# COEN6731 Distributed Software Systems

## Week 2: Coordination, Agreement, and Paxos

Gengrui (Edward) Zhang, PhD
Web: gengruizhang.com

# Today's outline

The consensus problem

Network assumptions

Failure assumptions

Paxos

# The consensus problem

Let's go to the beach!

Let's go get some food!

Let's go see a movie!

# The consensus problem

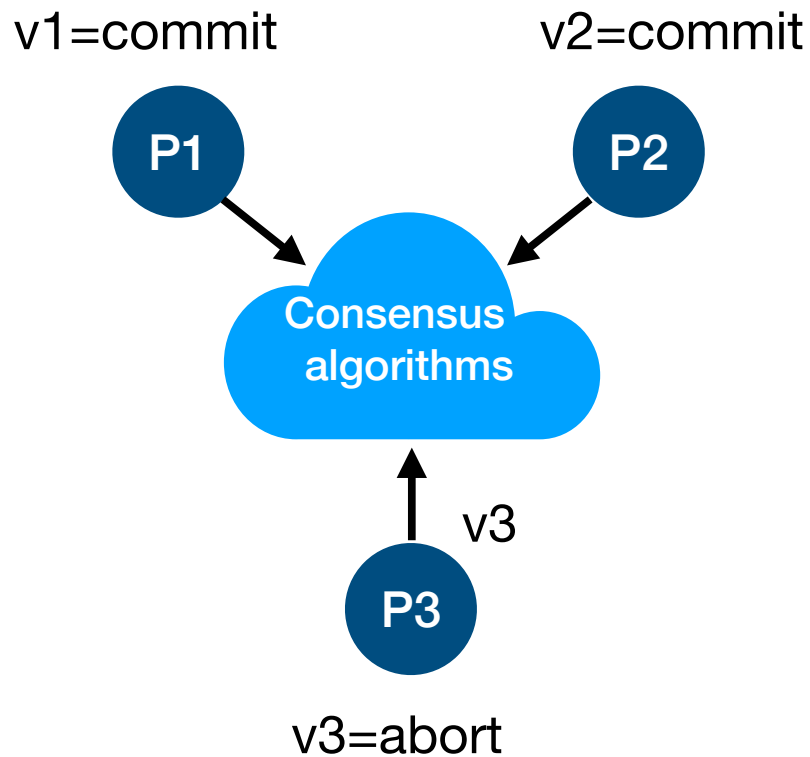Let's go to the beach!

Let's go get some food!

Let's go see a movie!
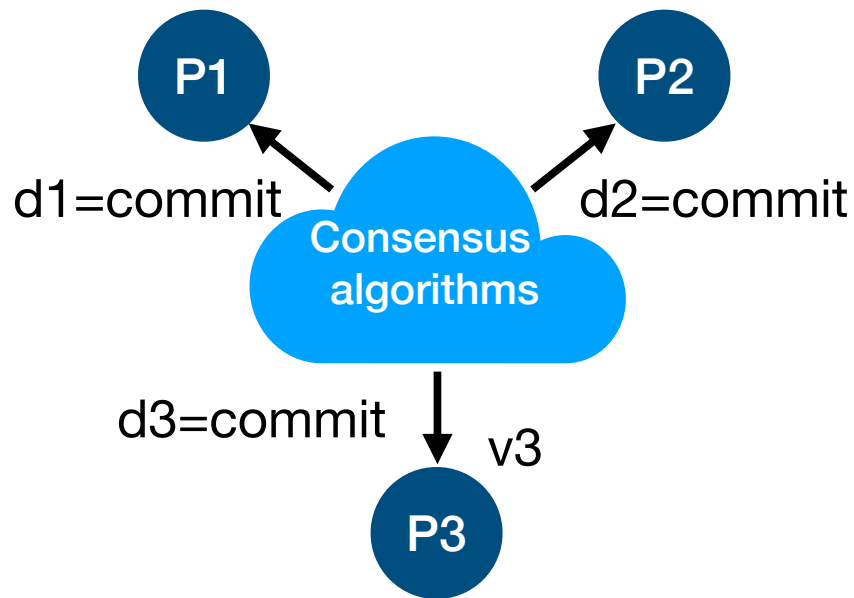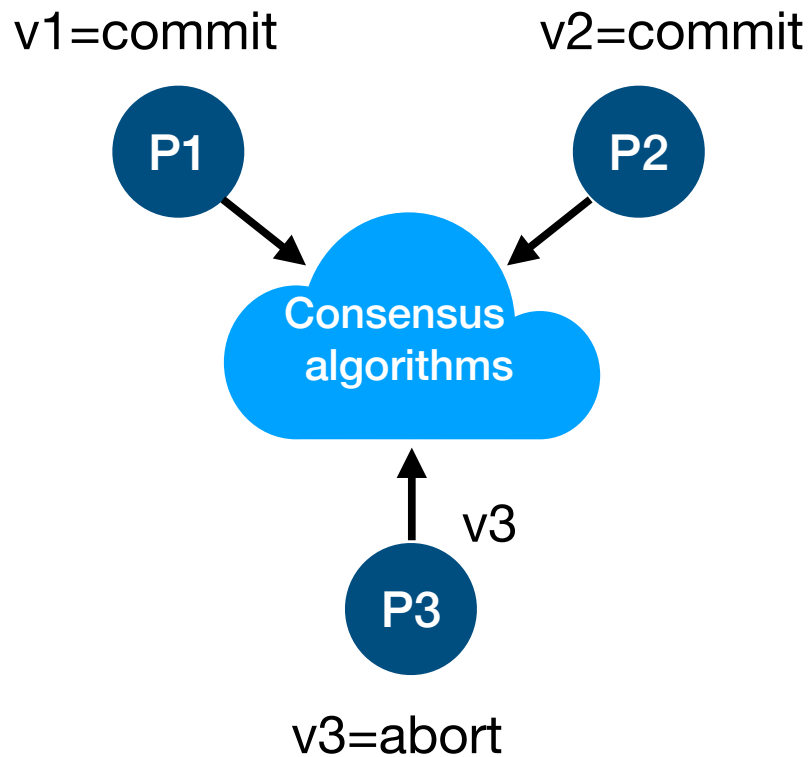
What's important in reaching agreement?

