

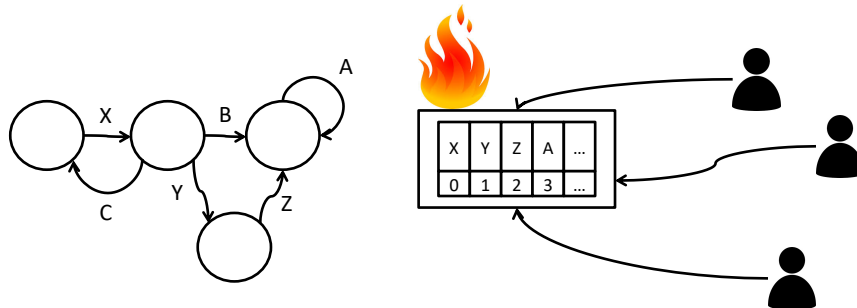
## Paxos

- Distributed consensus protocol
  - Group decision making by majority voting
    - $F+1$  out of  $2F+1$  can make the decision
    - Can handle  $F$  node failures
    - Can handle network failures
    - Can handle network delays (reordering)
  - Once decision is made the decision does not change (write-once register concept)
- Paxos: reaching consensus on ONE case
- Multi-Paxos: extension to multiple cases

1

## Original Intended Use of Paxos

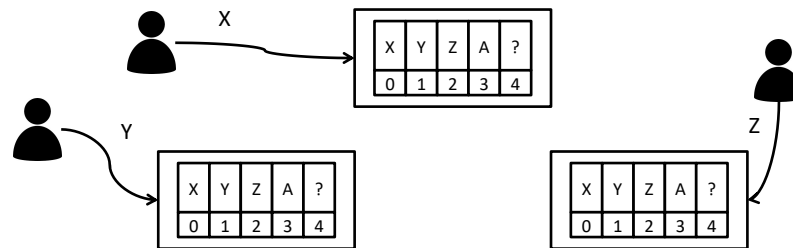
- State machine approach
  - Gets user input then changes state and emits output
  - Record of input can represent current state



2

## Original Intended Use of Paxos

- State machine approach with replication
  - All replicas should be in the same state
  - All inputs should be given in the same order
- Making decisions for one value is not easy



3

## Sources of Paxos Confusion

- The part-time parliament
  - Described in a form of a story of Aegean island of Paxos
  - Annotation from the journal:
    - “The author appears to be an archeologist with only a passing interest in computer science.”
  - Paper submitted in 1990, accepted in 1998
- Paxos made simple, 2001
  - Abstract:
    - “The Paxos algorithm, when presented in plain English, is very simple.”
  - Maybe too simplified

