

MC714 - Sistemas Distribuídos

Introdução ao algoritmo de Paxos

Islene Calciolari Garcia

Instituto de Computação - Unicamp

Primeiro Semestre de 2015

Sumário

Tutorial OPODIS2013

- Único servidor

- Dois servidores e ordem das mensagens

- Crash dos servidores

- Partições de rede

- Quanto azar: crash e partições

- Mais azar: vários crashes

From 2PC to Paxos

Paxos Explained from Scratch

Hein Meling and Leander Jehl



University of
Stavanger

hein.meling@uis.no

International Conference On Principles of Distributed Systems

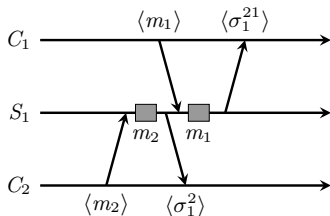
What is Paxos and why is it Relevant?

- Fault tolerant consensus protocol
- Used to order client requests in a fault tolerant server
 - For example a fault tolerant resource manager
- Used in production systems: Chubby, ZooKeeper, and Spanner
- It is always safe

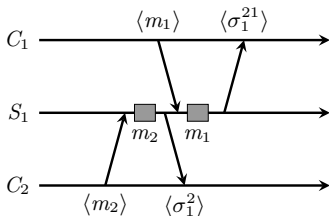
Objectives and Approach

- Explain Paxos
 - Using visual aids
 - In a step-wise manner
 - With minimal changes in each step
- Objective
 - Understand why it works and why the solution is necessary
 - (no focus on how to implement or formally prove it)
- Approach
 - Use a simple client/server system as base
 - To build fault tolerant server (replicated state machine)
 - Construct Multi-Paxos
 - Decompose Multi-Paxos into Paxos

A Stateful Service: *SingleServer*



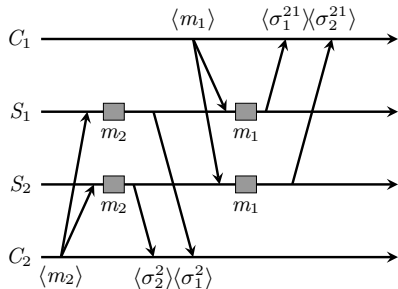
A Stateful Service: *SingleServer*



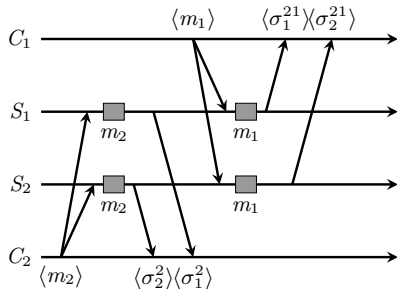
- Client C_2 sees: σ^2
- Client C_1 sees: σ^{21}
- Corresponds to execution sequence: $m_2 m_1$

We Want to Make the Service Fault Tolerant!

Fault Tolerance with Two Servers



Fault Tolerance with Two Servers

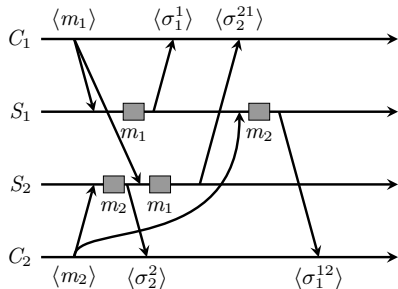


- Client C_2 sees: σ^2
- Client C_1 sees: σ^{21}
 - σ^2 is a prefix of σ^{21}
- Corresponds to execution sequence: $m_2 m_1$

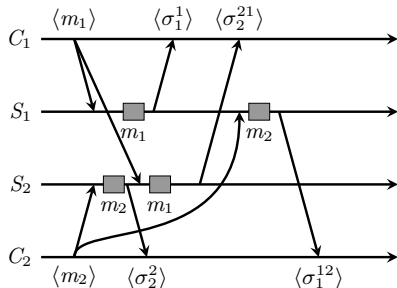
Deterministic State Machine

- The service is implemented as a deterministic state machine
- Thus processing requests results in unique state transitions:
 - Therefore $\sigma_1^2 = \sigma_2^2$ and $\sigma_1^{21} = \sigma_2^{21}$.
- Clients can detect and suppress identical replies

Fault Tolerance with Two Servers: Whoops!



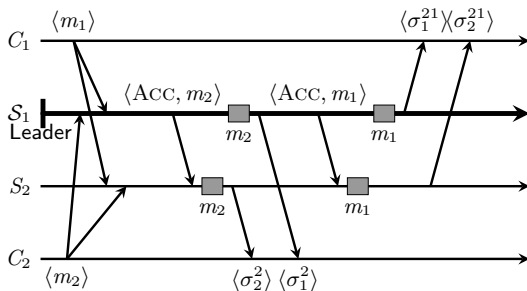
Fault Tolerance with Two Servers: Whoops!