1. Agree on the activities
2. Agree on the order of activities

**The happened-before relation of activities**

# Consensus in distributed systems

v1=commit          v2=commit



v3

v3=abort

# Consensus in distributed systems

v1=commit    v2=commit

P1    P2

Consensus algorithms

v3

P3

v3=abort

d1=commit    P1    P2    d2=commit

Consensus algorithms

d3=commit    v3

P3

# Formally, the consensus problem

- To reach consensus, every process $p_i$ begins in the **undecided** state and **proposes** a single value $v_i$, drawn from a set $D$ $(i = 1, 2, \ldots, N)$.

- Processes communicate with one another, exchanging values.

- Each process then sets the value of a **decision variable**, $d_i$.

- After that, each process enters the **decided** state, where $d_i$ $(i = 1, 2, \ldots, N)$ do not change

# Formally, the consensus problem

- To reach consensus, every process $p_i$ begins in the **undecided** state and **proposes** a single value $v_i$, drawn from a set $D$ ($i = 1, 2, \ldots, N$).

- Processes communicate with one another, exchanging values.

- Each process then sets the value of a **decision variable**, $d_i$.

- After that, each process enters the **decided** state, where $d_i$ ($i = 1, 2, \ldots, N$) do not change

In short, all correct processes commit the **same value** in the **same order**

# Today's outline

# System model: network synchrony
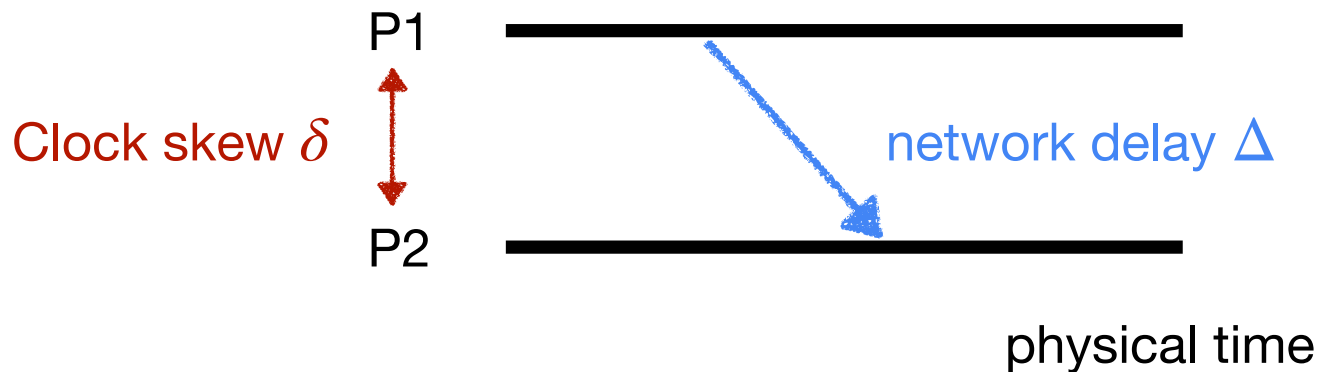
- Synchronous

- Asynchronous

- Partially synchronous

P1

**Clock skew $\delta$**     **network delay $\Delta$**

P2

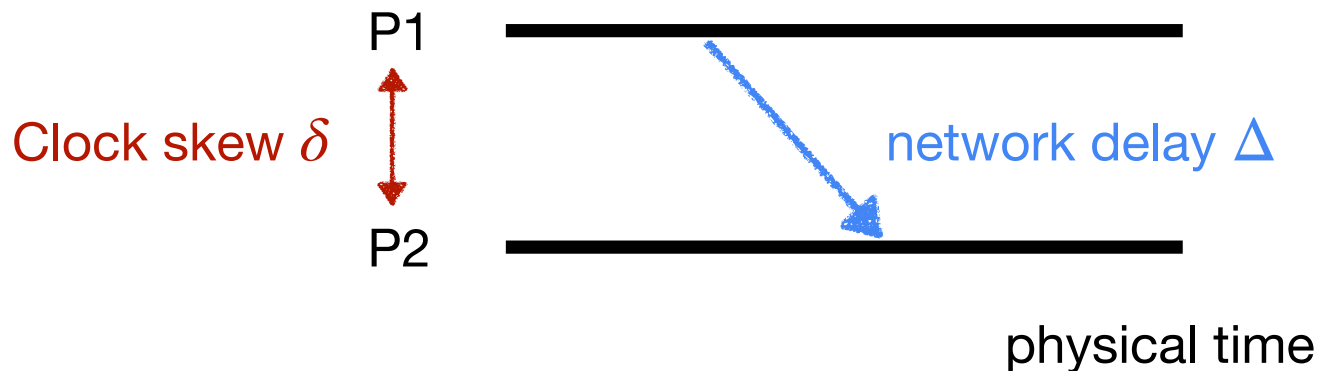physical time

# System model: network synchrony

- Synchronous

- Asynchronous

- Partially synchronous

<div style="background-color: yellow">

**Synchronous:**

Both $\delta$ and $\Delta$ have a fixed upper bound

</div>

P1

Clock skew $\delta$

network delay $\Delta$

P2

physical time

# System model: network synchrony

- Synchronous

- Asynchronous

- Partially synchronous

P1 ————————————————

Clock skew $\delta$        network delay $\Delta$

P2 ————————————————

physical time

# System model: network synchrony

- Synchronous

- Asynchronous

- Partially synchronous

P1 ——————————————

Clock skew $\delta$      network delay $\Delta$

P2 ——————————————

physical time

# System model: network synchrony

- Synchronous

- Asynchronous

- Partially synchronous