4

## Sources of Paxos Confusion

- Paxos made simple paper

- Clearly described

- Acceptor
- Proposer

Most  
Important  
Core

- Vaguely described

- Leader
- Learner
- Membership change
- Unique proposal number

Everybody invents  
their own mechanisms

5

## Modular Paxos

- Focus on the essence

- Implement only clearly described parts

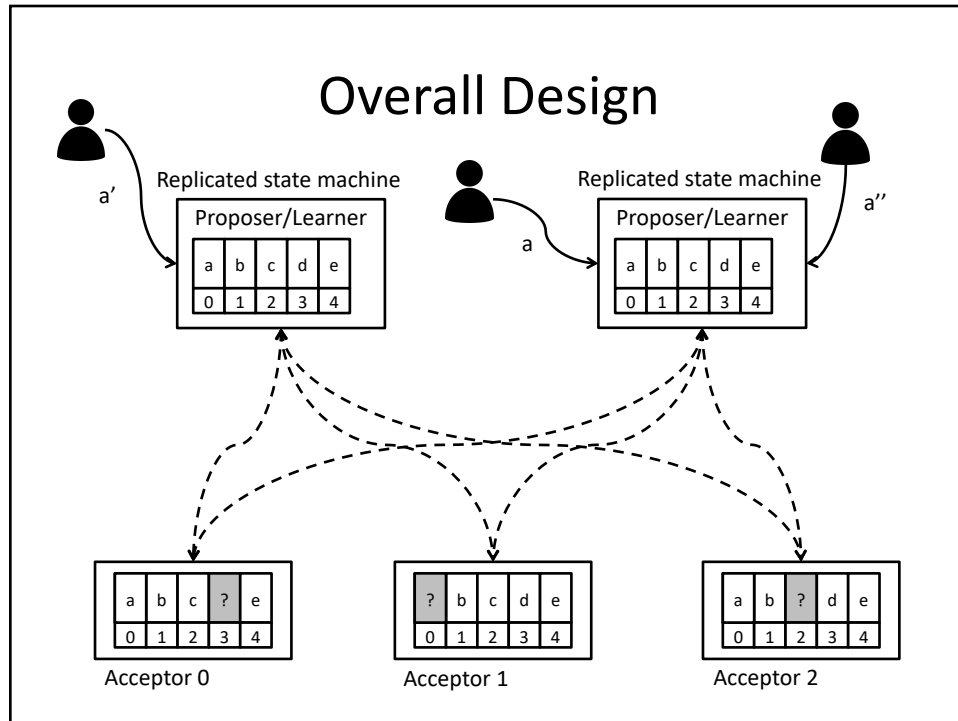
- Acceptor
- Proposer
- Learner implemented as proposer

- Ignore optimizations

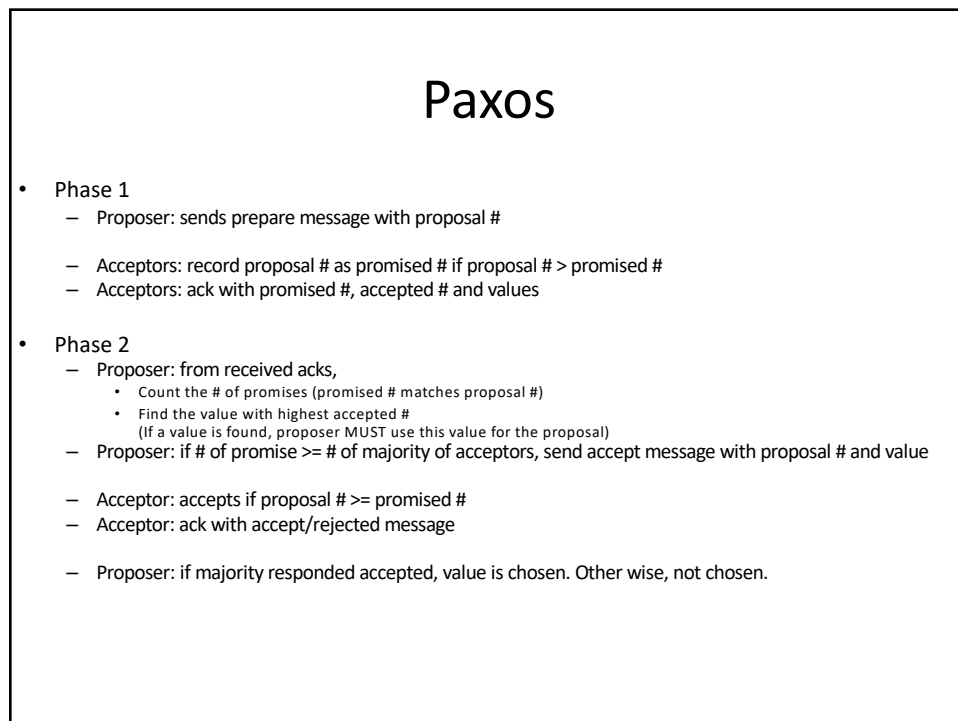
- Performance may be moderate
- Liveness guarantee can be an issue

\* Optimizations can be added later as needed

6



7



8

# Paxos Algorithm

## Proposer/Acceptor

### Proposer

$Pp$ : proposal #;  $Vp$ : Value; ( $Ap$ : Accepted #)

1. Prepare Request
  1. Send ( $Pp$ )
2. Handle prepare ack
  1. If ( $Aa, Va$ ) exists and  $Aa > Ap$  then  $Ap = Aa$  and  $Vp = Va$
  2. If  $Pp == Pa$  then  $Promised\_cnt++$

