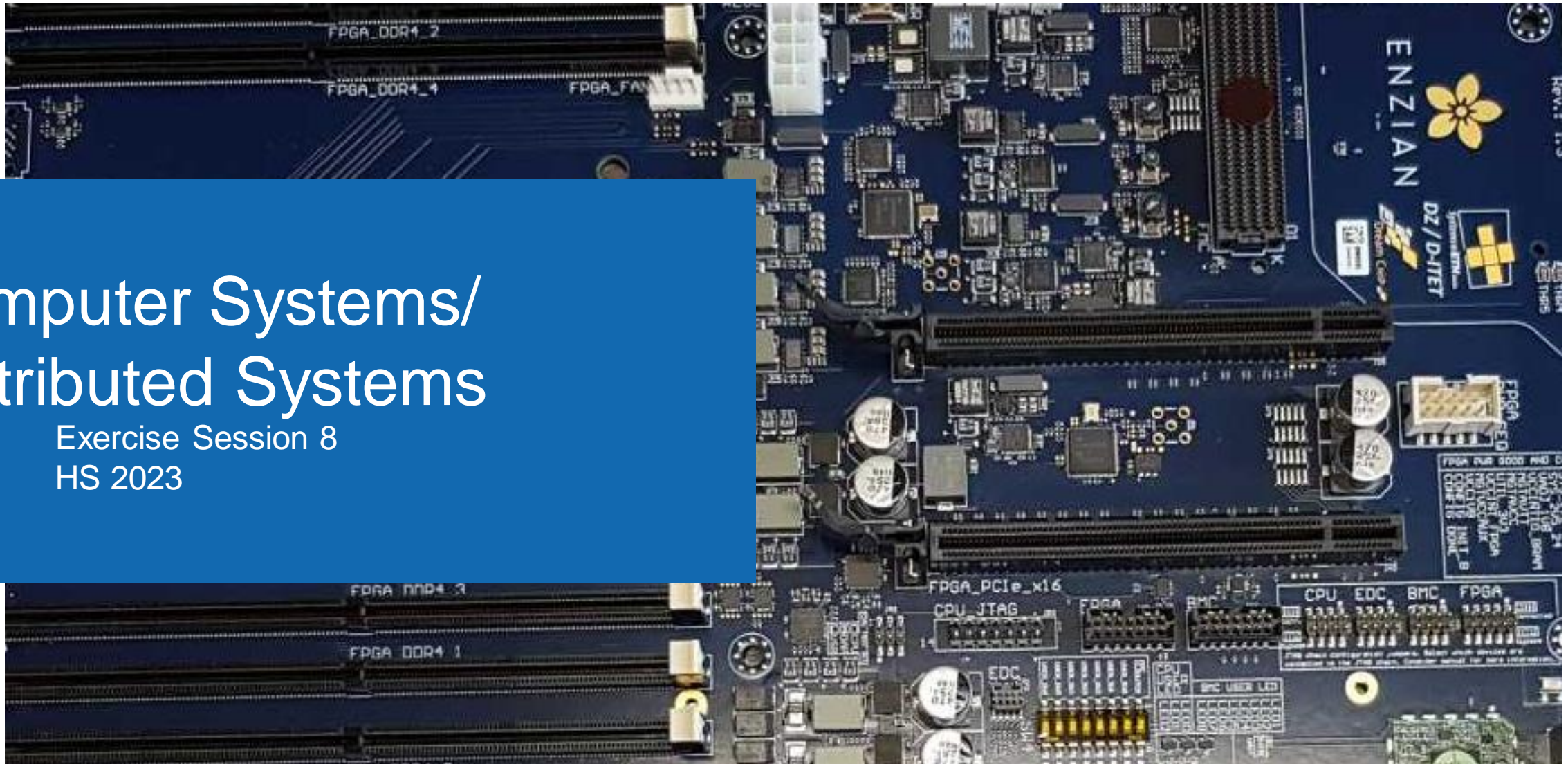




# Computer Systems/ Distributed Systems

Exercise Session 8  
HS 2023





## Program

### 1. Lecture Recap

a) Introduction: Distributed Systems

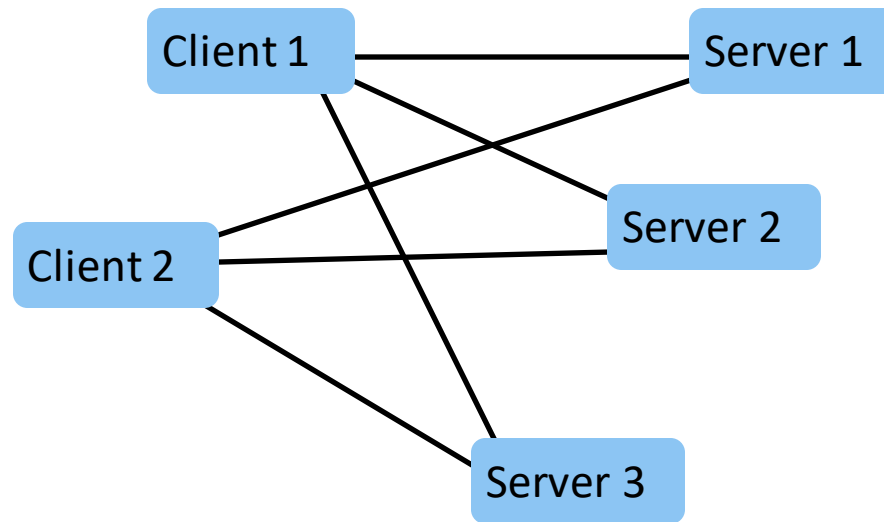
b) Fault Tolerance and Paxos

c) Consensus

### 2. Quiz

### 3. Assignment Preview

## Set-Up



**Node:** single actor in a distributed system

Can be both client or server



## Challenges

- Messages can get lost
- Nodes may crash
- Messages can have varying delays



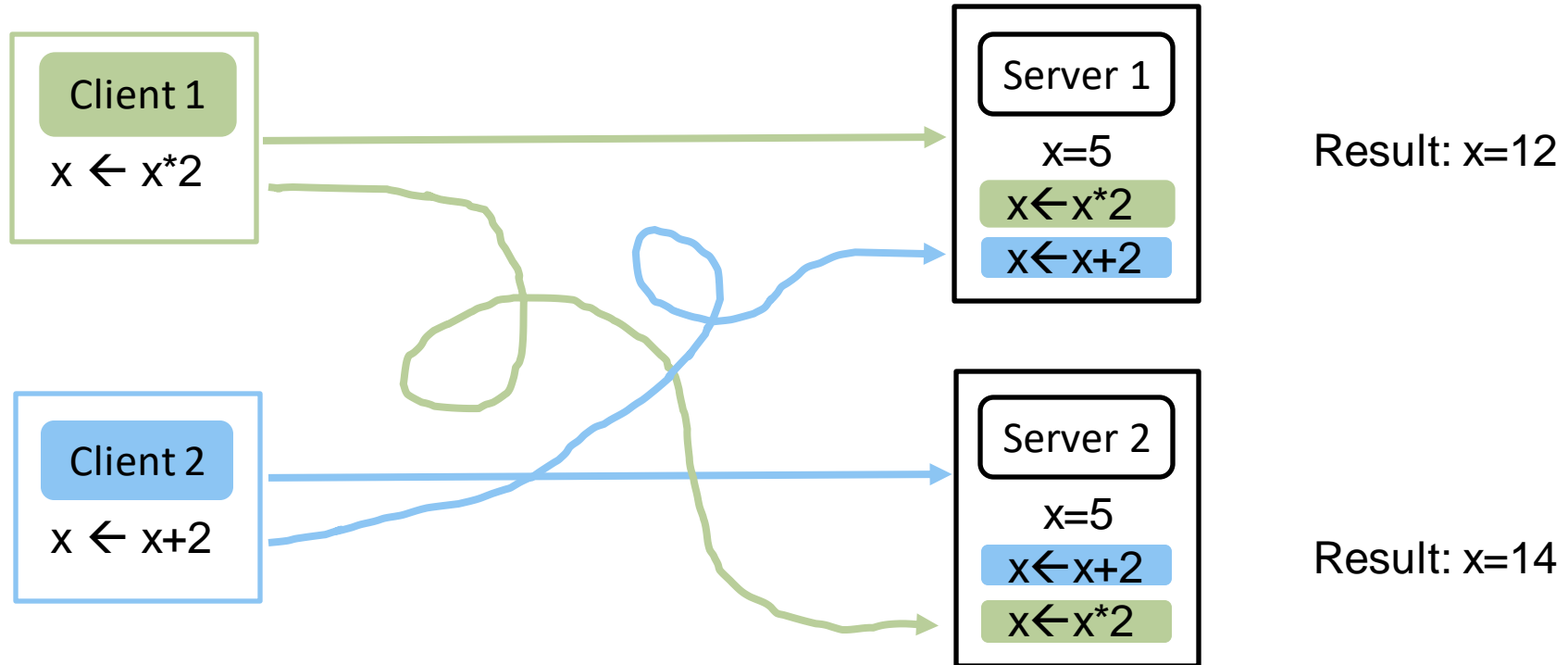
## Challenges

- Messages can get lost
- Nodes may crash
- Messages can have varying delays

## First Goal: State Replication

- All servers execute the *same commands* in the *same order*.

# Why do we want State Replication?





## First Approaches

Server sends acknowledgment message

- Reasonable with one client
- Inconsistent state with multiple clients and servers



## First Approaches

Server sends acknowledgment message

- Reasonable with one client
- Inconsistent state with multiple clients and servers