P1

Clock skew $\delta$

P2

network delay $\Delta$

physical time

# System model: network synchrony
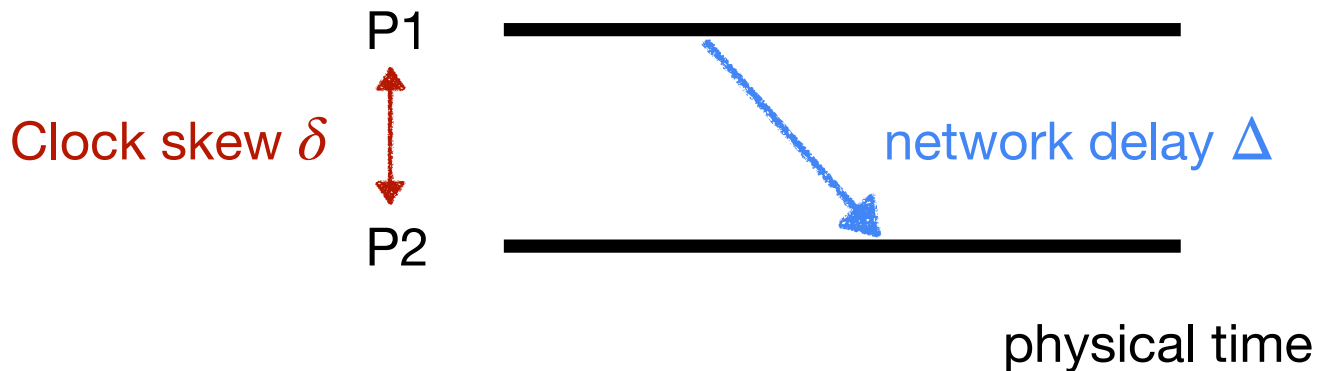
- Synchronous

- Asynchronous

- Partially synchronous

**Partially synchronous:**

Communication among servers can have a global stabilization time (GST), unknown to processors.

1. $\delta$ and $\Delta$ both exist but unknown, or

2. $\delta$ and $\Delta$ are known after GST

P1

Clock skew $\delta$

P2

network delay $\Delta$

physical time

# Let's design a simple consensus algorithm

- Assume processes cannot fail

- Synchronous network

- We'd like to have:

**Termination:** Eventually each correct process sets its decision variable

**Agreement:** Decision value of all correct processes is the same; if $p_i$ and $p_j$ are correct and ahem entered the decided state, then $d_i = d_j (i, j = 1, 2, \ldots, N)$

**Integrity/Validity:** If the correct processes all proposed the same value, then any correct process in the decided state has chosen that value.

# Service properties

**Termination:** Eventually each correct process sets its decision variable

**Agreement:** Decision value of all correct processes is the same; if $p_i$ and $p_j$ are correct and ahem entered the decided state, then $d_i = d_j(i, j = 1,2,\ldots,N)$

**Integrity/Validity:** If the correct processes all proposed the same value, then any correct process in the decided state has chosen that value.

# Service properties

**Termination:** Eventually each correct process sets its decision variable

**Agreement:** Decision value of all correct processes is the same; if $p_i$ and $p_j$ are correct and ahem entered the decided state, then $d_i = d_j (i, j = 1, 2, \ldots, N)$

**Integrity/Validity:** If the correct processes all proposed the same value, then any correct process in the decided state has chosen that value.

Something cannot happen

**Safety**

No two correct nodes decide differently

# Service properties

**Termination:** Eventually each correct process sets its decision variable

**Agreement:** Decision value of all correct processes is the same; if $p_i$ and $p_j$ are correct and ahem entered the decided state, then $d_i = d_j (i, j = 1, 2, \ldots, N)$

**Integrity/Validity:** If the correct processes all proposed the same value, then any correct process in the decided state has chosen that value.

