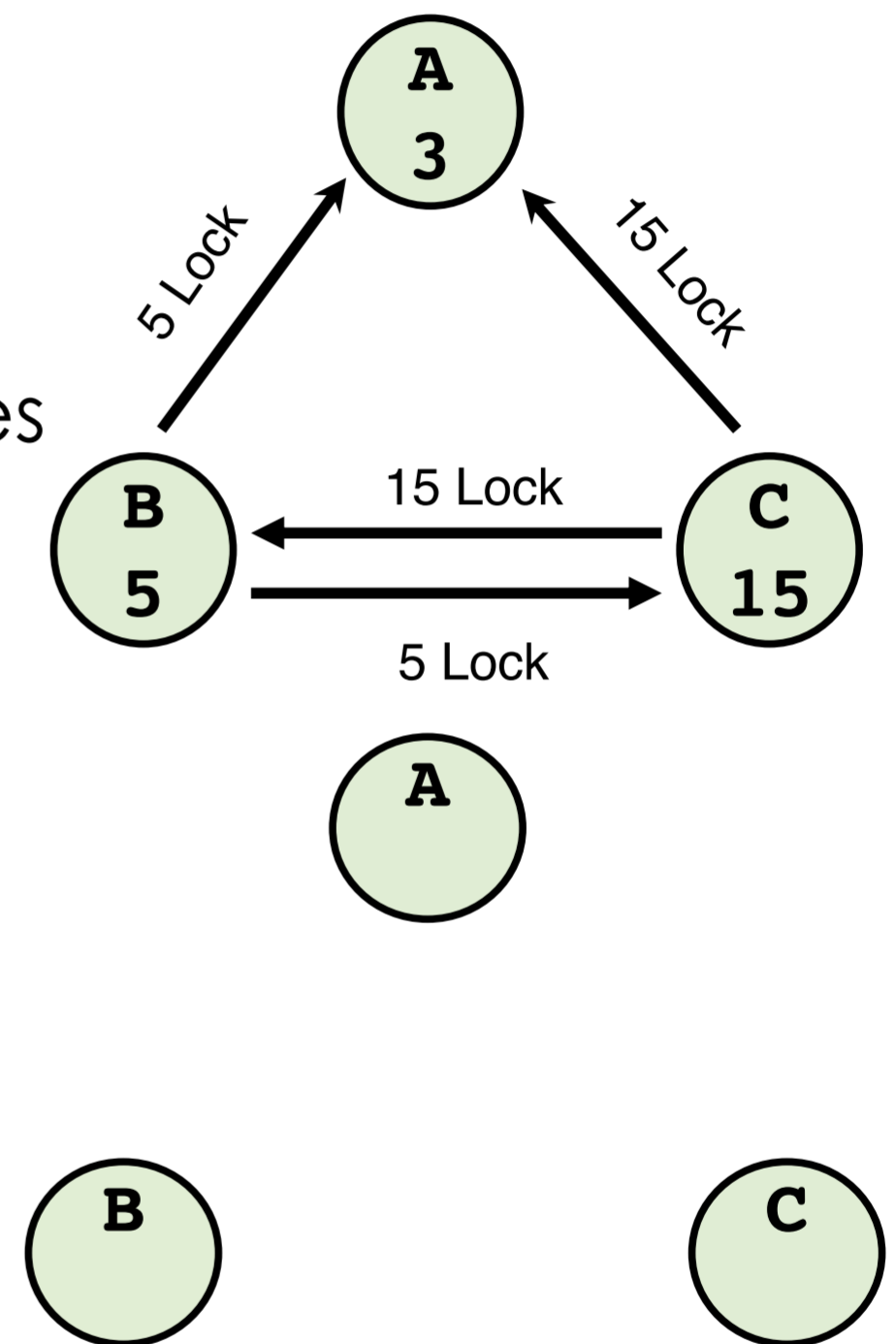Grant

Lock

B
wants lock

C
wants lock

# DISTRIBUTED APPROACH

- Use Lamport Clocks to order lock requests across nodes
- Send Lock message with ++clock
  - Wait for OKs from all nodes
- When receiving Lock msg:
  - Update clock following Lamport's rules
  - Send OK if not interested
  - If I want the lock:
    - Send OK if request's clock is smaller than own
    - Else, put request in queue
- When done with a lock:
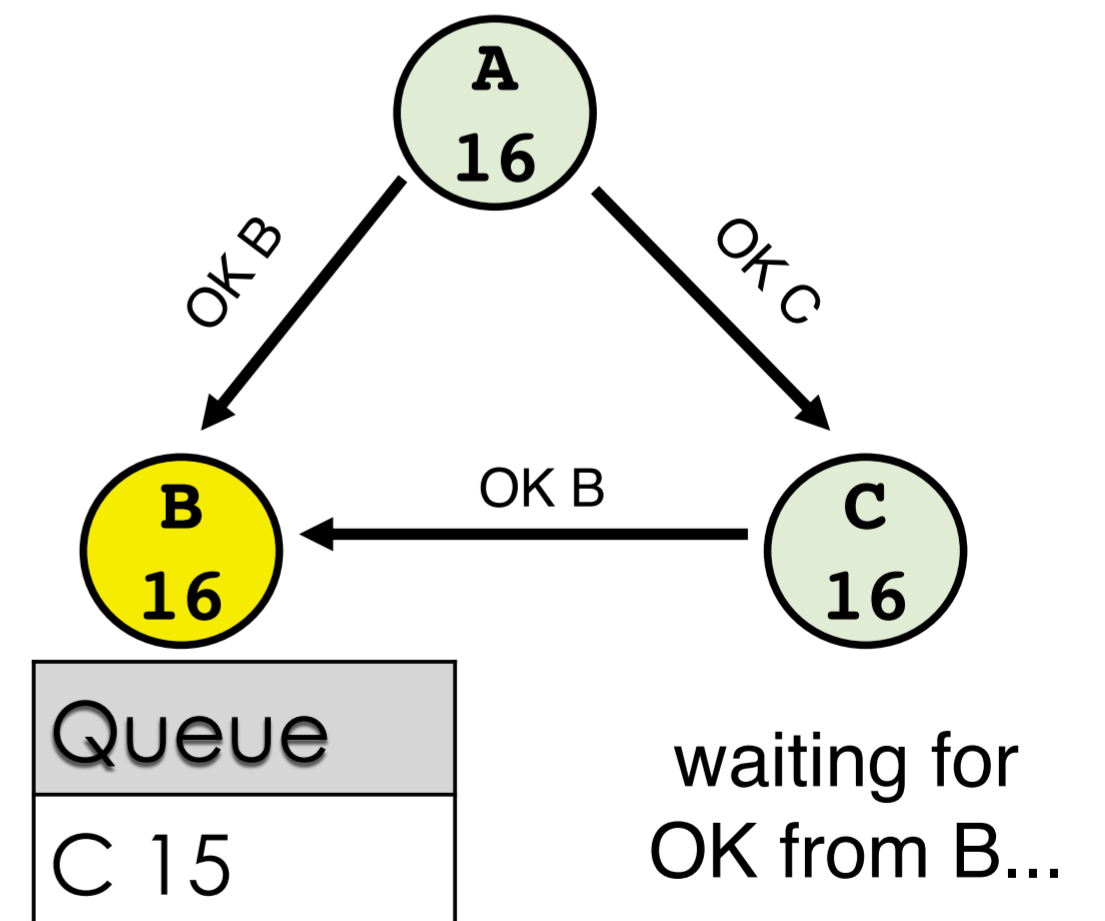  - Send OK to anybody in queue

# DISTRIBUTED APPROACH

- Use Lamport Clocks to order lock requests across nodes
- Send Lock message with ++clock
  - Wait for OKs from all nodes
- When receiving Lock msg:
  - Update clock following Lamport's rules
  - Send OK if not interested
  - If I want the lock:
    - Send OK if request's clock is smaller than own
    - Else, put request in queue
- When done with a lock:
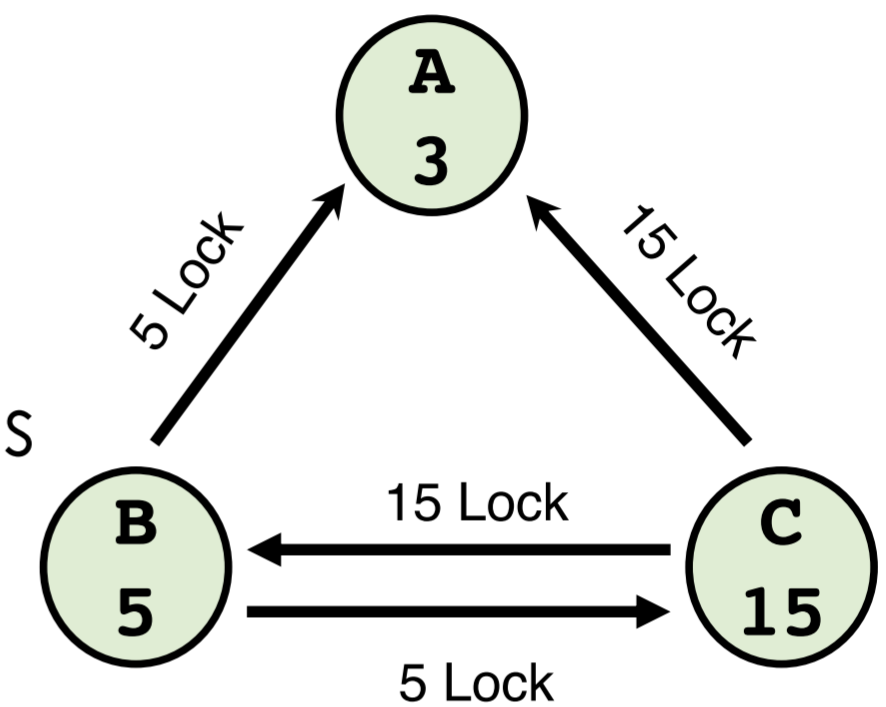  - Send OK to anybody in queue

# DISTRIBUTED APPROACH

- Use Lamport Clocks to order lock requests across nodes
- Send Lock message with ++clock
  - Wait for OKs from all nodes
- When receiving Lock msg:
  - Update clock following Lamport's rules
  - Send OK if not interested
  - If I want the lock:
    - Send OK if request's clock is smaller than own
    - Else, put request in queue
- When done with a lock:
  - Send OK to anybody in queue

A 3

5 Lock

15 Lock

B 5

15 Lock

C 15

5 Lock

A 16

OK B

OK C

B 16

OK B

C 16

Queue

C 15

waiting for OK from B...