Serializer – all commands go through one node which orders them

- Single point of failure





## First Approaches

Server sends acknowledgment message

- Reasonable with one client
- Inconsistent state with multiple clients and servers

Serializer – all commands go through one node which orders them

- Single point of failure

Two-Phase Protocol – ask for locks, execute once acquired all locks

- Breaks down if we even have just one node failure
- How to avoid deadlocks?



## Paxos – Main Ideas

### 1. Tickets

- “Weak lock”
- Can be overwritten by later tickets
- Reissuable
- Expiration



## Paxos – Main Ideas

### 1. Tickets

- “Weak lock”
- Can be overwritten by later tickets
- Reissuable
- Expiration

### 2. Require majority

- Ensures only single command gets accepted



## Paxos – Main Ideas

### 1. Tickets

- “Weak lock”
- Can be overwritten by later tickets
- Reissuable
- Expiration

### 2. Require majority

- Ensures only single command gets accepted

### 3. Servers inform clients about their stored command

- Client can switch to supporting this command



## Paxos – Main Ideas

### 1. Tickets

- “Weak lock”
- Can be overwritten by later tickets
- Reissuable
- Expiration

### 2. Require majority

- Ensures only single command gets accepted

### 3. Servers inform clients about their stored command

- Client can switch to supporting this command

Good video with slightly different terminology:

[https://www.youtube.com/watch?v=d7nAGI\\_NZPk](https://www.youtube.com/watch?v=d7nAGI_NZPk)

**Algorithm 7.13 Paxos**

Client (Proposer)

Server (Acceptor)

*Initialization* .....

$c$   $\leftarrow$  command to execute  
 $t = 0$   $\leftarrow$  ticket number to try

$T_{\max} = 0$   $\leftarrow$  largest issued ticket

$C = \perp$   $\leftarrow$  stored command

$T_{\text{store}} = 0$   $\leftarrow$  ticket used to store  $C$

*Phase 1* .....

- 1:  $t = t + 1$
- 2: Ask all servers for ticket  $t$

- 3: if  $t > T_{\max}$  then
- 4:    $T_{\max} = t$
- 5:   Answer with  $\text{ok}(T_{\text{store}}, C)$
- 6: end if

*Phase 2* .....