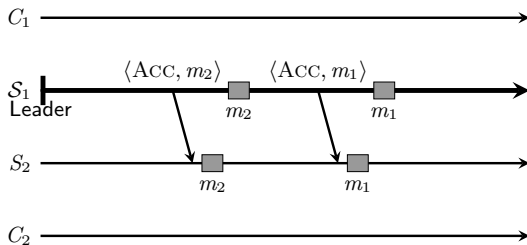
- Client C_2 sees: $\sigma^2 \sigma^{12}$
 - σ^2 is not a prefix of σ^{12}
- Client C_1 sees: $\sigma^1 \sigma^{21}$
 - σ^1 is not a prefix of σ^{21}
- Corresponds to execution sequence at
 - S_1 : $m_1 m_2$
 - S_2 : $m_2 m_1$

We Need to Order Client Requests!

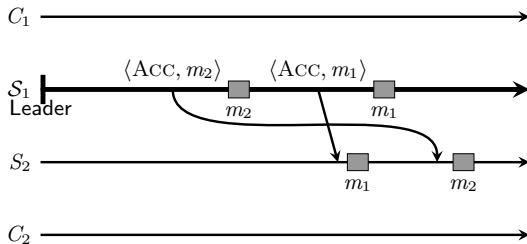
Let's Designate a Leader to Order Requests



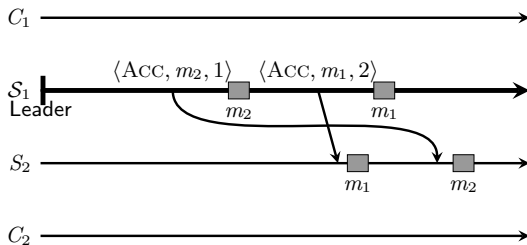
Without Clients



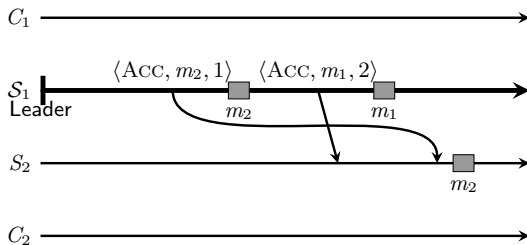
Problem: Also Accept Messages can be Reordered



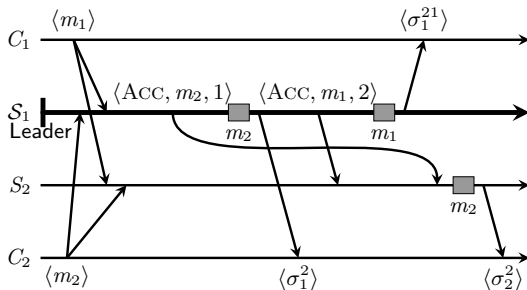
Add Sequence Numbers



Discard Out-of-Order Messages



Now with Clients



Clients Observe The Same Server States as Before

- Client C_2 sees: σ^2
- Client C_1 sees: σ^{21}
- However, S_2 didn't execute m_1
 - Q: What to do?

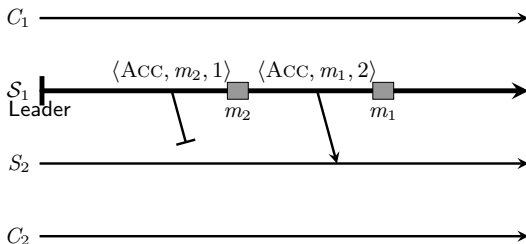
Clients Observe The Same Server States as Before

- Client C_2 sees: σ^2
- Client C_1 sees: σ^{21}
- However, S_2 didn't execute m_1
 - Q: What to do?
 - A1: Buffer

Clients Observe The Same Server States as Before

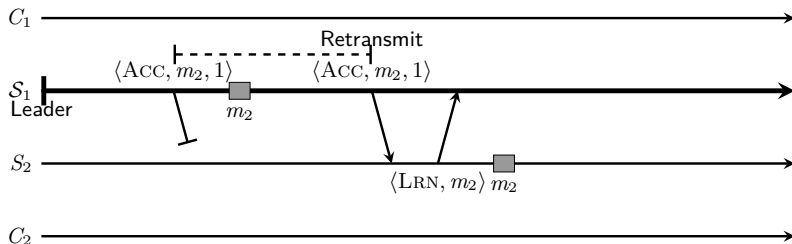
- Client C_2 sees: σ^2
- Client C_1 sees: σ^{21}
- However, S_2 didn't execute m_1
 - Q: What to do?
 - A1: Buffer
 - A2: Retransmission mechanism

Problem: Message Loss – S_2 Won't Execute Anything

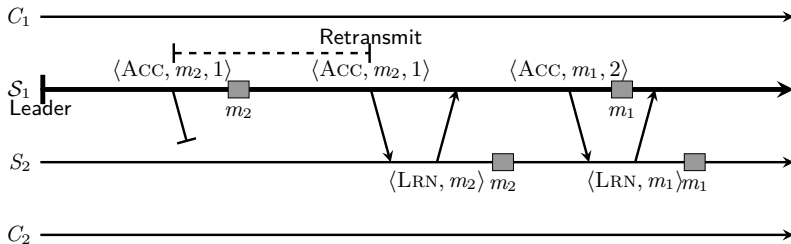


We Need a Retransmission Mechanism!