# COMPARISON

- Messages per lock acquire and release
  - Centralized:
  - Distributed:

- Delay before entry
  - Centralized:
  - Distributed:

- Problems
  - Centralized:
  - Distributed:

# COMPARISON

- Messages per lock acquire and release
  - Centralized: 2+1=3
  - Distributed: 2(n-1)
- Delay before entry
  - Centralized: 2
  - Distributed: 2(n-1) in parallel
- Problems
  - Centralized: Coordinator crashes
  - Distributed: anybody crashes

> Is the distributed approach better in any way?

Prof. Tim Wood & Prof. Roozbeh Haghnazar

# DISTRIBUTED SYSTEMS ARE HARD

- Going from centralized to distributed can be..

- Slower
  - If everyone needs to do more work

- More error prone
  - 10 nodes are 10x more likely to have a failure than one

- Much more complicated
  - If you need a complex protocol
  - If nodes need to know about all others

Often we need more than just a way to lock a resource!

Prof. Tim Wood & Prof. Roozbeh Haghnazar

# WHAT IS THE MEANING OF CONSENSUS

- Consensus is defined by Merriam-Webster as,
  - *general agreement,*
  - *group solidarity of belief or sentiment.*
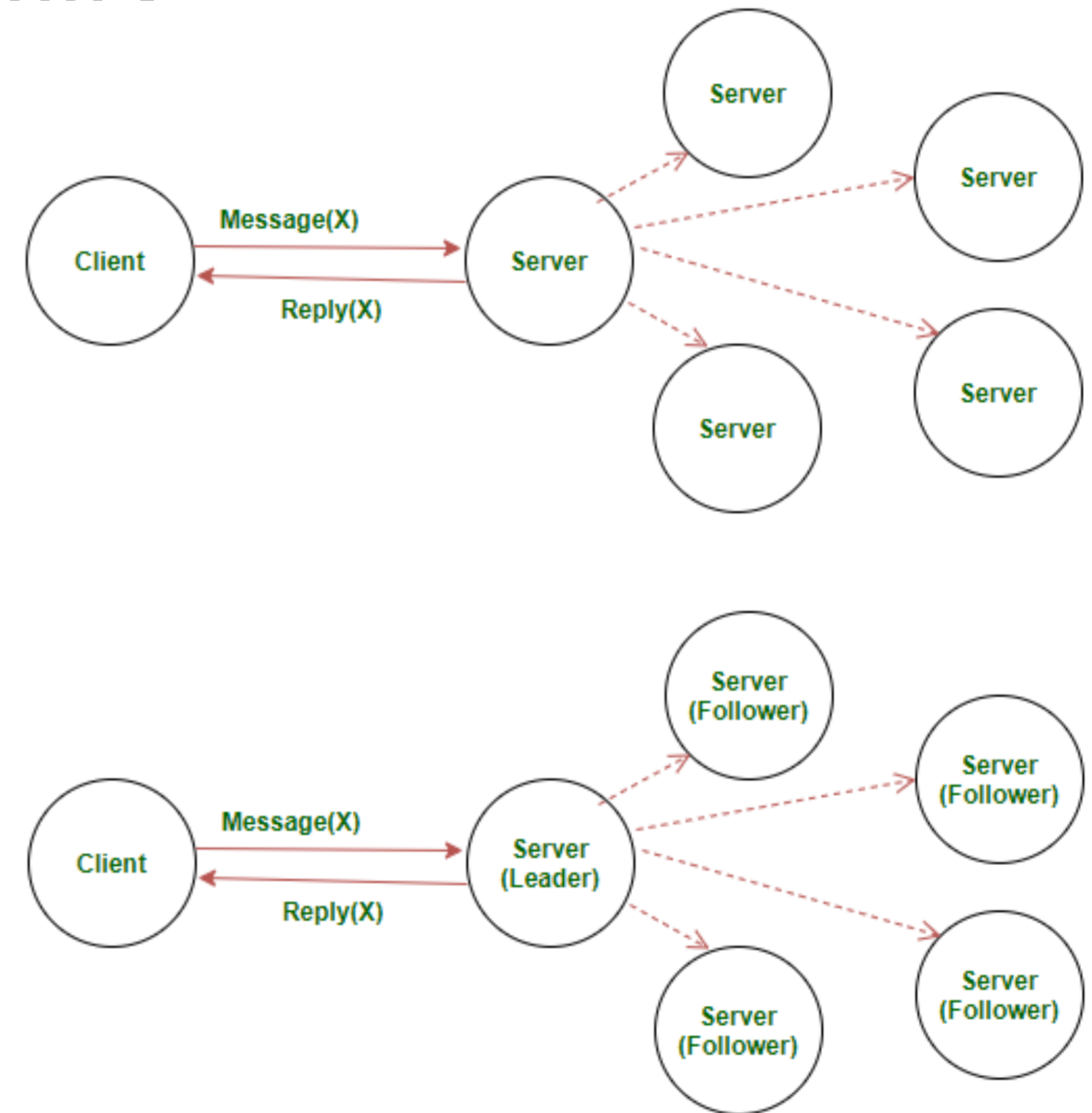
# WHY CONSENSUS?

When you sent a request to a server it answers you easily



- If server fails, there is no backup
- If the number of requests increase dramatically the server won't be able to respond

# WHY CONSENSUS?

- Symmetric :- Any of the multiple servers can respond to the client and all the other servers are supposed to sync up with the server that responded to the client's request, and

- Asymmetric :- Only the elected leader server can respond to the client. All other servers then sync up with the leader server.

# WHY CONSENSUS?

While this creates a system that is devoid of corruption from a single source, it still creates a major problem.

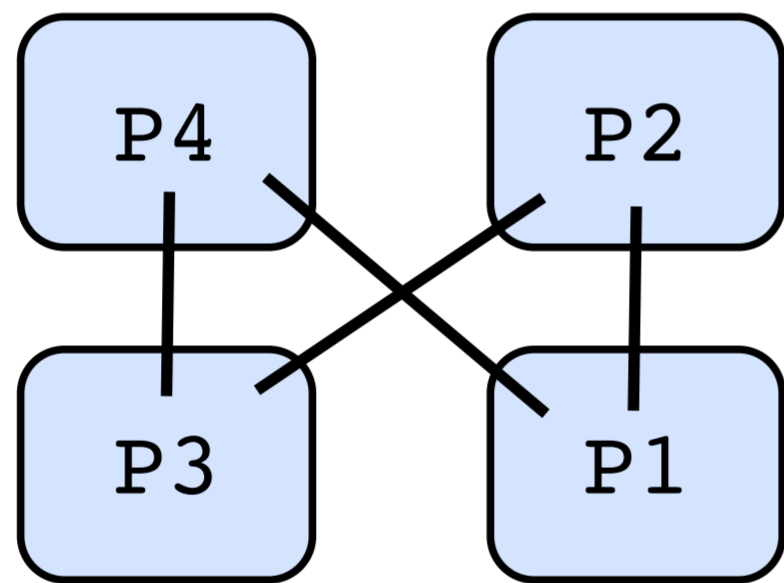- How are any decisions made?
- How does anything get done?

Prof. Tim Wood & Prof. Roozbeh Haghnazar
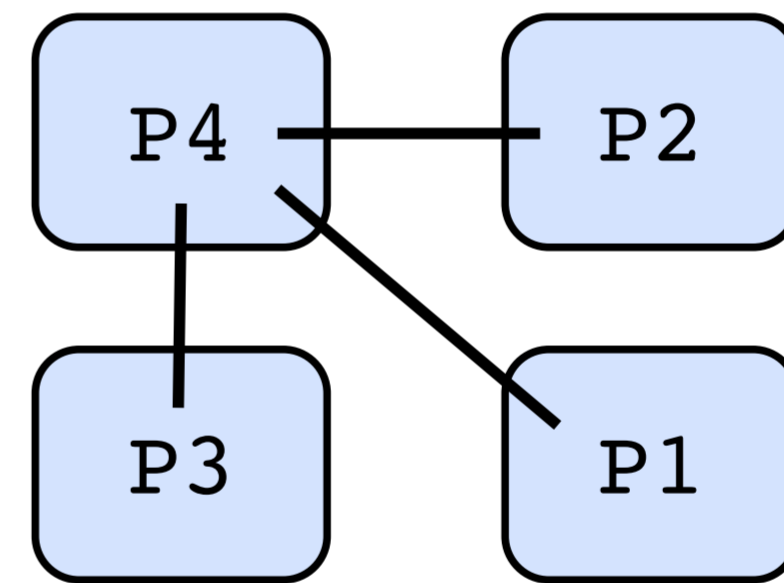
# CONSENSUS OBJECTIVES

- Therefore, objectives of a consensus mechanism are:
  - **Agreement seeking**: A consensus mechanism should bring about as much agreement from the group as possible.
  - **Collaborative**: All the participants should aim to work together to achieve a result that puts the best interest of the group first.
  - **Cooperative**: All the participants shouldn't put their own interests first and work as a team more than individuals.
  - **Egalitarian**: A group trying to achieve consensus should be as egalitarian as possible. What this basically means that each and every vote has equal weight. One person's vote can't be more important than another's.
  - **Inclusive**: As many people as possible should be involved in the consensus process. It shouldn't be like normal voting where people don't really feel like voting because they believe that their vote won't have any weight in the long run.
  - **Participatory**: The consensus mechanism should be such that everyone should actively participate in the the overall process.