Something cannot happen

**Safety**

No two correct nodes decide differently

**Liveness**

Nodes eventually decide

Something must happen

# Today's outline

The consensus problem

Network assumptions

**Failure assumptions**

Paxos

# Faults...

# Family of faults

- Crash faults

- Omission faults

  - Send omission

  - Receive omission

- Timing faults

# Family of faults

- Crash faults

- Omission faults

  - Send omission

  - Receive omission

- Timing faults

S1 ✖ &lt;k, v&gt;   S2

# Family of faults

- Crash faults

- Omission faults

  - Send omission

  - Receive omission

- Timing faults
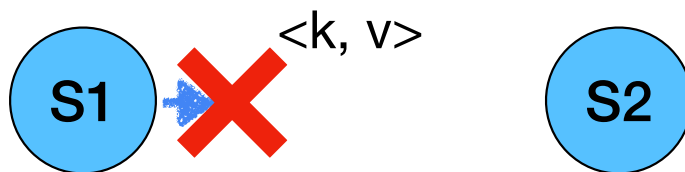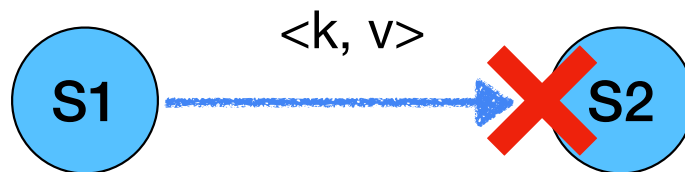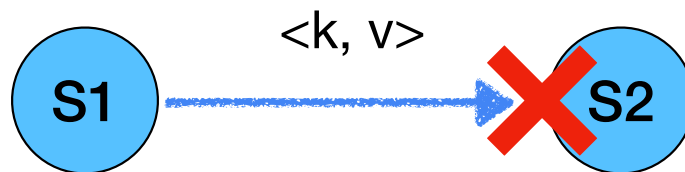
$<k, v>$

S1 ✕ S2

# Family of faults

- Crash faults

- Omission faults

  - Send omission

  - Receive omission

- Timing faults

# Family of faults

- Crash faults

- Omission faults

  - Send omission

  - Receive omission

- Timing faults

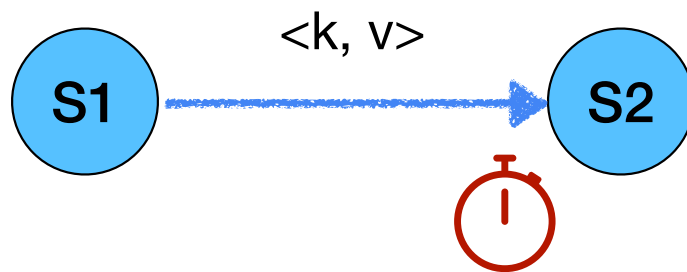<k, v>

S1 → S2

**if timer.timeout:**
**proceed without v**

# Family of faults

- Crash faults

- Omission faults

  - Send omission

  - Receive omission

- Timing faults

Worst thing that can happen:
S2 does not have the value

<k, v>

S1 →→→→ S2

# Family of faults

- Crash faults

- Omission faults

  - Send omission

  - Receive omission

- Timing faults

Worst thing that can happen: S2 does not have the value

$$S1 \xrightarrow{<k, v>} S2$$

Cli → S2

```
Read(k)
no failure? v : ∅
```
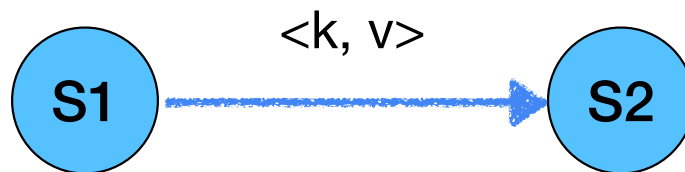
# Family of faults

- Crash faults
- Omission faults
  - Send omission
  - Receive omission
- Timing faults

Worst thing that can happen:
S2 does not have the value

<k, v>

S1 → S2

Cli

`Read(k)`
`no failure? v : ∅`

20

# Family of faults

**Benign faults**

- Crash faults
- Omission faults
  - Send omission
  - Receive omission
- Timing faults

**Byzantine faults**

- Any arbitrary behaviour, e.g.,
  - Stop responding
  - Send erroneous values

S1 → S2  <k, v>

# Family of faults

**Benign faults**

- Crash faults
- Omission faults
  - Send omission
  - Receive omission
- Timing faults

- Any arbitrary behaviour, e.g.,
  - Stop responding
  - Send erroneous values

**Byzantine faults**

**Any arbitrary behaviour**

<k, v>

S1 → S2

Cli    Cli    Cli

`Read(k)`
`response:` $\varnothing$

`Read(k)`
`response:` $x$

`Read(k)`
`response:` $y$

# Family of faults

**Benign faults**

- Crash faults
- Omission faults
  - Send omission
  - Receive omission
- Timing faults