3. Accept request
  1. If  $Promised\_cnt \geq majority$  then  
Send ( $Pp, Vp$ )

4. Handle accept ack
  1. If  $Accepted$  then  $Accepted\_cnt++$

5. Decide success/failure
  1. If  $Accepted\_cnt \geq Majority$  then **Success**  
else **Failure**

### Acceptor (total= $2F+1$ ; majority = $F+1$ )

$Pa$ : Promised #;  $Aa$ : Accepted #;  $Va$ : Value

1. Handle prepare
  1. If  $Pp > Pa$  then  $Pa = Pp$
  2. Ack with ( $Pa, (Aa, Va)$ )
2. Handle accept
  1. If  $Pp \geq Pa$  then  $Aa = Pp, Va = Vp$
  2. Ack with (**Accepted/Rejected**)

9

# Implementation

### Proposer client (This will become Write)

```
// 1. Prepare request and get ack
For (number of acceptors) {
    Send (Pp)
    RecvAck(ack)
}
// 2. Handle prepare ack
Foreach (acks) {
    If (Aa, Va) exists and Aa > Ap then Ap = Aa and Vp = Va
    If Pp == Pa then Promised_cnt++
}
// 3. Accept request and get ack
If Promised_cnt >= majority then
for (number of acceptors) {
    Send (Pp, Vp)
    RecvAck(ack)
}
else return;
// 4. Handle accept ack
Foreach (acks) {
    If Accepted then Accepted_cnt++
}
// 5. Decide success/failure
If Accepted_cnt >= Majority then Success
else Failure
```

### Acceptor server: main

```
While (true) {
    AcceptConnections(socket);
    Thread_create(thread_body, args);
}
```

### Acceptor server: thread\_body

```
While (Receive(message)) {
    Switch (message.type) {
    Case Prepare:
        If Pp > Pa then Pa = Pp
        Ack with Send(Pa, (Aa, Va))
    Case Accept:
        If Pp >= Pa then Aa = Pp, Va = Vp
        Ack with Send(Accepted/Rejected)
    }
}
```

10

## Paxos Algorithm

### Learner

**Learner (Almost same as proposer)**

**$Pp$ : proposal #; ( $Vp$ : Value;  $Ap$ : Accepted #)**

1. Prepare Request
  1. Send ( $Pp$ )
2. Handle prepare ack
  1. If ( $Aa, Va$ ) exists and  $Aa > Ap$  then  $Ap = Aa$  and  $Vp = Va$
  2. If  $Pp == Pa$  then  $Promised\_cnt++$
3. Accept request
  1. If  $Promised\_cnt \geq majority$  then Send ( $Pp, Vp$ )
4. Handle accept ack
  1. If **Accepted** then  $Accepted\_cnt++$
5. Decide success/failure
  1. If  $Accepted\_cnt \geq Majority$  then **Success** else **Failure**

**Acceptor (total=  $2F+1$ ; majority =  $F+1$ )**

**$Pa$ : Promised #;  $Aa$ : Accepted #;  $Va$ : Value**

1. Handle prepare
  1. If  $Pp > Pa$  then  $Pa = Pp$
  2. Ack with ( $Pa, (Aa, Va)$ )

- 2.5 Learn value
    1. For all ( $Aa, Va$ ) pairs find the majority pair
    2. If found then return the majority **Value** else read **Failed**