Prof. Tim Wood & Prof. Roozbeh Haghnazar

# DISTRIBUTED ARCHITECTURES

- Purely distributed / decentralized architectures are difficult to run correctly and efficiently (decentralized locking was pretty bad!)
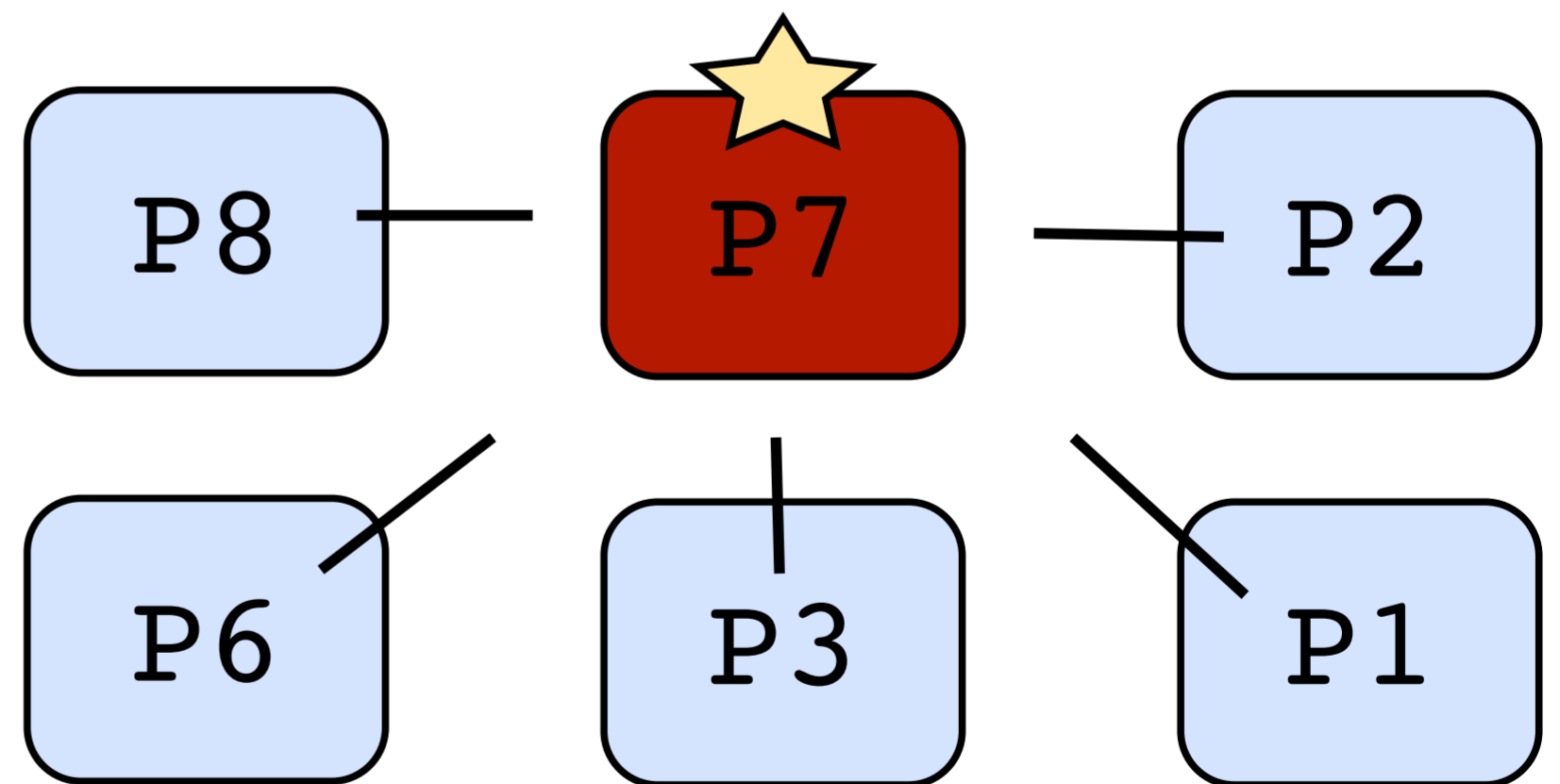


**Decentralized**



**Centralized**

- Can we mix the two?

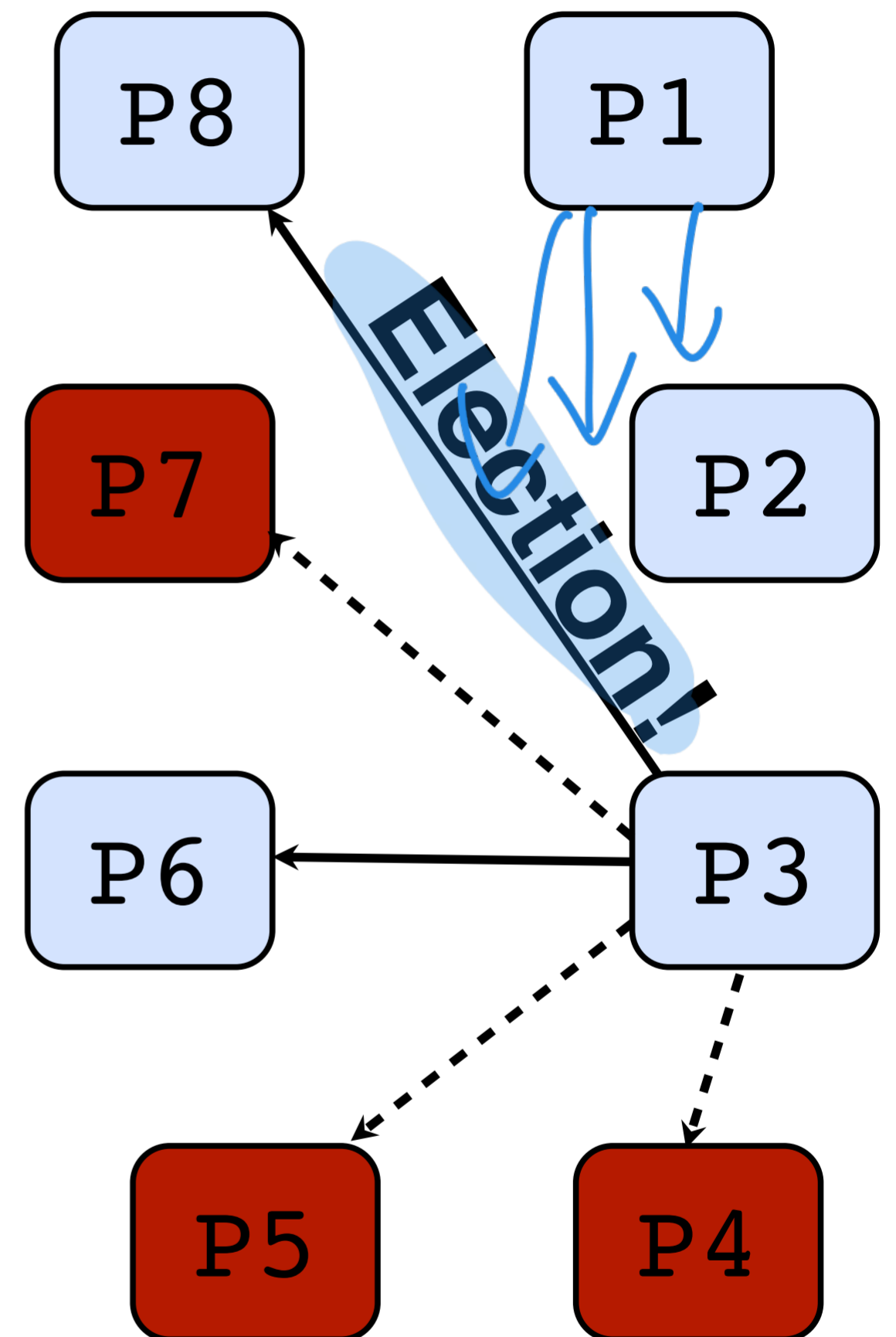Prof. Tim Wood & Prof. Roozbeh Haghnazar

# ELECTIONS

- Appoint a central coordinator
  - But allow them to be replaced in a safe, distributed way

- Must be able to handle simultaneous elections
  - Reach a consistent result

- Who should win?

# BULLY ALGORITHM

- The biggest (ID) wins
- Any process P can initiate an election
- P sends **Election** messages to all process with higher Ids and awaits **OK** messages
- If it receives an OK, it drops out and waits for an **I won**
- If a process receives an **Election** msg, it returns an **OK**...

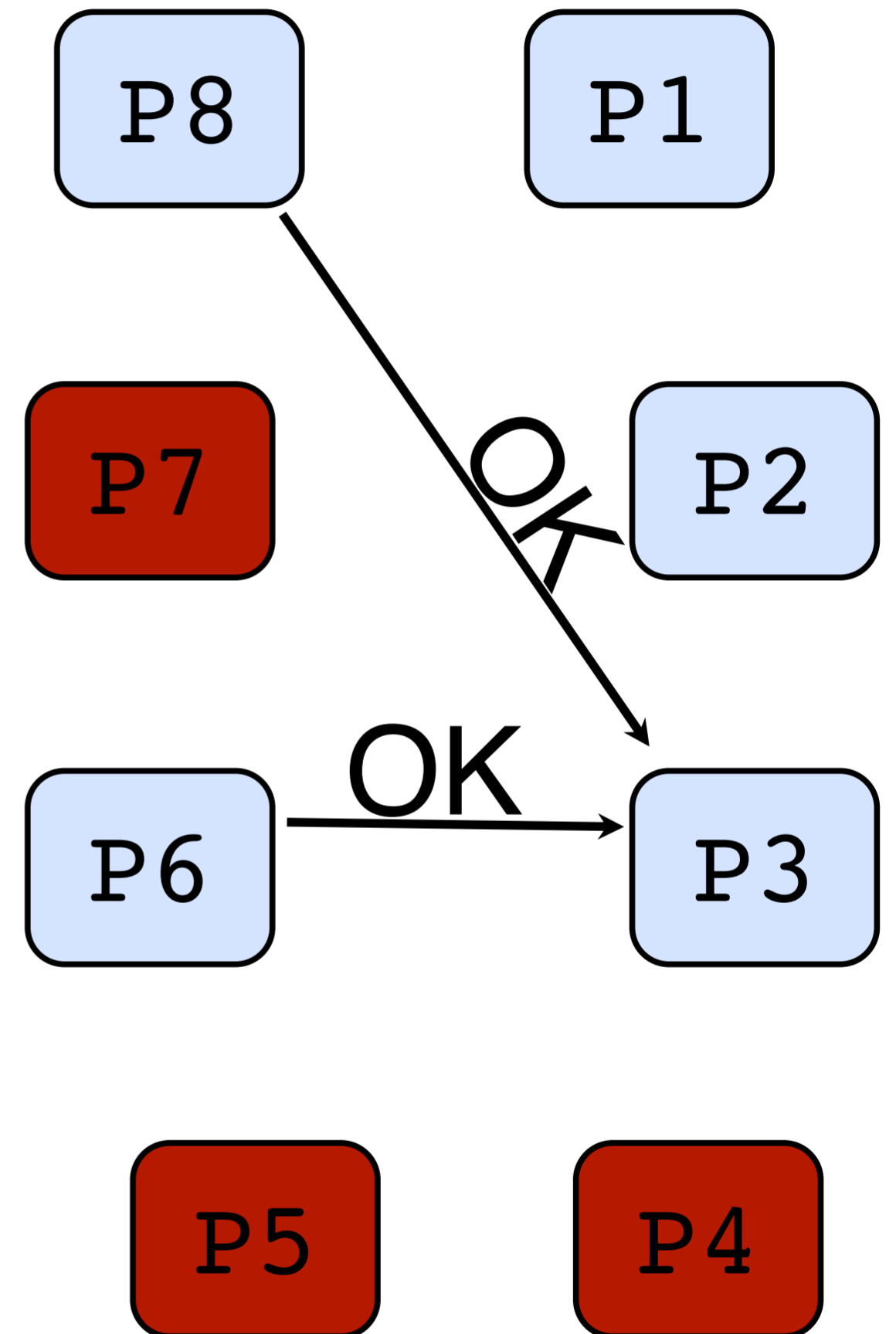Prof. Tim Wood & Prof. Roozbeh Haghnazar

# Bully Algorithm

- The biggest (ID) wins
- Any process P can initiate an election
- P sends **Election** messages to all process with higher Ids and awaits **OK** messages
- If it receives an OK, it drops out and waits for an **I won**
- If a process receives an **Election** msg, it returns an **OK**...