- Any arbitrary behaviour, e.g.,
  - Stop responding
  - Send erroneous values

**Byzantine faults**

**Worst thing that can happen:
Any behaviour that can do the most harm**

**Any arbitrary behaviour**

S1 — <k, v> → S2

Cli → S2
Read(k)
response: $\varnothing$

Cli → S2
Read(k)
response: $x$

Cli → S2
Read(k)
response: $y$

# Family of faults: summary

Fault tolerance

**Benign faults**

- Crash faults
- Omission faults
  - Send omission
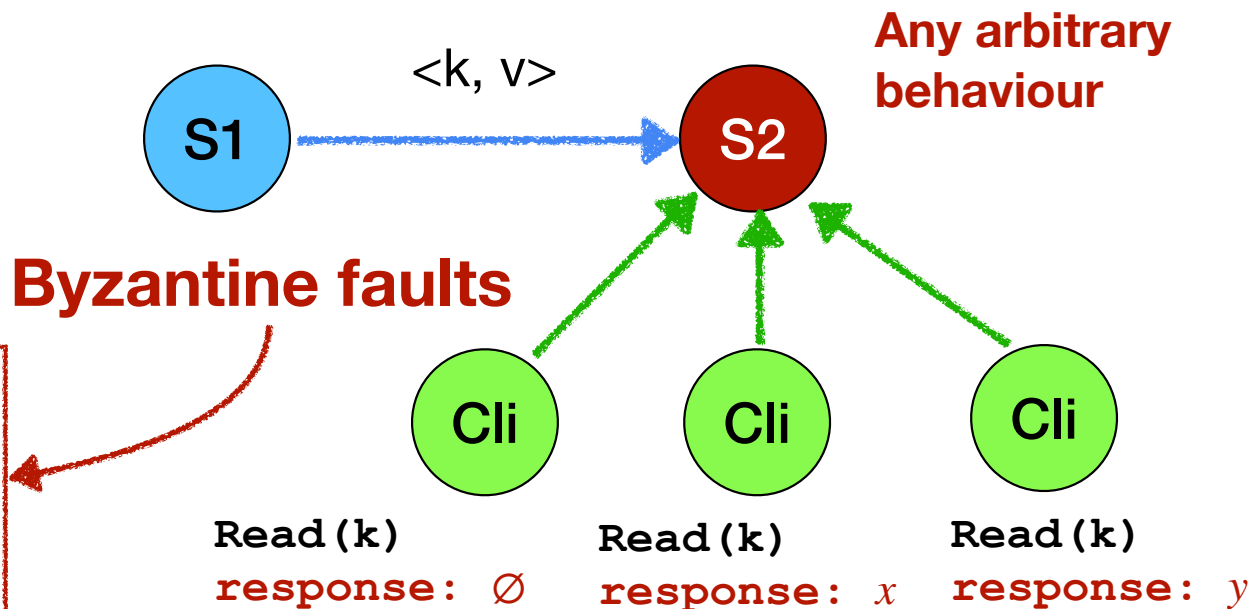  - Receive omission
- Timing faults

- Crash fault tolerance (CFT) algorithms
  - Paxos, ViewStamped Replication, Raft [ATC'13]
- Applications (everything distributed):
  - File systems: HDFS and GFS
  - Databases: Google Spanner and etcd
  - Coordination: Chubby and Zookeeper

**Byzantine faults**

- Any arbitrary behaviour, e.g.,
  - Stop responding
  - Send erroneous values

- Byzantine fault tolerance (BFT) algorithms
  - PBFT [OSDI'99], HotStuff [PODC'21], Pompe [OSDI'22]
- Applications (safety critical):
  - Unreliable hardware: Airplanes
  - Blockchains: Facebook Diem, Microsoft CCF

# Algorithms we will talk about

- Paxos:                                                    <- Today's topic

  - How to choose a value under **benign failures**

- Raft [ATC'14]:

  - How to replicate log under **benign failures**?

- PBFT [OSDI'99]:

  - How to replicate log under **Byzantine (arbitrary) failures**?

# Today's outline

The consensus problem

Network assumptions

Failure assumptions

**Paxos**

# Paxos

- Papers:

    - Lamport L. The part-time parliament[J]. ACM Transactions on Computer Systems (TOCS), 1998, 16(2): 133-169.

    - Lamport L. Paxos made simple[J]. ACM Sigact News, 2001, 32(4): 18-25.

- System model

    - Asynchronous

    - CFT: tolerating benign faults (non-Byzantine)

# Fundamental #1: Server roles

**Proposers** (leader)

- receive client requests

- propose received requests

- coordinate consensus process for its proposed requests

**Acceptors** (follower)

- respond to requests from proposers

- validate states of requests

- store chosen values and state of the process

# Fundamental #1: Server roles

**Proposers** (leader)

- receive client requests
- propose received requests
- coordinate consensus process for its proposed requests

**Acceptors** (follower)

- respond to requests from proposers
- validate states of requests
- store chosen values and state of the process

**Learners** (subscriber)

- want to know which value is chosen
- subscribe to acceptors
  - one or a few learners communicate with acceptors
  - propagate the message among learners

# Fundamental #1: Server roles

**Proposers** (leader)

- receive client requests
- propose received requests
- coordinate consensus process for its proposed requests

**Acceptors** (follower)

- respond to requests from proposers
- validate states of requests
- store chosen values and state of the process

**Learners** (subscriber)

- want to know which value is chosen
- subscribe to acceptors
  - one or a few learners communicate with acceptors
  - propagate the message among learners

According to the application that uses Paxos, a server can be a proposer, an acceptor, or both

# Fundamental #2: Proposals

- Each proposal has a unique number (proposal number)

  - Higher number take a priority over lower numbers

  - Similar to Lamport clock, proposers can increase a proposal number

# Fundamental #3: Phases

**Prepare phase
(Phase 1)**

**Phase 1.** (a) A proposer selects a proposal number $n$ and sends a *prepare* request with number $n$ to a majority of acceptors.

(b) If an acceptor receives a *prepare* request with number $n$ greater than that of any *prepare* request to which it has already responded, then it responds to the request with a promise not to accept any more proposals numbered less than $n$ and with the highest-numbered proposal (if any) that it has accepted.

**Accept phase
(Phase 2)**

**Phase 2.** (a) If the proposer receives a response to its *prepare* requests (numbered $n$) from a majority of acceptors, then it sends an *accept* request to each of those acceptors for a proposal numbered $n$ with a value $v$, where $v$ is the value of the highest-numbered proposal among the responses, or is any value if the responses reported no proposals.

(b) If an acceptor receives an *accept* request for a proposal numbered $n$, it accepts the proposal unless it has already responded to a *prepare* request having a number greater than $n$.

# Proposers

(1) Choose new proposal number $n$.
(2) Broadcast Prepare($n$) to all servers.

(4) When responses received from majority,
if any *acceptedValue* returned, replace
value with *acceptedValue* for highest
*acceptedProposal*.