and  $Vp$  contains a value

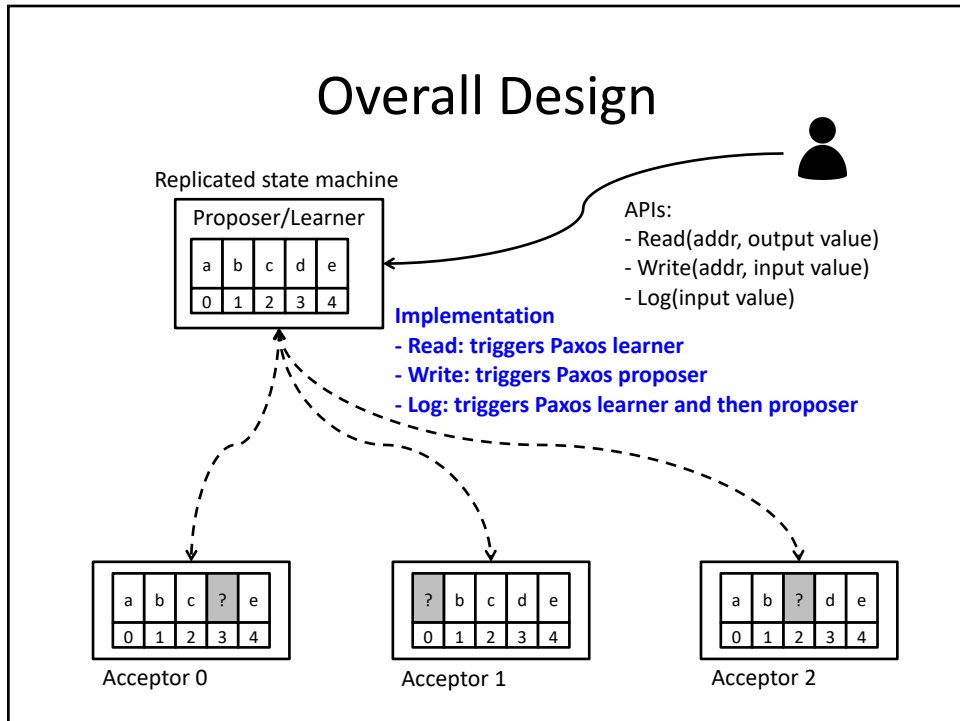
1. If  $Pp \geq Pa$  then  $Aa = Pp, Va = Vp$
2. Ack with **Accepted/Rejected**

11

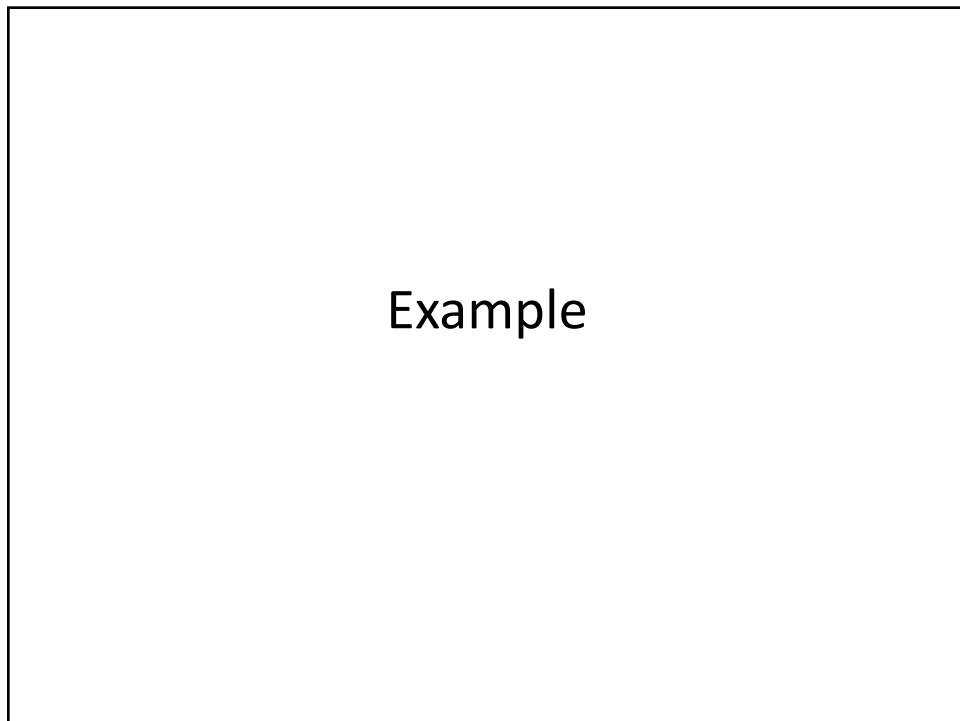
## Extending to Multi-Paxos

- Create an array of acceptor states
- Add an address to every message
  - Send( $Pp$ )                   -> Send(**addr**,  $Pp$ )
  - Send( $Pp, Vp$ )               -> Send(**addr**,  $Pp, Vp$ )
  - Ack( $Pa, Aa, Va$ )           -> Ack(**addr**,  $Pa, Aa, Va$ )
  - Ack(Accept/Reject) -> Ack(**addr**, Accept/Reject)
- Add log order writing