P8   P1

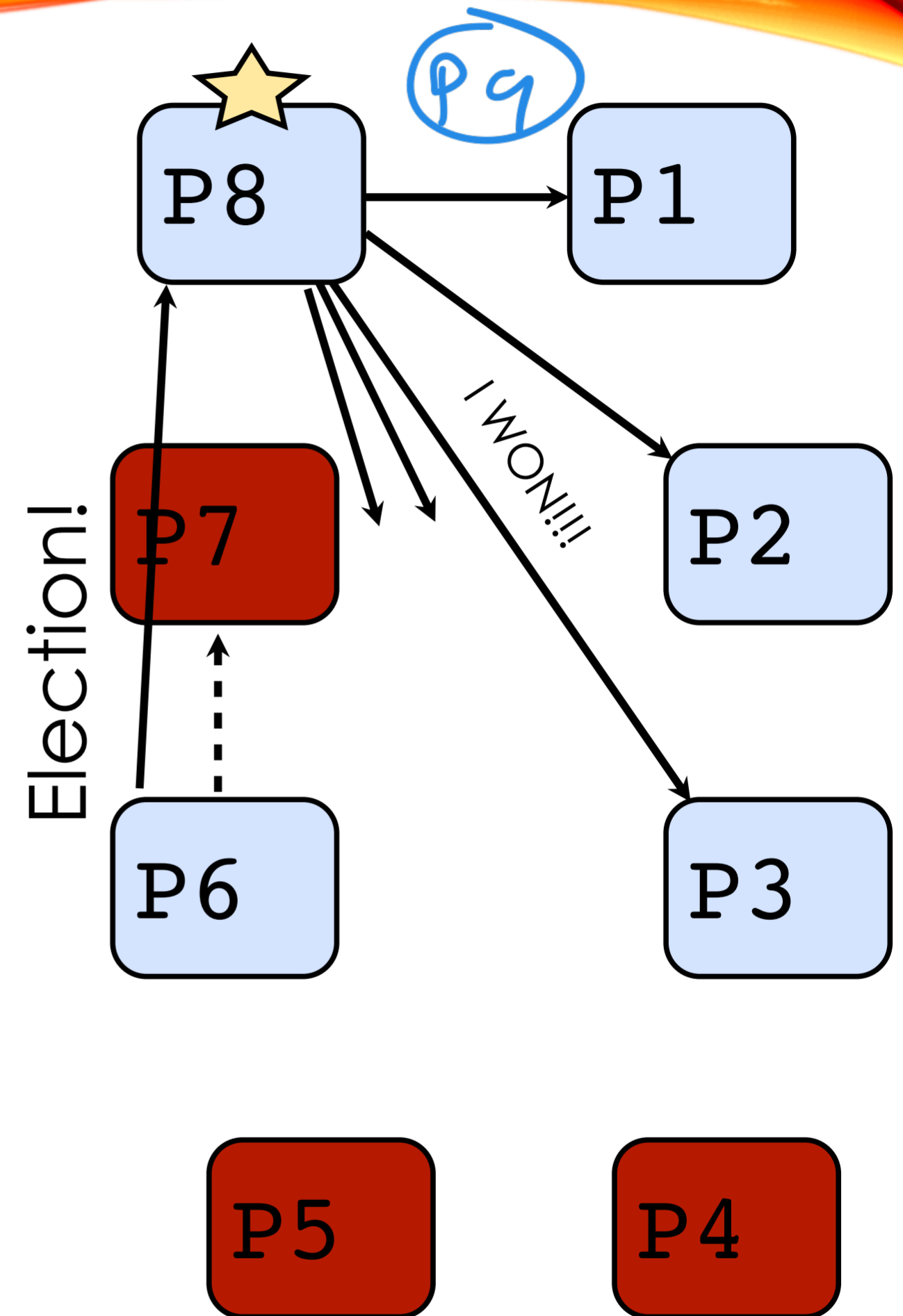P7   OK   P2

P6   OK   P3

P5   P4

# Bully Algorithm

- The biggest (ID) wins

- Any process P can initiate an election

- P sends **Election** messages to all process with higher Ids and awaits **OK** messages

- If it receives an OK, it drops out and waits for an **I won**

- If a process receives an **Election** msg, it returns an **OK** and starts **another** election

- If no **OK** messages, P becomes leader and sends **I won** to all process with lower Ids

- If a process receives a **I won**, it treats sender as the leader

# Ring Algorithm

P8    P1

- Any other ideas?

P7    P2

P6    P3

P5    P4

# RING ALGORITHM

- **Initiator** sends an **Election** message around the ring

- Add your ID to the message

- When Initiator receives message again, it announces the winner
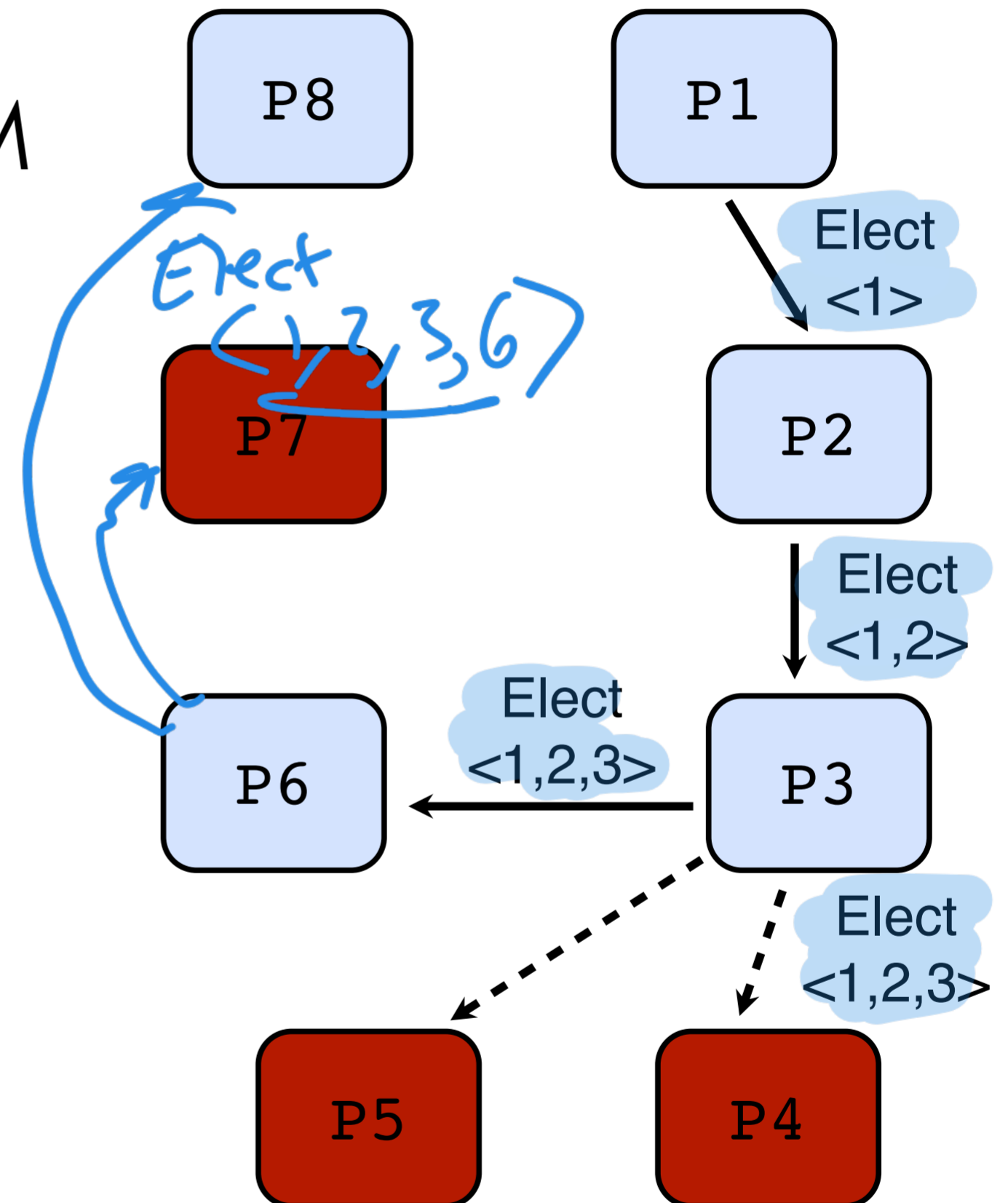
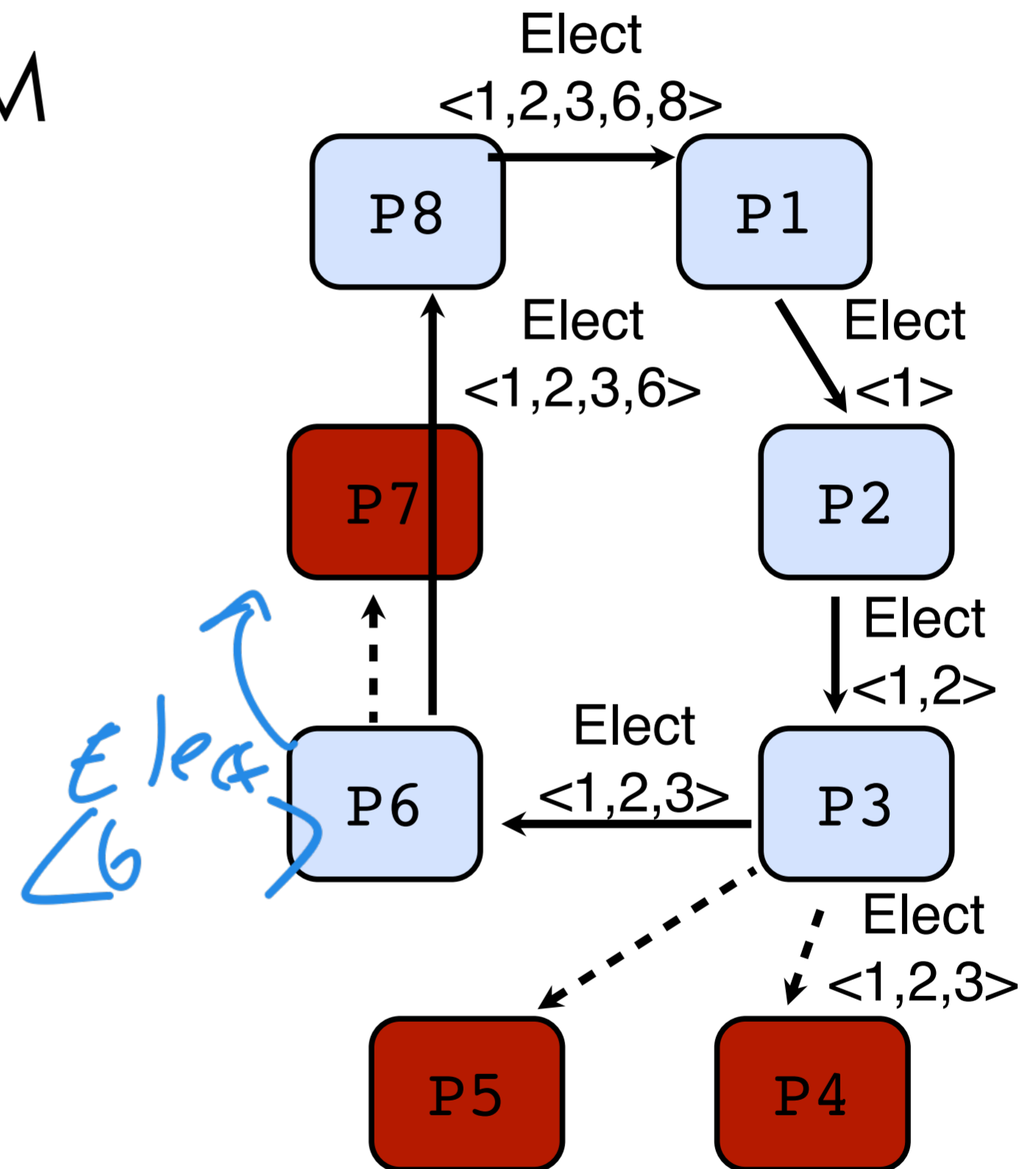- What happens if multiple elections occur at the same time?