(5) Broadcast Accept($n$, *value*) to all servers

(7) When responses received from majority:
-> Any rejections (result > n) : go to (1)
-> Otherwise, value is chosen

# Acceptors
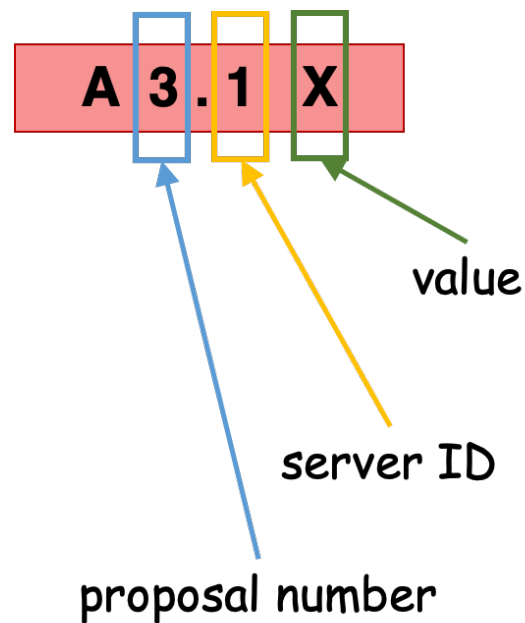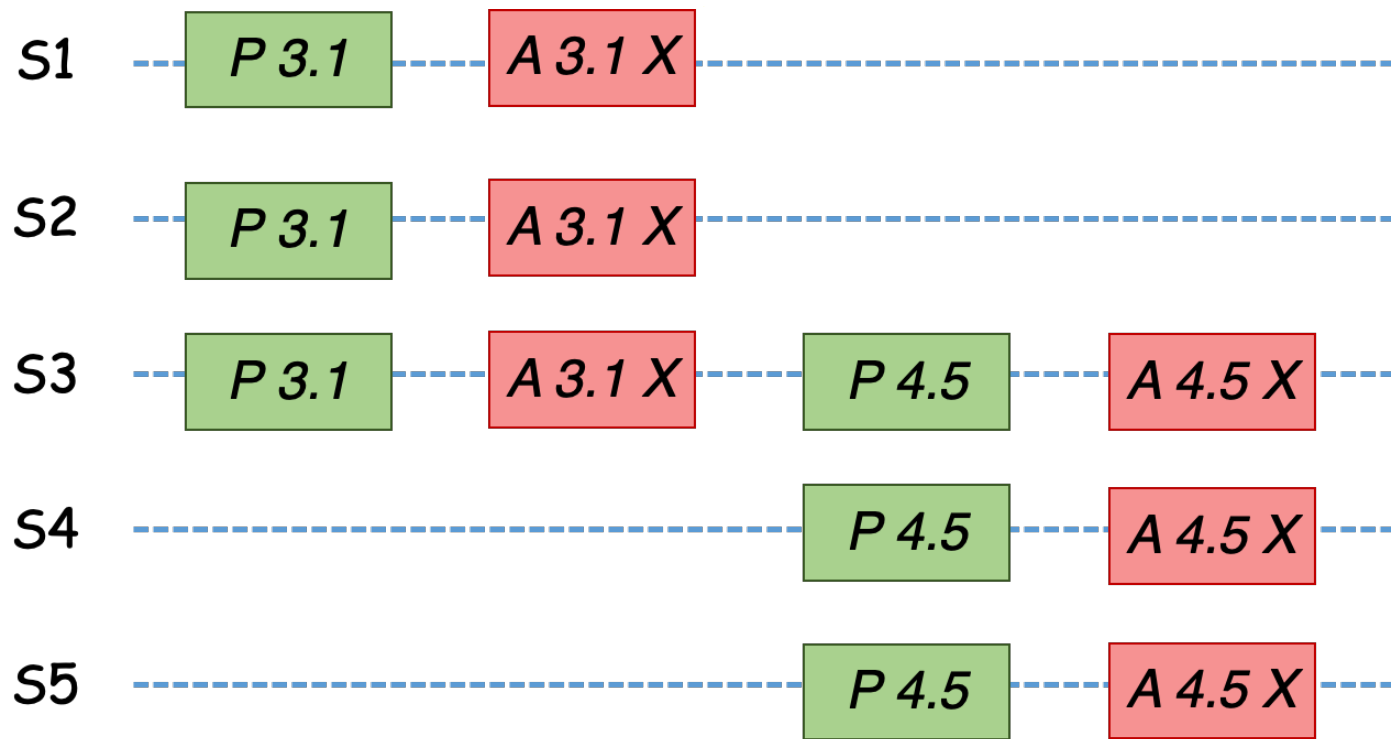
(3) Respond to Prepare($n$):
-> If n > *minProposal*, then *minProposal* = n
-> Return (*acceptedProposal*, *acceptedValue*)

(6) Respond to Accept($n$, *value*):
-> If $n$ >= *minProposal* then
*acceptedProposal* = *minProposal* = $n$;
*acceptedValue* = *value*;
-> Return (minProposal)