12



13



14

## Paxos

- Proposer ( $P_p = 3, V_p = A$ )
- Acceptor 1 ( $P_a = -1, A_a = -1, V_a = \text{Nil}$ )
- Acceptor 2 ( $P_a = -1, A_a = -1, V_a = \text{Nil}$ )
- Acceptor 3 ( $P_a = -1, A_a = -1, V_a = \text{Nil}$ )

15

## Paxos

- Proposer ( $P_p = 3, V_p = A$ )  
Prepare( $P_p=3$ )
- Acceptor 1 ( $P_a = -1, A_a = -1, V_a = \text{Nil}$ )
- Acceptor 2 ( $P_a = -1, A_a = -1, V_a = \text{Nil}$ )
- Acceptor 3 ( $P_a = -1, A_a = -1, V_a = \text{Nil}$ )

16



## Paxos

- Proposer ( $P_p = 3, V_p = A$ )
- Acceptor 1 ( $P_a = 3, A_a = -1, V_a = Nil$ )
- Acceptor 2 ( $P_a = 3, A_a = -1, V_a = Nil$ )
- Acceptor 3 ( $P_a = 3, A_a = -1, V_a = Nil$ )

17

## Paxos

- Proposer ( $P_p = 3, V_p = A$ )
- Acceptor 1 ( $P_a = 3, A_a = -1, V_a = Nil$ )  
Ack( $P_a=3, A_a=-1, V_a=Nil$ )
- Acceptor 2 ( $P_a = 3, A_a = -1, V_a = Nil$ )  
Ack( $P_a=3, A_a=-1, V_a=Nil$ )
- Acceptor 3 ( $P_a = 3, A_a = -1, V_a = Nil$ )  
Ack( $P_a=3, A_a=-1, V_a=Nil$ )

18

## Paxos

- Proposer ( $P_p = 3, V_p = A$ )  
Accept( $P_p=3, V_p=A$ )
- Acceptor 1 ( $P_a = 3, A_a = -1, V_a = \text{Nil}$ )
- Acceptor 2 ( $P_a = 3, A_a = -1, V_a = \text{Nil}$ )
- Acceptor 3 ( $P_a = 3, A_a = -1, V_a = \text{Nil}$ )

19

## Paxos

- Proposer ( $P_p = 3, V_p = A$ )
- Acceptor 1 ( $P_a = 3, A_a = 3, V_a = A$ )
- Acceptor 2 ( $P_a = 3, A_a = 3, V_a = A$ )
- Acceptor 3 ( $P_a = 3, A_a = 3, V_a = A$ )