# RING ALGORITHM

- **Initiator** sends an **Election** message around the ring

- Add your ID to the message

- When Initiator receives message again, it announces the winner

- What happens if multiple elections occur at the same time?

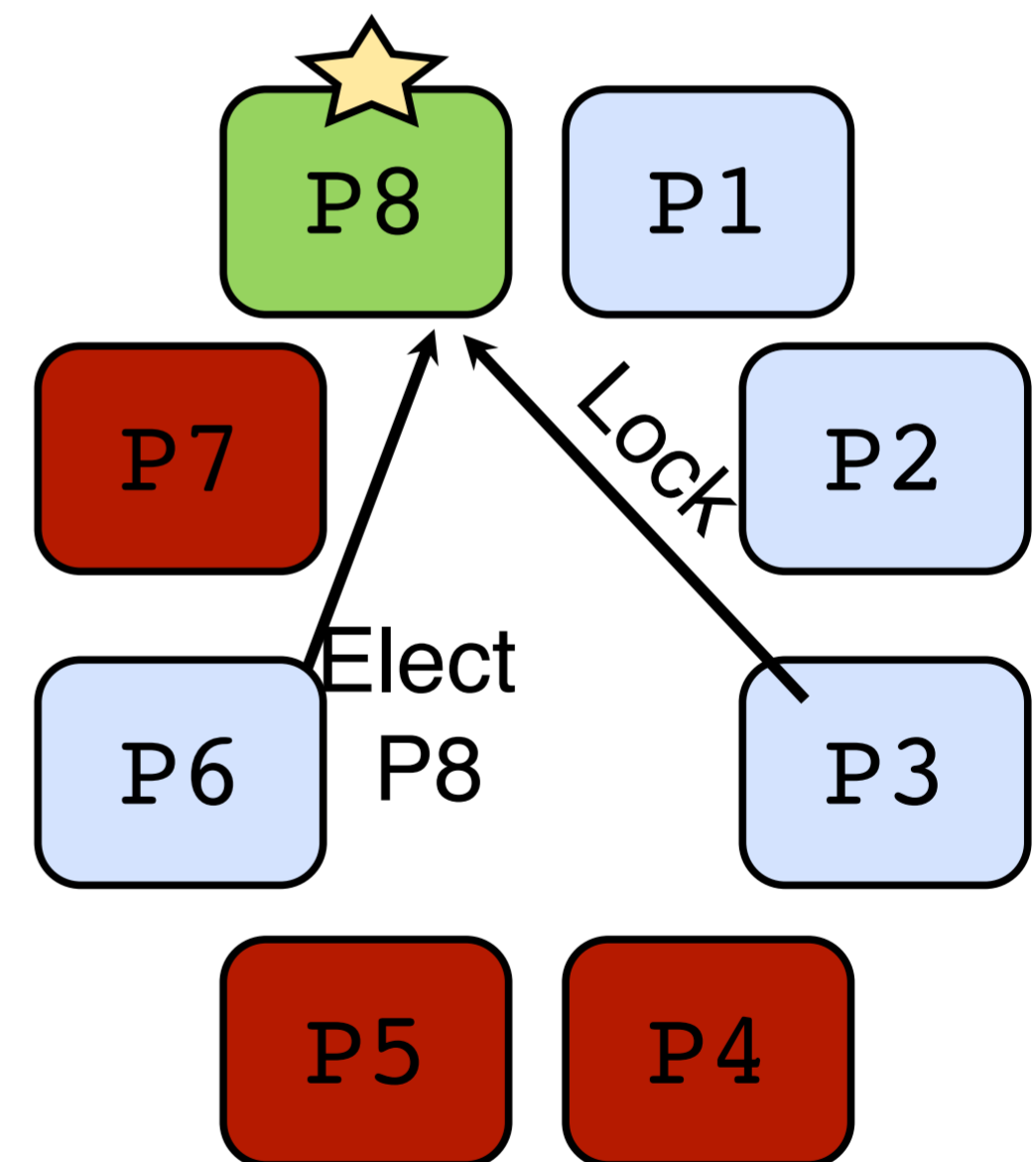# COMPARISON

- Number of messages sent to elect a leader:

- Bully Algorithm
  - Worst case: lowest ID node initiates election
    - Triggers n-1 elections at every other node = O(n^2) messages
  - Best case: Immediate election after n-2 messages

- Ring Algorithm
  - Always 2(n-1) messages
  - Around the ring, then notify all

Prof. Tim Wood & Prof. Roozbeh Haghnazar

# Elections + Centralized Locking

- Elect a leader

- Let them make all the decisions about locks

- What kinds of failures can we handle?
  - Leader/non-leader?
  - Locked/unlocked?
  - During election?

This can be the basis for **consensus-**based distributed systems!

P8

P1

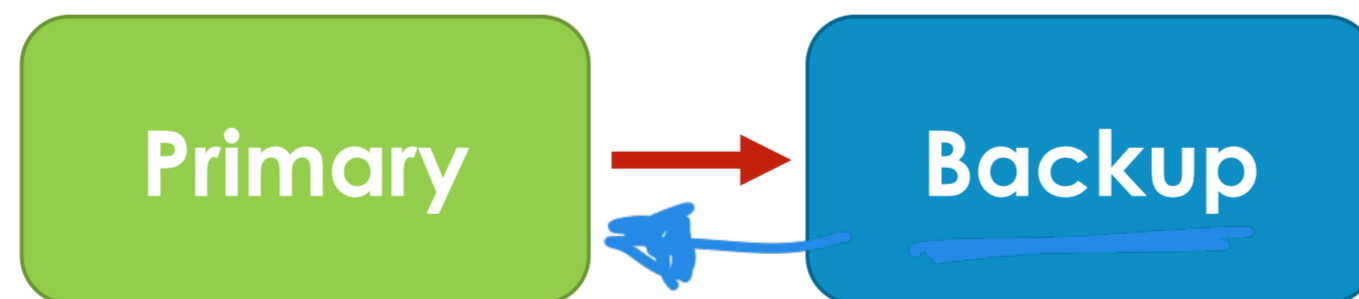P7

P2

Lock

P6

Elect
P8

P3

P5

P4

# CHUBBY: GOOGLE'S LOCK SERVICE

- Google services are composed of many thousands of nodes

- Need a way to coordinate data and access to shared resources!
  - Used by Google File System, BigTable, etc

- Chubby: lock service for loosely coupled distributed systems
  - Focuses on availability and reliability (not performance)
  - Scales to ~10,000 servers per Chubby Cell

- See paper at OSDI 2006 by Mike Burrows for full details!

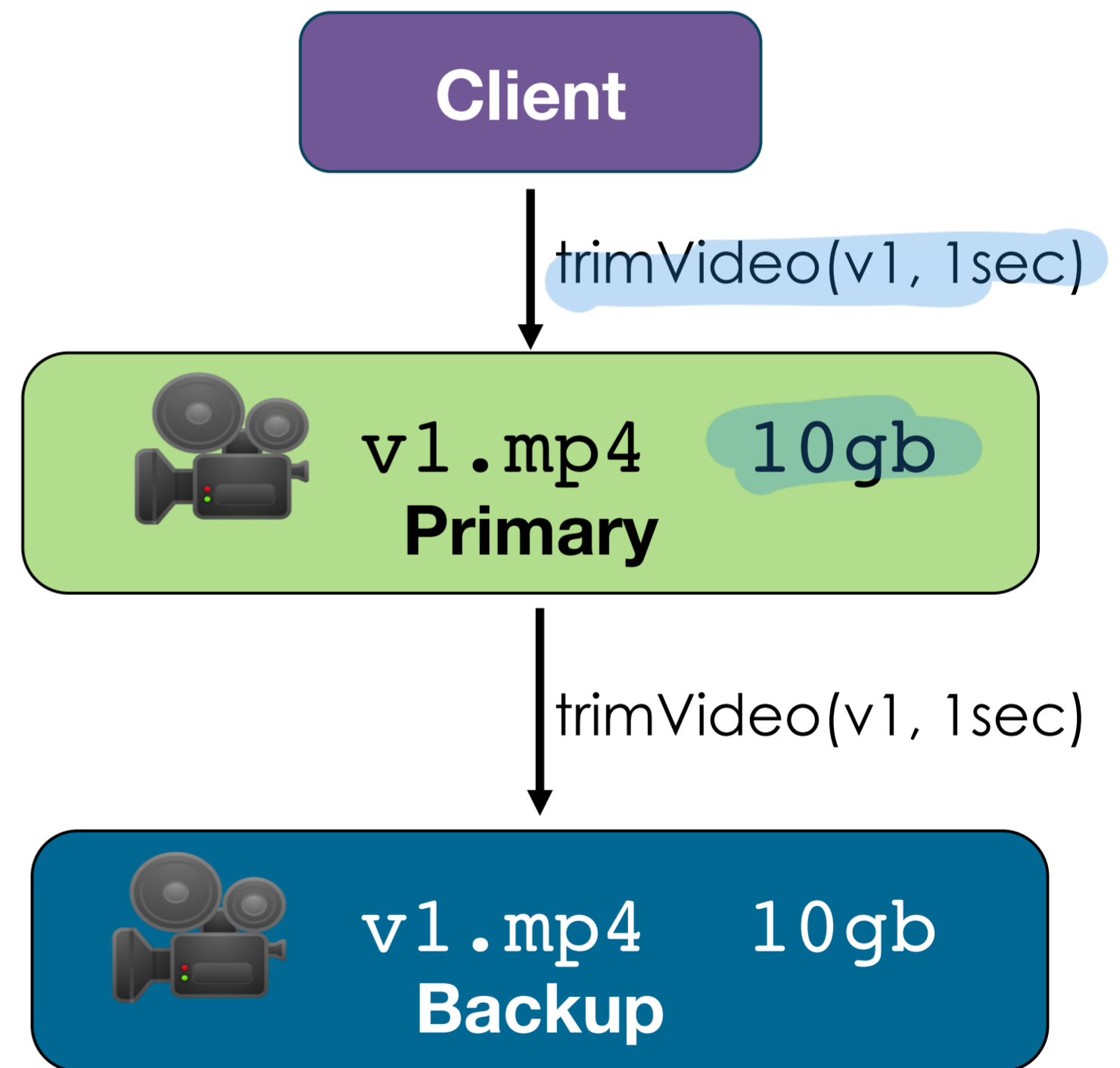| | |
|---|---|
| time since last fail-over | 18 days |
| fail-over duration | 14s |
| active clients (direct) | 22k |
| additional proxied clients | 32k |
| files open | 12k |
|    naming-related | 60% |
| client-is-caching-file entries | 230k |
| distinct files cached | 24k |
| names negatively cached | 32k |
| exclusive locks | 1k |
| shared locks | 0 |
| stored directories | 8k |
|    ephemeral | 0.1% |
| stored files | 22k |
|    0-1k bytes | 90% |
|    1k-10k bytes | 10% |
|    > 10k bytes | 0.2% |
|    naming-related | 46% |
|    mirrored ACLs & config info | 27% |
|    GFS and Bigtable meta-data | 11% |
|    ephemeral | 3% |
| RPC rate | 1-2k/s |
|    KeepAlive | 93% |
|    GetStat | 2% |
|    Open | 1% |
|    CreateSession | 1% |
|    GetContentsAndStat | 0.4% |
|    SetContents | 680ppm |
|    Acquire | 31ppm |