- 7: if a majority answers ok then
- 8:   Pick  $(T_{\text{store}}, C)$  with largest  $T_{\text{store}}$
- 9:   if  $T_{\text{store}} > 0$  then
- 10:      $c = C$
- 11:   end if
- 12:   Send  $\text{propose}(t, c)$  to same majority
- 13: end if

- 14: if  $t = T_{\max}$  then
- 15:    $C = c$
- 16:    $T_{\text{store}} = t$
- 17:   Answer success
- 18: end if

*Phase 3* .....

- 19: if a majority answers success then
- 20:   Send  $\text{execute}(c)$  to every server
- 21: end if

Clients asks for a specific ticket  $t$ .

If client receives majority of tickets, it proposes a command.

If a majority of servers store the command, the client notifies all servers to execute the command.

Clients can restart Phase 1 at any time.

Server only issues ticket  $t$  if  $t$  is the highest ticket requested so far.

When a server receives a proposal, and the ticket of the client is still valid, the server stores the command and notifies the client.

## Consensus

We want...

1. Agreement:

All (correct) nodes decide on the same value.

## Consensus

We want...

1. Agreement:

All (correct) nodes decide on the same value.

2. Termination:

All (correct) nodes terminate (violated by Paxos).





## Consensus

We want...

1. Agreement:

All (correct) nodes decide on the same value.

2. Termination:

All (correct) nodes terminate (violated by Paxos).

3. Validity:

The decision value is the input value of at least one node.



## Consensus

We want...

1. Agreement:

All (correct) nodes decide on the same value.

2. Termination:

All (correct) nodes terminate (violated by Paxos).

3. Validity:

The decision value is the input value of at least one node.

## Impossibility

Consensus cannot be solved **deterministically** in the asynchronous model!



## Randomized Consensus

Easy cases:

- All inputs are equal (all 0 or 1)
- Almost all input values equal



## Randomized Consensus

Easy cases:

- All inputs are equal (all 0 or 1)
- Almost all input values equal

Otherwise:

- Choose a *random* value locally → Expected time  $O(2^n)$  until all agree (once)



## 8.4 Randomized Consensus

---

### Algorithm 8.15 Randomized Consensus (Ben-Or)

---

```

1:  $v_i \in \{0, 1\}$            $\triangleleft$  input bit
2: round = 1
3: decided = false
4: Broadcast myValue( $v_i$ , round)
5: while true do
    Propose
6:   Wait until a majority of myValue messages of current round arrived
7:   if all messages contain the same value  $v$  then
8:     Broadcast propose( $v$ , round)
9:   else
10:    Broadcast propose( $\perp$ , round)
11:   end if
12:   if decided then
13:     Broadcast myValue( $v_i$ , round+1)
14:     Decide for  $v_i$  and terminate
15:   end if
    Adapt
16:   Wait until a majority of propose messages of current round arrived
17:   if all messages propose the same value  $v$  then
18:      $v_i = v$ 
19:     decided = true
20:   else if there is at least one proposal for  $v$  then
21:      $v_i = v$ 
22:   else
23:     Choose  $v_i$  randomly, with  $Pr[v_i = 0] = Pr[v_i = 1] = 1/2$ 
24:   end if
25:   round = round + 1
26:   Broadcast myValue( $v_i$ , round)
27: end while

```

---

Majority has seen a majority →

At least someone has seen majority →

No majority seen →



## Ben-Or: Consensus Proof

### Validity:

If all nodes start with the same value, then all proposals are for the same value.  
Thus, the algorithm terminated within one round, deciding on the common value.

If some nodes start with 0 and some start with 1, then both outcomes are legal.



## Ben-Or: Consensus Proof

**Agreement: (need to show: if one node decides  $\rightarrow$  all nodes decide on the same value)**

In a single round  $r$ :

- Nodes only decide after having received a proposal.
  - Note, that a proposal required a majority, therefore a proposal in round  $r$  can only occur for one value.
- $\rightarrow$  In any round  $r$ , all nodes decide on at most one identical value.