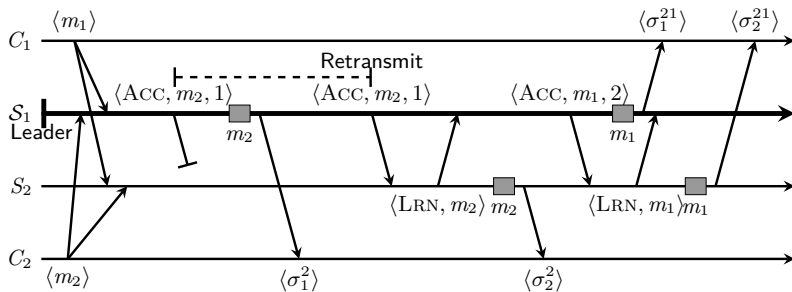
A Learn Stops Retransmission



Don't Send New Accept Until Learn



With Clients



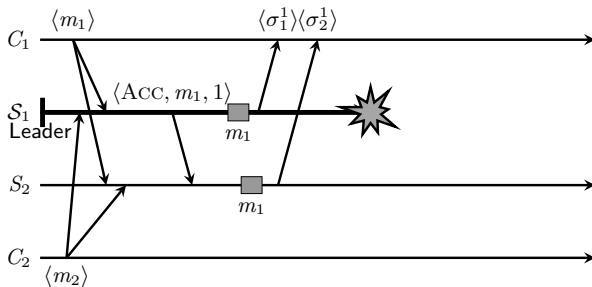
- A leader
 - To decide the order of client requests
 - By sending an accept message to S_2
- Sequence numbers
 - To cope with message reordering
- Retransmission mechanism
 - To cope with message loss
 - Leader only sends next accept when learn from S_2
 - Allows leader to *make progress*, as long as messages are not lost infinitely often

- A leader
 - To decide the order of client requests
 - By sending an accept message to S_2
- Sequence numbers
 - To cope with message reordering
- Retransmission mechanism
 - To cope with message loss
 - Leader only sends next accept when learn from S_2
 - Allows leader to *make progress*, as long as messages are not lost infinitely often

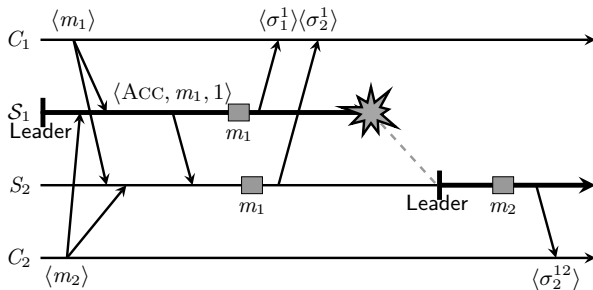
Combination of mechanisms:
RetransAccept protocol

What About Server Crashes?