# STATE MACHINE REPLICATION (SMR)

- We can think of an application as a state machine
  - A program is just **data** that is updated based on **operations** -> **state**

- Consensus means that all distributed nodes should be in the same state!
  - If a node fails, it should not disrupt the system
  - When a node recovers it should be able to "catch up"

Primary → Backup
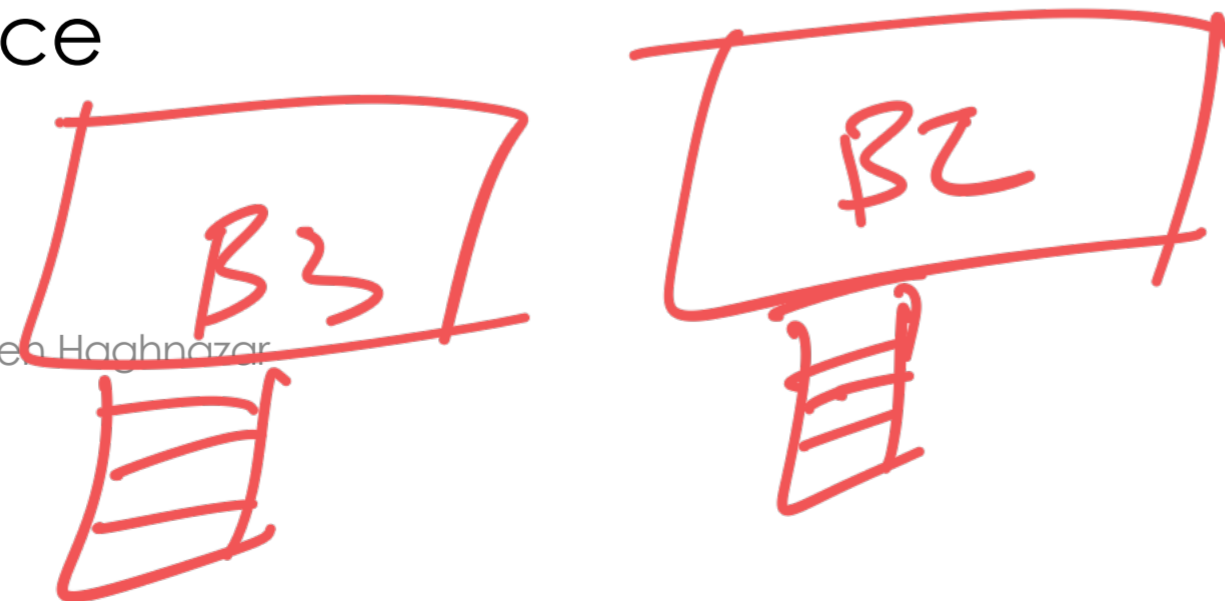
# Distributed Video Editing SMR

- Sometimes **data** is big!

- Replicate the **operation** to be performed, not the data!

- Treat like a state machine
  - Incoming requests just perform some operation on that data
  - If all replicas perform same operations, they will end in the same state

- If **Primary** fails, switch to **Backup**

Prof. Tim Wood & Prof. Roozbeh Haghnazar

**Client**

trimVideo(v1, 1sec)

v1.mp4    10gb
**Primary**

trimVideo(v1, 1sec)

v1.mp4    10gb
**Backup**

# Hash Table SMR

x = ?

- SMR creates a **replicated log** of actions to be performed
  - E.g., updates to the value stored by a key

- Primary orders incoming requests to form the log

- Actions must be deterministic

- We can keep adding more backup replicas to improve fault tolerance



C-1    C-2    C-3

set(x=3)    inc(x)    set(x=99)

| Log |
|-----|
| x=3 |
| x=99 |
| x++ |

Hash Table    x= 110
**Primary**

set(x=3)
set(x=99)
inc(x)

| Log |
|-----|
| x=3 |
| x=99 |
| inc(x) |

Hash Table    x=100
**Backup**

B3    B2

# SMR FAILURES?

- What to do on a failure?

- How many failures can we handle?



C-1  C-2  C-3

set(x=3)  inc(x)  set(x=99)

**Log**

| |
|---|
| x=3 |
| x=99 |
| |

Hash Table     x=
**Primary**

set(x=3)
set(x=99)

**Log**

| |
|---|
| x=3 |
| x=99 |
| |

Hash Table     x=
**Backup** Primary

# HANDLING FAILURES

- F = number of nodes which can crash at one time

- # of nodes needed must depend on f!

**f=1, f+1=2**
replicas

| Log |
|-----|
| x=1 |
| x=5 |

1: Primary

**Client**

2: Backup

| Log |
|-----|
| x=2 |
| x=5 |
| x=10 |

## What failure scenarios can happen?

# HANDLING FAILURES

- F = number of nodes which can crash at one time
- # of nodes needed must depend on f!

f=1, f+1
replicas

| Log |
| --- |
| |
| |
| |

**1: Primary**

**Client**

| Log |
| --- |
| |
| |
| |

**2: Backup**

Can't resync state if failure
"flip flops" between nodes!

f=1 **f+2 = 3**
replicas

| Log |
| --- |
| x=1 |
| |
| |

**1: Primary**

Client

| Log |
| --- |
| x=1 |
| x=2 |
| |

**2: Backup**

| Log |
| --- |
| x=1 |
| x=2 |
| |

**3: Backup**

Fixed?

# HANDLING FAILURES

- F = number of nodes which can crash at one time
- # of nodes needed must depend on f!