## Ben-Or: Consensus Proof

**Agreement: (need to show: if one node decides  $\rightarrow$  all nodes decide on the same value)**

In a single round  $r$ :

- Nodes only decide after having received a proposal.
  - Note, that a proposal required a majority, therefore a proposal in round  $r$  can only occur for one value.
- $\rightarrow$  In any round  $r$ , all nodes decide on at most one identical value.

In the first round after a node decided:

- Deciding node received  $> n/2$  proposals.  $\rightarrow$  All nodes received  $\geq 1$  proposal.
- Will adapt their own value to proposal, and broadcast value in round  $r$ .
- As all nodes broadcast the same value, they will propose same value in round  $r+1$ .
- At the latest in round  $r+2$  nodes receive  $> n/2$  proposals.





## Ben-Or: Consensus Proof

### Termination:

Trivial case: all nodes start with the same value

→ Termination after one round.

In the worst case: no node receives all identical majorities, and all repeatedly choose a random value. The probability of all nodes getting the same value is  $2^{-n}$ , thus we expect all nodes to send the same “my value” after  $2^n$  runs.



## Randomized Consensus

Easy cases:

- All inputs are equal (all 0 / all 1)
- Almost all input values equal

Otherwise:

- Choose a *random* value locally → Expected time  $O(2^n)$  until all agree (once)
- Wouldn't it be useful if the nodes could all toss the *same* coin? → Shared Coin



# Shared Coin

## 8.5 Shared Coin

---

**Algorithm 8.22** Shared Coin (code for node  $u$ )

---

- 1: Choose local coin  $c_u = 0$  with probability  $1/n$ , else  $c_u = 1$
  - 2: Broadcast `myCoin( $c_u$ )`
  - 3: Wait for  $n - f$  coins and store them in the local coin set  $C_u$
  - 4: Broadcast `mySet( $C_u$ )`
  - 5: Wait for  $n - f$  coin sets
  - 6: **if** at least one coin is 0 among all coins in the coin sets **then**
  - 7:   return 0
  - 8: **else**
  - 9:   return 1
  - 10: **end if**
-



## Shared Coin

- The algorithm stays exactly the same, except the standard coin flip is replaced by a call to the shared coin algorithm.



## Shared Coin

- The algorithm stays exactly the same, except the standard coin flip is replaced by a call to the shared coin algorithm.
- Proofs for validity and agreement still hold (since it is still the same algorithm).



## Shared Coin

- The algorithm stays exactly the same, except the standard coin flip is replaced by a call to the shared coin algorithm.
- Proofs for validity and agreement still hold (since it is still the same algorithm).
- The proof for termination has to be changed slightly to account for the changed probability that all coins will give the same result.



## Shared Coin

- The algorithm stays exactly the same, except the standard coin flip is replaced by a call to the shared coin algorithm.
- Proofs for validity and agreement still hold (since it is still the same algorithm).
- The proof for termination has to be changed slightly to account for the changed probability that all coins will give the same result.
- Runtime: From exponential down to constant!



## Shared Coin

- The algorithm stays exactly the same, except the standard coin flip is replaced by a call to the shared coin algorithm.
- Proofs for validity and agreement still hold (since it is still the same algorithm).
- The proof for termination has to be changed slightly to account for the changed probability that all coins will give the same result.
- Runtime: From exponential down to constant!
- Can only tolerate  $f < n/3$  crash failures, not  $f < n/2$ .





## Quiz Paxos

1. How does a node in Paxos know if a majority answered with ok?
2. Does the Paxos algorithm in the script achieve state replication?
3. How many nodes could crash so the Paxos still works?
4. Does Paxos solve consensus?



## Quiz Paxos

1. How does a node in Paxos know if a majority answered with ok?  
Each node needs to know the number  $n$  of servers in the system.
2. Does the Paxos algorithm in the script achieve state replication?
3. How many nodes could crash so the Paxos still works?
4. Does Paxos solve consensus?



## Quiz Paxos

1. How does a node in Paxos know if a majority answered with ok?  
Each node needs to know the number  $n$  of servers in the system.
2. Does the Paxos algorithm in the script achieve state replication?  
No, it only shows agreement on a single command, for several commands we would need to restart the system.
3. How many nodes could crash so the Paxos still works?
4. Does Paxos solve consensus?