20

## Paxos

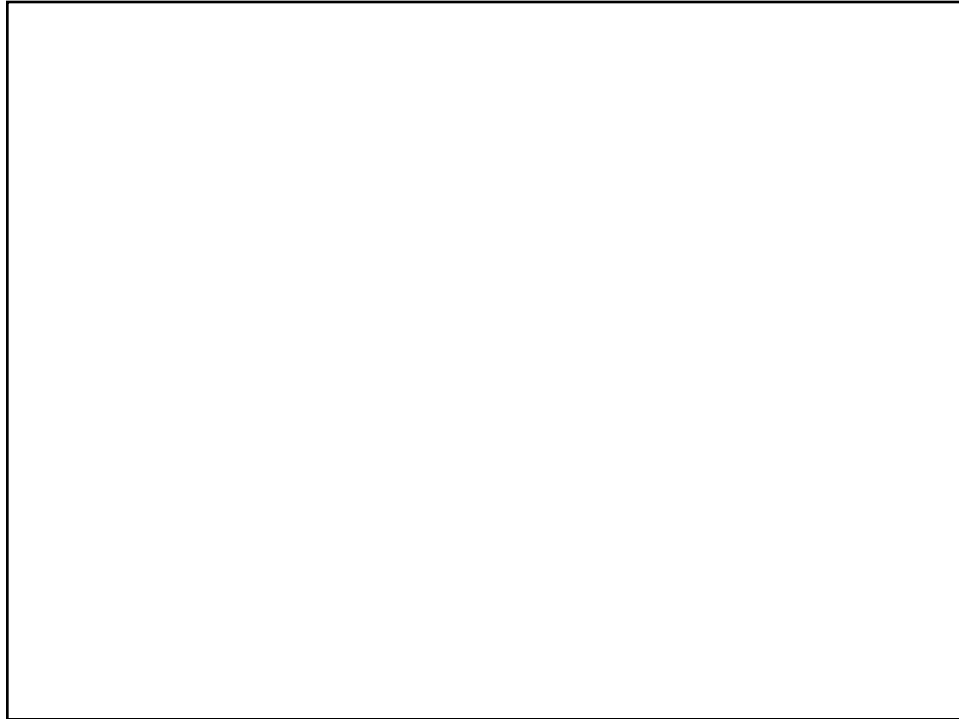
- Proposer ( $P_p = 3, V_p = A$ )
- Acceptor 1 ( $P_a = 3, A_a = 3, V_a = A$ )  
Ack(Accepted)
- Acceptor 2 ( $P_a = 3, A_a = 3, V_a = A$ )  
Ack(Accepted)
- Acceptor 3 ( $P_a = 3, A_a = 3, V_a = A$ )  
Ack(Accepted)

21

## Paxos

- Proposer ( $P_p = 3, V_p = A$ )  
Value Chosen!
- Acceptor 1 ( $P_a = 3, A_a = 3, V_a = A$ )
- Acceptor 2 ( $P_a = 3, A_a = 3, V_a = A$ )
- Acceptor 3 ( $P_a = 3, A_a = 3, V_a = A$ )

22



23

## Paxos

- Proposer 1 ( $P_p = 3, V_p = A$ )
- Acceptor 1 ( $P_a = -1, A_a = -1, V_a = \text{Nil}$ )
- Acceptor 2 ( $P_a = -1, A_a = -1, V_a = \text{Nil}$ )
- Acceptor 3 ( $P_a = -1, A_a = -1, V_a = \text{Nil}$ )

24

## Paxos

- Proposer 1 ( $P_p = 3, V_p = A$ )  
Prepare( $P_p=3$ )
- Acceptor 1 ( $P_a = -1, A_a = -1, V_a = Nil$ )
- Acceptor 2 ( $P_a = -1, A_a = -1, V_a = Nil$ )  
Network down
- Acceptor 3 ( $P_a = -1, A_a = -1, V_a = Nil$ )

25

## Paxos

- Proposer 1 ( $P_p = 3, V_p = A$ )
- Acceptor 1 ( $P_a = 3, A_a = -1, V_a = Nil$ )
- Acceptor 2 ( $P_a = -1, A_a = -1, V_a = Nil$ )  
Network down
- Acceptor 3 ( $P_a = 3, A_a = -1, V_a = Nil$ )

26

## Paxos

- Proposer 1 ( $P_p = 3, V_p = A$ )
- Acceptor 1 ( $P_a = 3, A_a = -1, V_a = Nil$ )  
Ack( $P_a=3, A_a=-1, V_a=Nil$ )
- Acceptor 2 ( $P_a = -1, A_a = -1, V_a = Nil$ )  
Network down
- Acceptor 3 ( $P_a = 3, A_a = -1, V_a = Nil$ )  
Ack( $P_a=3, A_a=-1, V_a=Nil$ )

27

## Paxos

- Proposer 1 ( $P_p = 3, V_p = A$ )  
Accept( $P_p=3, V_p=A$ )
- Acceptor 1 ( $P_a = 3, A_a = -1, V_a = Nil$ )  
Network down
- Acceptor 2 ( $P_a = -1, A_a = -1, V_a = Nil$ )  
Network down
- Acceptor 3 ( $P_a = 3, A_a = -1, V_a = Nil$ )