Acceptors must record *minProposal*, *acceptedProposal*, and *acceptedValue* on stable storage
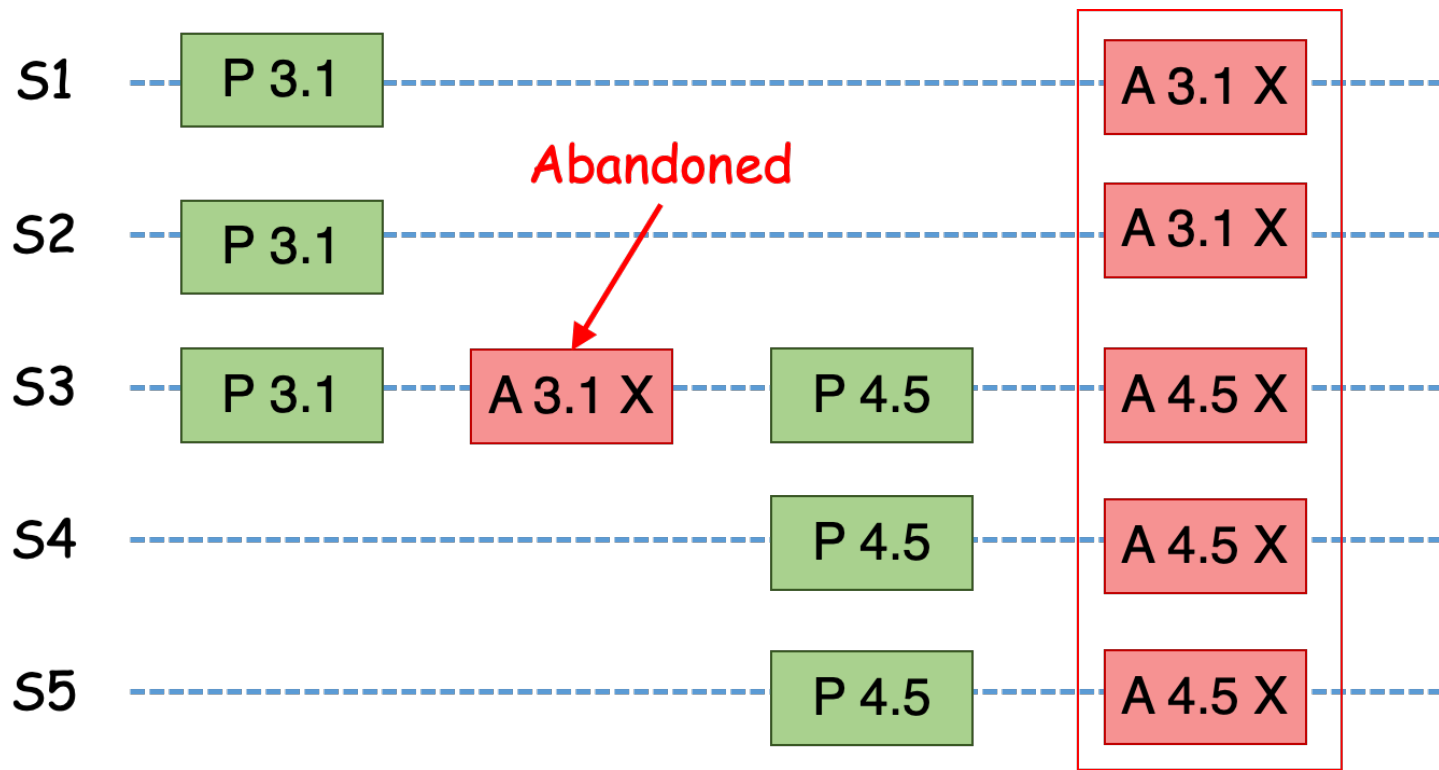(disk).

# Value chosen in different proposal numbers

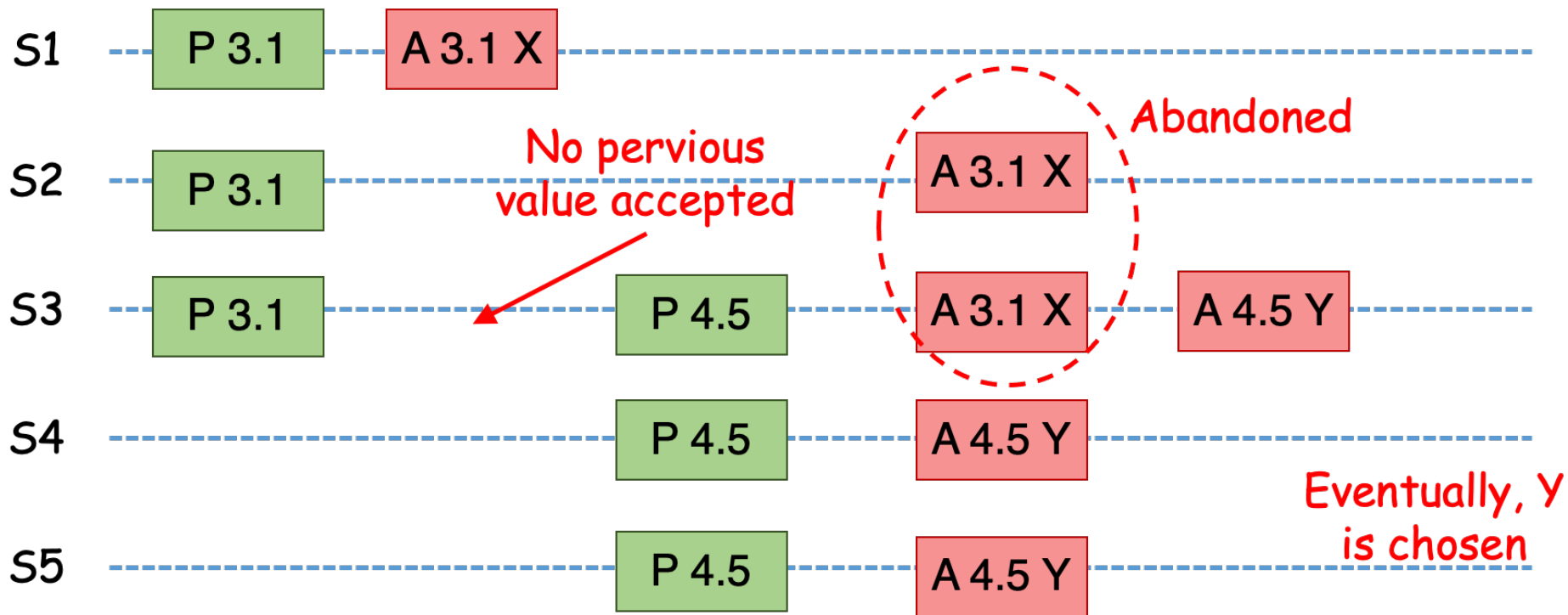A proposer "learns" a already chosen value