## Quiz Paxos

1. How does a node in Paxos know if a majority answered with ok?  
Each node needs to know the number  $n$  of servers in the system.
2. Does the Paxos algorithm in the script achieve state replication?  
No, it only shows agreement on a single command, for several commands we would need to restart the system.
3. How many nodes could crash so the Paxos still works?  
Less than  $n/2$
4. Does Paxos solve consensus?



## Quiz Paxos

1. How does a node in Paxos know if a majority answered with ok?  
Each node needs to know the number  $n$  of servers in the system.
2. Does the Paxos algorithm in the script achieve state replication?  
No, it only shows agreement on a single command, for several commands we would need to restart the system.
3. How many nodes could crash so the Paxos still works?  
Less than  $n/2$
4. Does Paxos solve consensus?  
No, termination is not guaranteed.



## More quiz questions (choose the right answer)

1. State replication is trivial for fewer than 3 nodes? T/F
2. In Paxos, a new ticket can only be issued if all previous tickets have been returned. T/F
3. Which is not a property of consensus? Agreement Termination Tolerance Validity
4. A configuration includes all received messages but not the messages in transit. T/F
5. In a synchronous system, a message has a delay of \_\_\_ time units.  $1/n$ / $f-n$ / potentially infinite



## More quiz questions (choose the right answer)

1. State replication is trivial for fewer than 3 nodes? T/**F**
2. In Paxos, a new ticket can only be issued if all previous tickets have been returned. T/**F**
3. Which is not a property of consensus? Agreement Termination **Tolerance** Validity
4. A configuration includes all received messages but not the messages in transit. T/**F**
5. In a synchronous system, a message has a delay of \_\_\_ time units. **1**/n/f-n/ potentially infinite



# Assignment Preview

## 1.1 An Asynchronous Riddle

A hangman summons his 100 prisoners, announcing that they may meet to plan a strategy, but will then be put in isolated cells, with no communication. He explains that he has set up a switch room that contains a single switch. Also, the switch is not connected to anything, but a prisoner entering the room may see whether the switch is on or off (because the switch is up or down). Every once in a while the hangman will let one arbitrary prisoner into the switch room. The prisoner may throw the switch (on to off, or vice versa), or leave the switch unchanged. Nobody but the prisoners will ever enter the switch room. The hangman promises to let any prisoner enter the room from time to time, arbitrarily often. That is, eventually, each prisoner has been in the room at least once, twice, a thousand times or any number you want. At any time, any prisoner may declare “We have all visited the switch room at least once”. If the claim is correct, all prisoners will be released. If the claim is wrong, the hangman will execute his job (on all the prisoners). Which strategy would you choose...

- a) ...if the hangman tells them, that the switch is off at the beginning?
- b) ...if they don't know anything about the initial state of the switch?





# Assignment Preview

## 2.1 Consensus with Edge Failures

In the lecture we only discussed node failures, but we always assumed that edges (links) never fail. Let us now study the opposite case: Assume that all nodes work correctly, but up to  $f$  edges may fail.

Analogously to node failures, edges may fail at any point during the execution. We say that a failed edge does not forward any message anymore, and remains failed until the algorithm terminates. Assume that an edge always simultaneously fails completely, i.e., no message can be exchanged over that edge anymore in either direction.

We assume that the network is initially fully connected, i.e., there is an edge between every pair of nodes. Our goal is to solve consensus in such a way, that *all* nodes know the decision.

- a) What is the smallest  $f$  such that consensus might become impossible? (Which edges fail in the worst-case)
- b) What is the largest  $f$  such that consensus might still be possible? (Which edges fail in the best-case)
- c) Assume that you have a setup which guarantees you that the nodes always remain connected, but possibly many edges might fail. A very simple algorithm for consensus is the following: Every node learns the initial value of all nodes, and then decides locally. How much time might this algorithm require?

Assume that a message takes at most 1 time unit from one node to a direct neighbor.