**f=2, f+2 = 4 replicas**

**Client**

**1: Primary**

**2: Backup**

**3: Backup**

**4: Backup**

f=1, f+2 = 3 replicas

**1: Primary**   Log

**2: Backup**   Log

**3: Backup**   Log

**Fixed for f=2?**
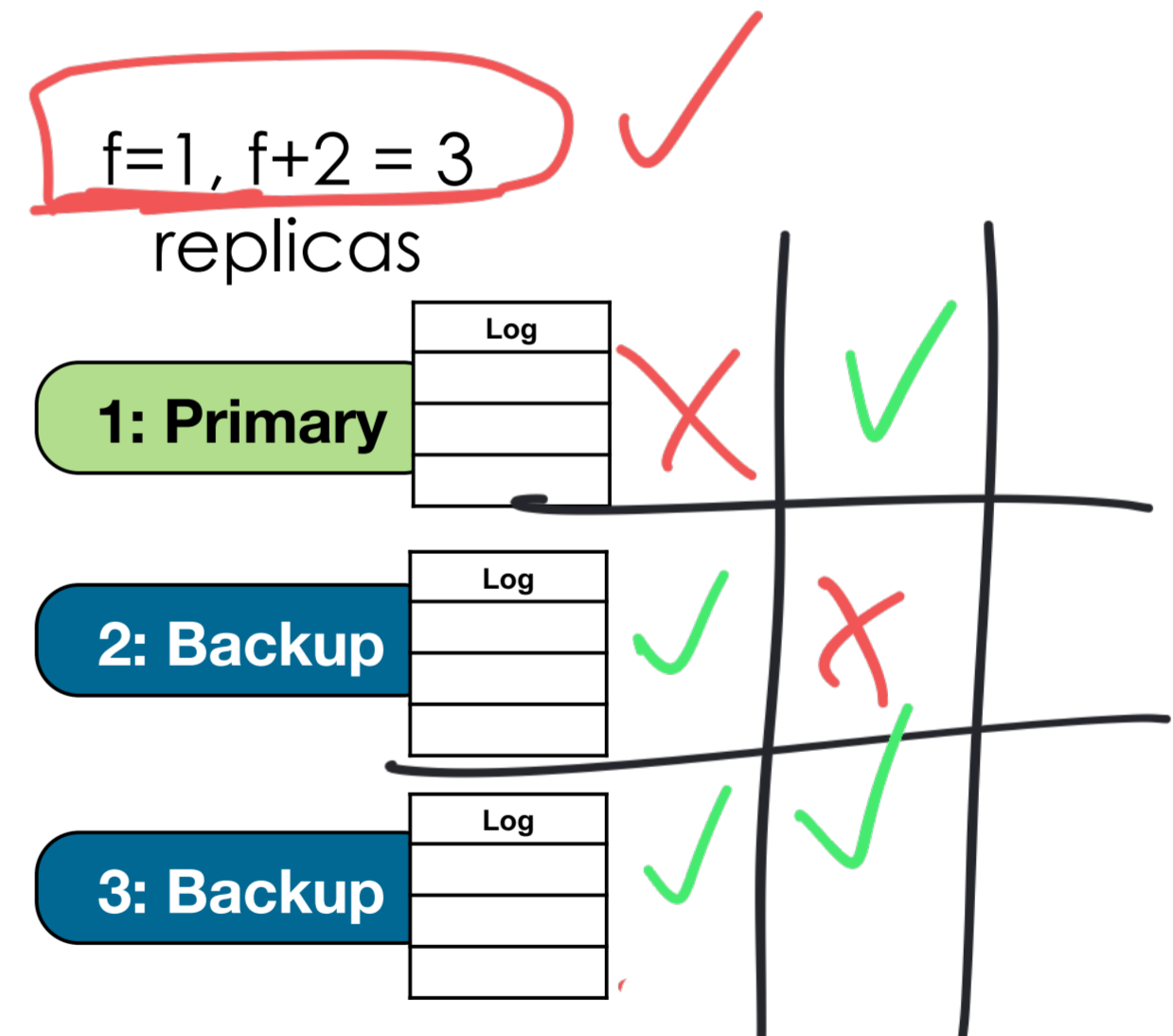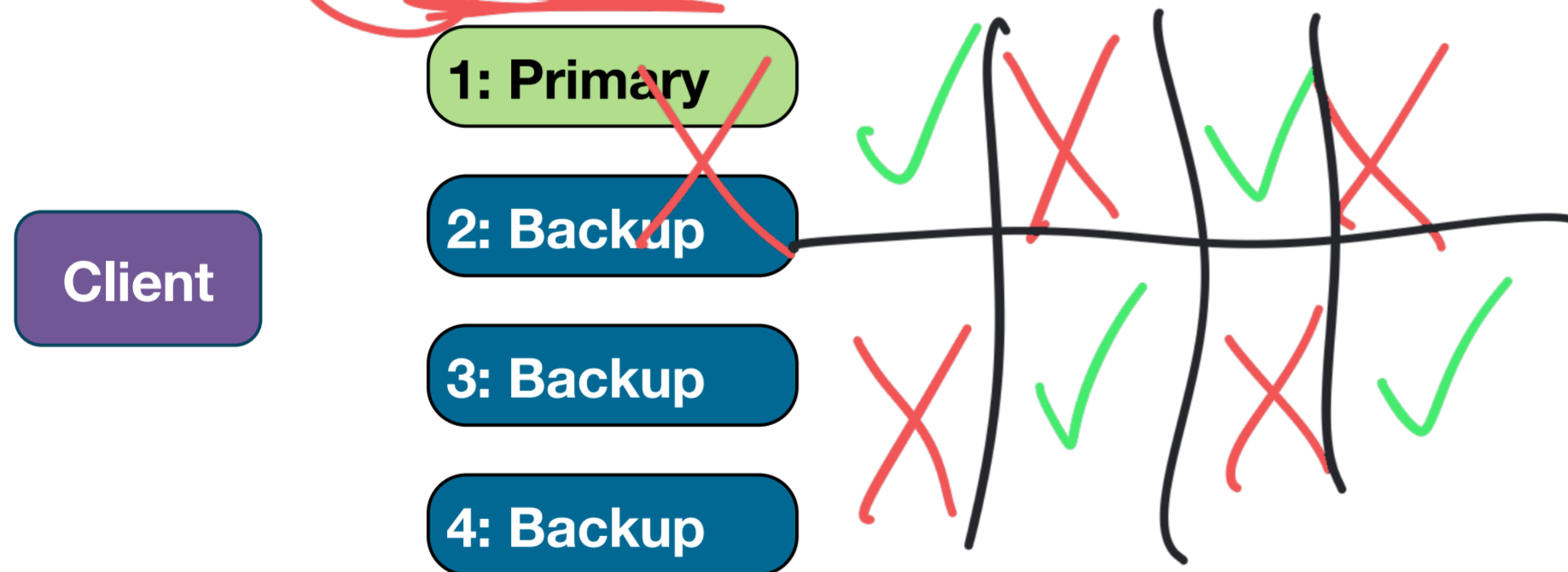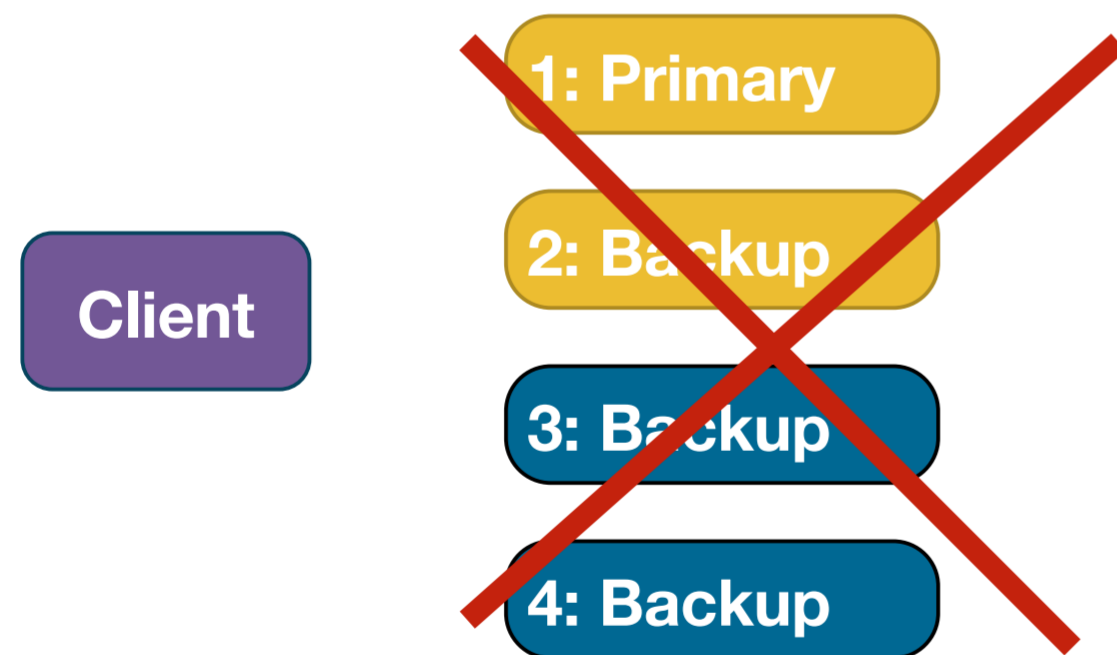
# HANDLING FAILURES

- F = number of nodes which can crash at one time
- # of nodes needed must depend on f!

**f=2**, **f+2 = 4 replicas**

Client

1: Primary

2: Backup

3: Backup

4: Backup

Can't resync state if failure
"flip flops" between **2** nodes!

f=1, f+2 = 3
replicas

Log

1: Primary

Log

2: Backup

Log

3: Backup

Fixed for f=2? No!

# HANDLING FAILURES

f=1 → 3

- F = number of nodes which can crash at one time
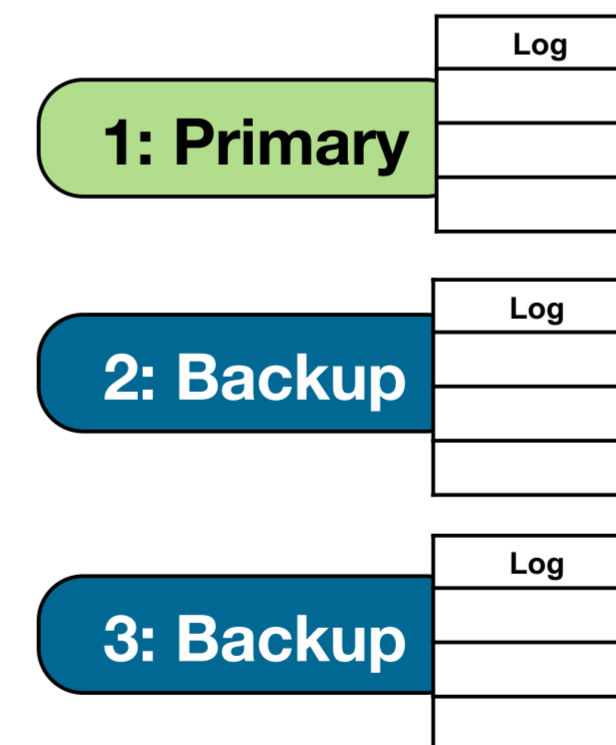- # of nodes needed must depend on f!

f=2, 2f+1 = 5 **replicas**

f=2, f+2 replicas

**Client**

1: Primary

2: Backup

3: Backup

4: Backup
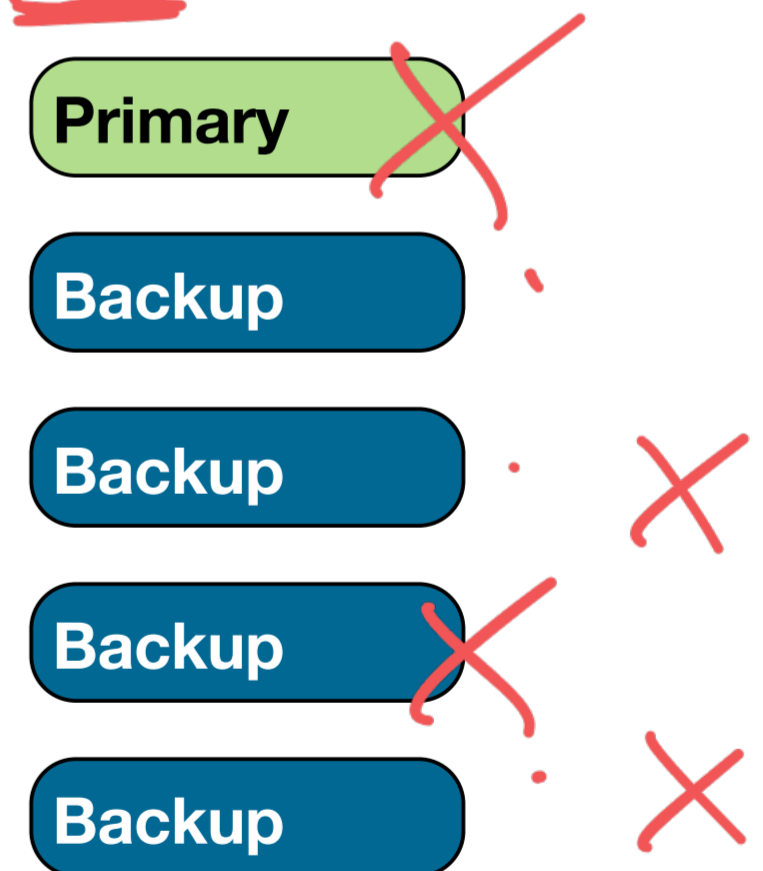
Can't resync state if failure "flip flops" between **2** nodes!