28

## Paxos

- Proposer 1 ( $P_p = 3, V_p = A$ )
- Acceptor 1 ( $P_a = 3, A_a = -1, V_a = \text{Nil}$ )  
Network down
- Acceptor 2 ( $P_a = -1, A_a = -1, V_a = \text{Nil}$ )  
Network down
- Acceptor 3 ( $P_a = 3, A_a = 3, V_a = A$ )

29

## Paxos

- Proposer 1 ( $P_p = 3, V_p = A$ )
- Acceptor 1 ( $P_a = 3, A_a = -1, V_a = \text{Nil}$ )  
Network down
- Acceptor 2 ( $P_a = -1, A_a = -1, V_a = \text{Nil}$ )  
Network down
- Acceptor 3 ( $P_a = 3, A_a = 3, V_a = A$ )  
Ack(Accepted)

30

## Paxos

- Proposer 1 ( $P_p = 3$ ,  $V_p = A$ )  
Failed!
- Acceptor 1 ( $P_a = 3$ ,  $A_a = -1$ ,  $V_a = \text{Nil}$ )  
Network down
- Acceptor 2 ( $P_a = -1$ ,  $A_a = -1$ ,  $V_a = \text{Nil}$ )  
Network down
- Acceptor 3 ( $P_a = 3$ ,  $A_a = 3$ ,  $V_a = A$ )

31

## Paxos

- Proposer 2 ( $P_p = 4$ ,  $V_p = B$ )
- Acceptor 1 ( $P_a = 3$ ,  $A_a = -1$ ,  $V_a = \text{Nil}$ )
- Acceptor 2 ( $P_a = -1$ ,  $A_a = -1$ ,  $V_a = \text{Nil}$ )
- Acceptor 3 ( $P_a = 3$ ,  $A_a = 3$ ,  $V_a = A$ )  
Network down

32



## Paxos

- Proposer 2 ( $P_p = 4$ ,  $V_p = B$ )  
Prepare( $P_p=4$ )
- Acceptor 1 ( $P_a = 3$ ,  $A_a = -1$ ,  $V_a = Nil$ )
- Acceptor 2 ( $P_a = -1$ ,  $A_a = -1$ ,  $V_a = Nil$ )
- Acceptor 3 ( $P_a = 3$ ,  $A_a = 3$ ,  $V_a = A$ )  
Network down

33

## Paxos

- Proposer 2 ( $P_p = 4$ ,  $V_p = B$ )
- Acceptor 1 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = Nil$ )
- Acceptor 2 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = Nil$ )
- Acceptor 3 ( $P_a = 3$ ,  $A_a = 3$ ,  $V_a = A$ )  
Network down

34

## Paxos

- Proposer 2 ( $P_p = 4$ ,  $V_p = B$ )
- Acceptor 1 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = Nil$ )  
Ack( $P_a=4$ ,  $A_a=-1$ ,  $V_a=Nil$ )
- Acceptor 2 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = Nil$ )  
Ack( $P_a=4$ ,  $A_a=-1$ ,  $V_a=Nil$ )
- Acceptor 3 ( $P_a = 3$ ,  $A_a = 3$ ,  $V_a = A$ )  
Network down

35

## Paxos

- Proposer 2 ( $P_p = 4$ ,  $V_p = B$ )  
Accept( $P_p=4$ ,  $V_p=B$ )
- Acceptor 1 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = Nil$ )  
Network down
- Acceptor 2 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = Nil$ )
- Acceptor 3 ( $P_a = 3$ ,  $A_a = 3$ ,  $V_a = A$ )  
Network down

36

## Paxos

- Proposer 2 ( $P_p = 4$ ,  $V_p = B$ )
- Acceptor 1 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = \text{Nil}$ )  
Network down
- Acceptor 2 ( $P_a = 4$ ,  $A_a = 4$ ,  $V_a = B$ )
- Acceptor 3 ( $P_a = 3$ ,  $A_a = 3$ ,  $V_a = A$ )  
Network down

37

## Paxos

- Proposer 2 ( $P_p = 4$ ,  $V_p = B$ )
- Acceptor 1 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = \text{Nil}$ )  
Network down
- Acceptor 2 ( $P_a = 4$ ,  $A_a = 4$ ,  $V_a = B$ )  
Ack(Accepted)
- Acceptor 3 ( $P_a = 3$ ,  $A_a = 3$ ,  $V_a = A$ )  
Network down