# Value chosen in different proposal numbers

## A proposer "learns" a not chosen value

A new value is chosen in different proposal numbers

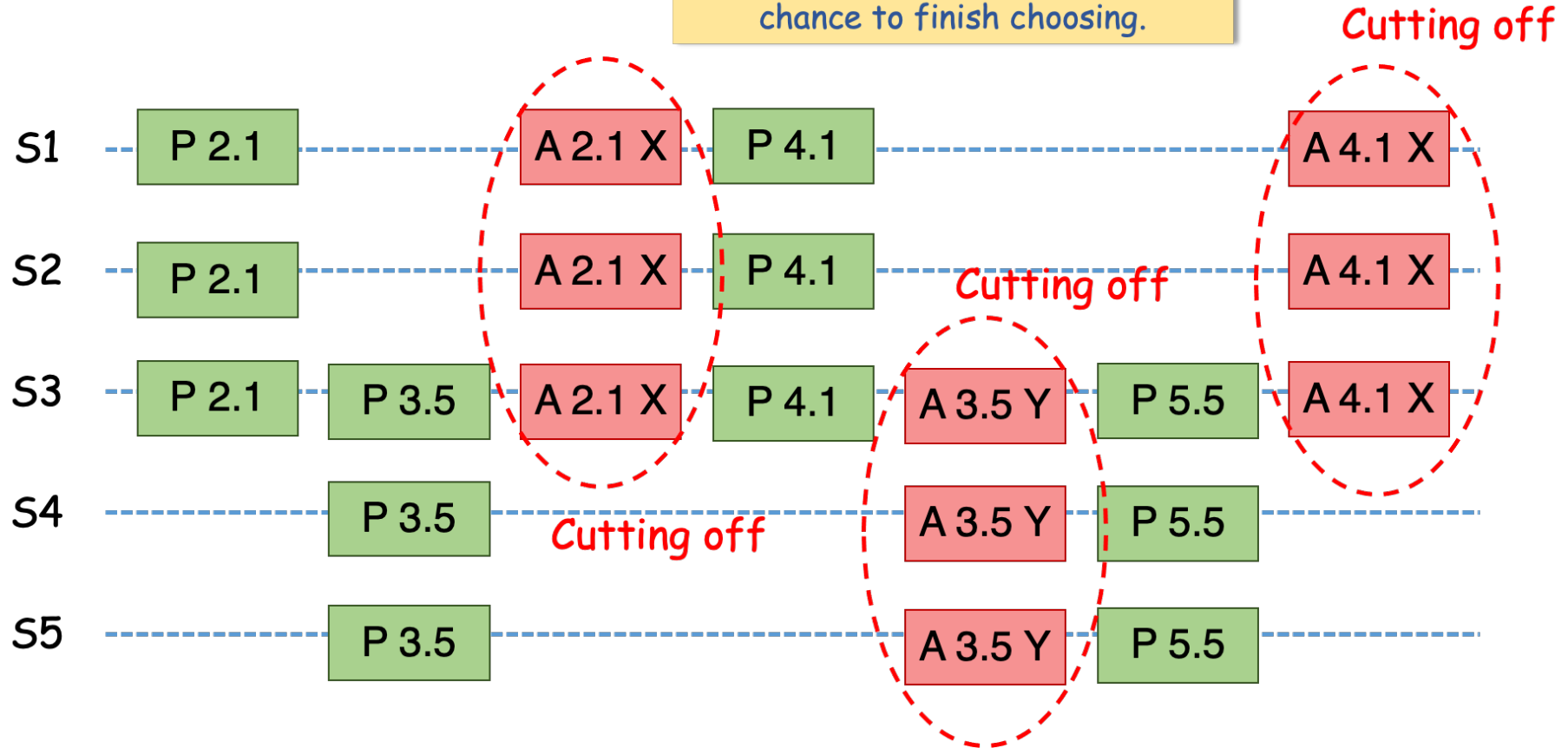A proposer does not see an unchosen value

# Livelock



Hint:
=> one solution:
    Randomized delay before restarting. Give other proposers a chance to finish choosing.

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| **S1** | P 2.1 | | A 2.1 X | P 4.1 | | | A 4.1 X |
| **S2** | P 2.1 | | A 2.1 X | P 4.1 | | | A 4.1 X |
| **S3** | P 2.1 | P 3.5 | A 2.1 X | P 4.1 | A 3.5 Y | P 5.5 | A 4.1 X |
| **S4** | | P 3.5 | | | A 3.5 Y | P 5.5 | |
| **S5** | | P 3.5 | | | A 3.5 Y | P 5.5 | |

Cutting off

# Summary of Paxos

- Anyone can be a proposer/leader

  - Advantages?

  - Disadvantages?

- Only proposer knows which value has been chosen

- If other servers want to know, must execute Paxos with their own proposal

- Competing proposers can cause a livelock

# Worksheet