Use **2f+1** replicas!
Insight: Always need a **majority** of nodes to stay alive!

Primary

Backup

Backup

Backup

Backup

# STATE MACHINE REPLICATION OVERVIEW

- Provides a generic **fault tolerance** mechanism
  - Application just needs to have well defined operations and a way to avoid non-determinism

- Primary orders requests into log

- Backups execute log in order

- Log allows out of date replicas to recover

- Need **2f+1** replicas to tolerate **f** failures

- But how do we pick who should be primary…?
  - Use an election algorithm!

**Optional HW 3:** Implement the Election algorithm used by the Raft SMR protocol

# CASE STUDY

- Two important challenges in BlockChain
  - How are any decisions made?
  - How does anything get done?

Prof. Tim Wood & Prof. Roozbeh Haghnazar

# DISTRIBUTED LEDGER TECH



Centralised Ledger

Distributed Ledger

Prof. Tim Wood & Prof. Roozbeh Haghnazar

# DIFFERENT TYPES OF DLT

- Blockchain
- Hashgraph
- DAG
- Holochain
- Tangle
- Radix (Tempo)

Prof. Tim Wood & Prof. Roozbeh Haghnazar

# HASHGRAPH

- It's so fast – 250000 transaction per second (Scalability characteristics in Distributed Systems)
- Being Time-Based and using Gossip protocol for consensus reduces the process and math complexity.
- In the level of security it is evaluating in the banking system level and it means it is a **Byzantine Fault Tolerance** system.
- Controlled Network (Consensus is easier)

Hashgraph

Node in the graph

Time

timestamp

transactions

hash2

hash1

A    B    C    D

Hashgraph Data Structure

Prof. Tim Wood & Prof. Roozbeh Haghnazar

# Tangle (IOTA)

- IOTA is an open-source distributed ledger and cryptocurrency designed for the Internet of things.

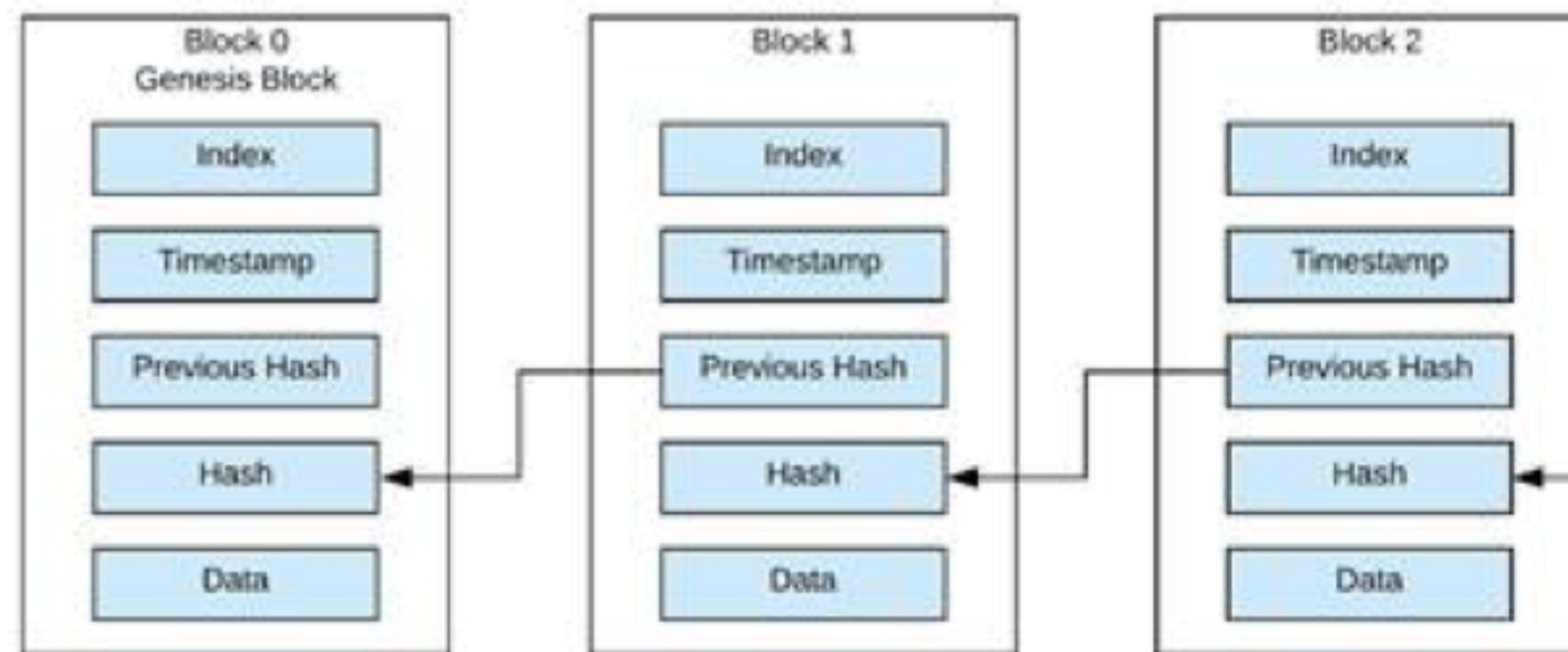- Uses DAG to store transactions on its ledger, motivated by a potentially higher scalability over blockchain based distributed ledgers for nano-Transactions between IOT devices.

- There are categories of participants,
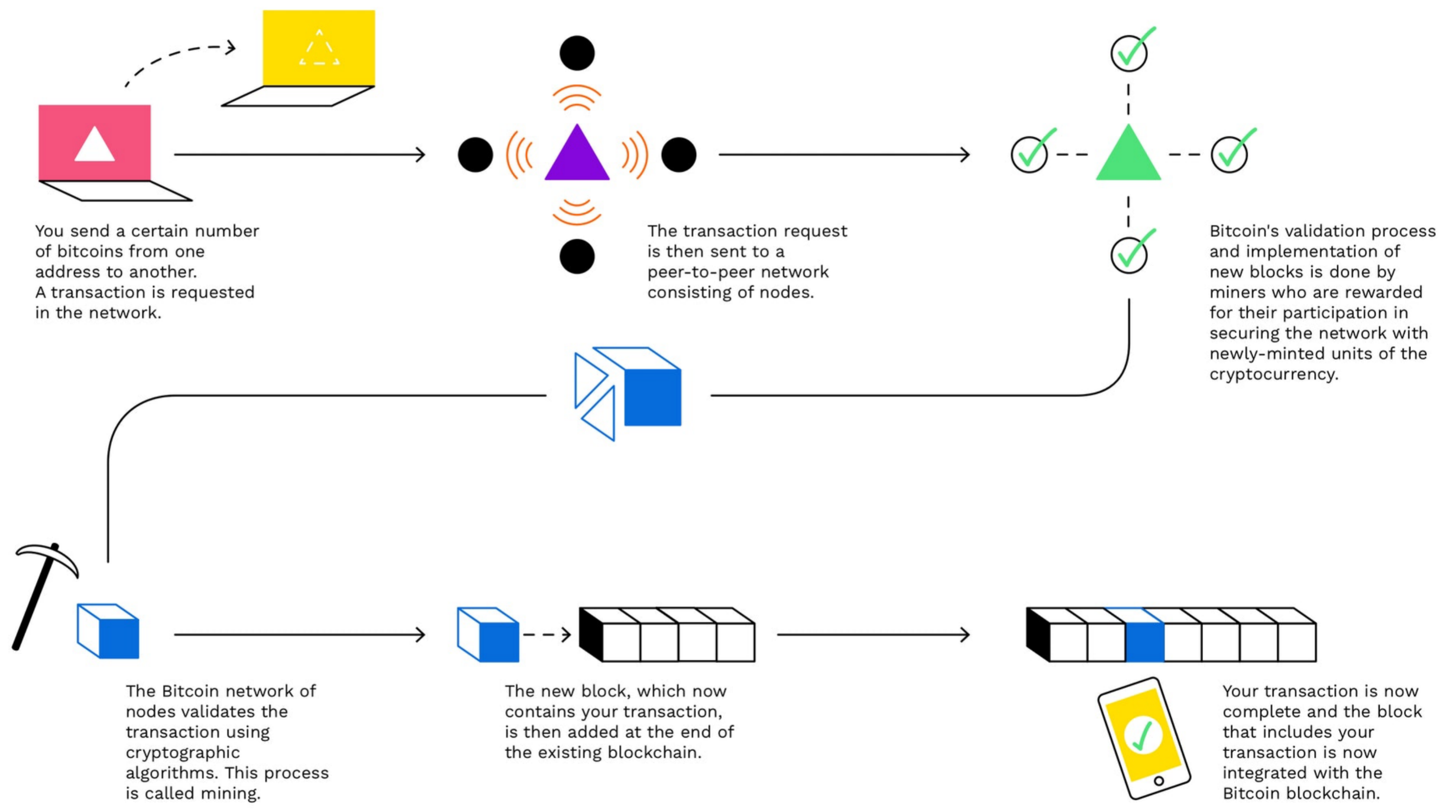  - Transaction creators
  - Transaction verifiers

# BLOCKCHAIN

- Unofficial definition: A blockchain is an unchangeable and sequence of records and transactions which is called **BLOCK**

- The blocks connects to each other with Hash Codes

- Each block contains an index, time stamp, list of transactions, evidence, and **last block hash** (which guarantees the unchangeability of the chain)

You send a certain number of bitcoins from one address to another. A transaction is requested in the network.

The transaction request is then sent to a peer-to-peer network consisting of nodes.

Bitcoin's validation process and implementation of new blocks is done by miners who are rewarded for their participation in securing the network with newly-minted units of the cryptocurrency.

The Bitcoin network of nodes validates the transaction using cryptographic algorithms. This process is called mining.

The new block, which now contains your transaction, is then added at the end of the existing blockchain.

Your transaction is now complete and the block that includes your transaction is now integrated with the Bitcoin blockchain.

Prof. Tim Wood & Prof. Roozbeh Haghnazar

# CONSENSUS IN BLOCKCHAIN

- A consensus mechanism enables the blockchain network to attain reliability and build a level of trust between different nodes, while ensuring security in the environment.
    - Proof of Work (PoW)
    - Proof of Stake (PoS)
    - Delegated Proof of Stake (DPoS)
    - Leased Proof of Stake (LPoS)
    - Direct Acyclic Graph (DAG)
    - Byzantine Fault Tolerance (BFT)
    - Practical Byzantine Fault Tolerance (PBFT)
    - Delegated Byzantine Fault Tolerance (DBFT)
    - Proof of Capacity (PoC)
    - Etc.

Prof. Tim Wood & Prof. Roozbeh Haghnazar