38

## Paxos

- Proposer 2 ( $P_p = 4$ ,  $V_p = B$ )  
Failed!
- Acceptor 1 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = \text{Nil}$ )  
Network down
- Acceptor 2 ( $P_a = 4$ ,  $A_a = 4$ ,  $V_a = B$ )
- Acceptor 3 ( $P_a = 3$ ,  $A_a = 3$ ,  $V_a = A$ )  
Network down

39

## Paxos

- Proposer 3 ( $P_p = 5$ ,  $V_p = C$ )
- Acceptor 1 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = \text{Nil}$ )  
Network down
- Acceptor 2 ( $P_a = 4$ ,  $A_a = 4$ ,  $V_a = B$ )
- Acceptor 3 ( $P_a = 3$ ,  $A_a = 3$ ,  $V_a = A$ )

40

## Paxos

- Proposer 3 ( $P_p = 5$ ,  $V_p = C$ )  
Prepare( $P_p=5$ )
- Acceptor 1 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = \text{Nil}$ )  
Network down
- Acceptor 2 ( $P_a = 4$ ,  $A_a = 4$ ,  $V_a = B$ )
- Acceptor 3 ( $P_a = 3$ ,  $A_a = 3$ ,  $V_a = A$ )

41

## Paxos

- Proposer 3 ( $P_p = 5$ ,  $V_p = C$ )
- Acceptor 1 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = \text{Nil}$ )  
Network down
- Acceptor 2 ( $P_a = 5$ ,  $A_a = 4$ ,  $V_a = B$ )
- Acceptor 3 ( $P_a = 5$ ,  $A_a = 3$ ,  $V_a = A$ )

42

## Paxos

- Proposer 3 ( $P_p = 5$ ,  $V_p = C$ )
- Acceptor 1 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = \text{Nil}$ )  
 Network down
- Acceptor 2 ( $P_a = 5$ ,  $A_a = 4$ ,  $V_a = B$ )  
 Ack( $P_a=5$ ,  $A_a=4$ ,  $V_a=B$ )
- Acceptor 3 ( $P_a = 5$ ,  $A_a = 3$ ,  $V_a = A$ )  
 Ack( $P_a=5$ ,  $A_a=3$ ,  $V_a=A$ )

43

## Paxos

- Proposer 3 ( $P_p = 5$ ,  $V_p = B$ )
- Acceptor 1 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = \text{Nil}$ )  
 Network down
- Acceptor 2 ( $P_a = 5$ ,  $A_a = 4$ ,  $V_a = B$ )
- Acceptor 3 ( $P_a = 5$ ,  $A_a = 3$ ,  $V_a = A$ )

44

## Paxos

- Proposer 3 ( $P_p = 5$ ,  $V_p = B$ )  
Accept( $P_p=5$ ,  $V_p=B$ )
- Acceptor 1 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = Nil$ )  
Network down
- Acceptor 2 ( $P_a = 5$ ,  $A_a = 4$ ,  $V_a = B$ )
- Acceptor 3 ( $P_a = 5$ ,  $A_a = 3$ ,  $V_a = A$ )

45

## Paxos

- Proposer 3 ( $P_p = 5$ ,  $V_p = B$ )
- Acceptor 1 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = Nil$ )  
Network down
- Acceptor 2 ( $P_a = 5$ ,  $A_a = 5$ ,  $V_a = B$ )  
Ack(Accepted)
- Acceptor 3 ( $P_a = 5$ ,  $A_a = 5$ ,  $V_a = B$ )  
Ack(Accepted)

46

## Paxos

- Proposer 3 ( $P_p = 5$ ,  $V_p = B$ )  
Value Chosen!
- Acceptor 1 ( $P_a = 4$ ,  $A_a = -1$ ,  $V_a = \text{Nil}$ )  
Network down
- Acceptor 2 ( $P_a = 5$ ,  $A_a = 5$ ,  $V_a = B$ )
- Acceptor 3 ( $P_a = 5$ ,  $A_a = 5$ ,  $V_a = B$ )