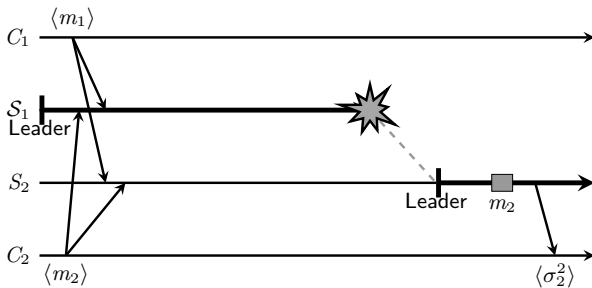
Crash



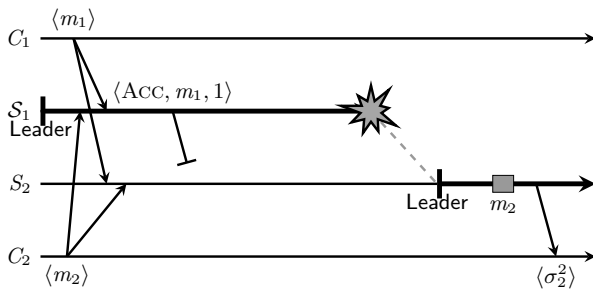
Crash: Leader Takeover



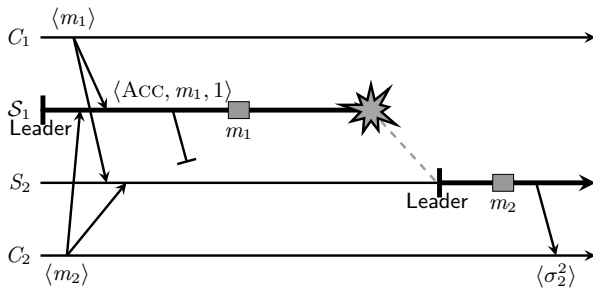
Single Server Rule: Case 1



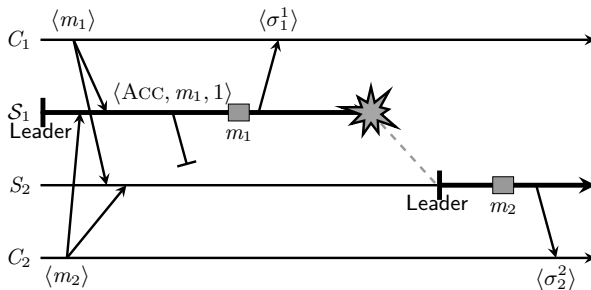
Single Server Rule: Case 2



Single Server Rule: Case 3



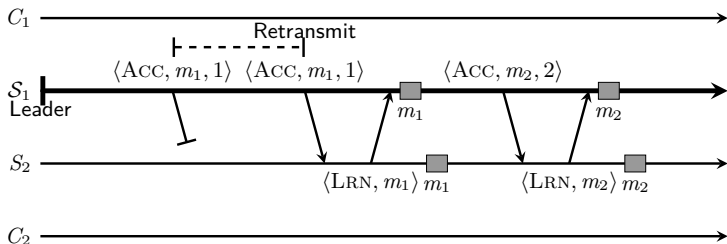
Single Server Rule: Case 4 – A Problem



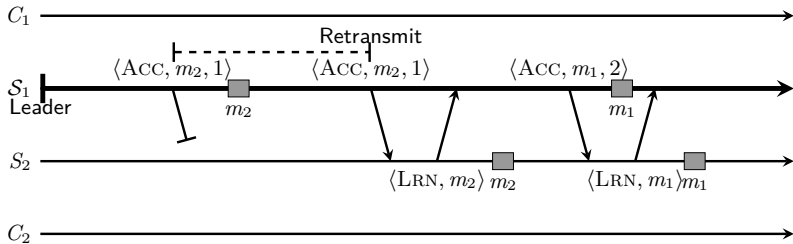
Single Server Rule: Case 4 – A Problem

- Imagine that (S_1, S_2) implements a fault tolerant resource manager, e.g. a lock service
- Both clients could have gotten the lock

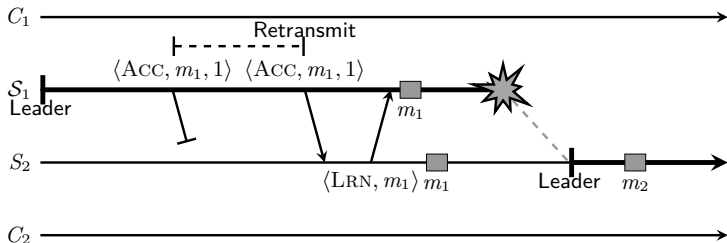
Solution: Leader Waits for Learn Before Executing



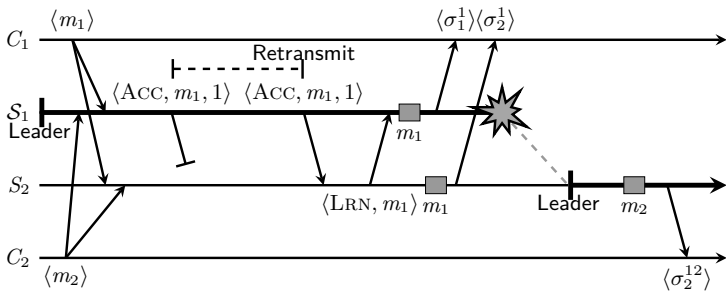
Recall Earlier Version



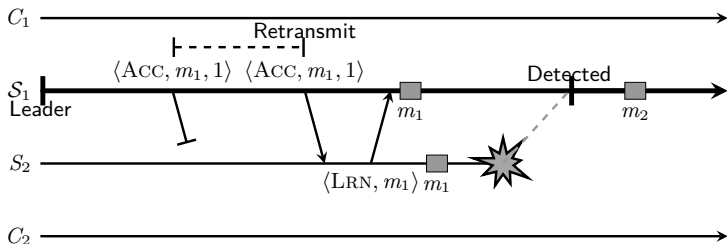
Now Leader Takeover is Safe



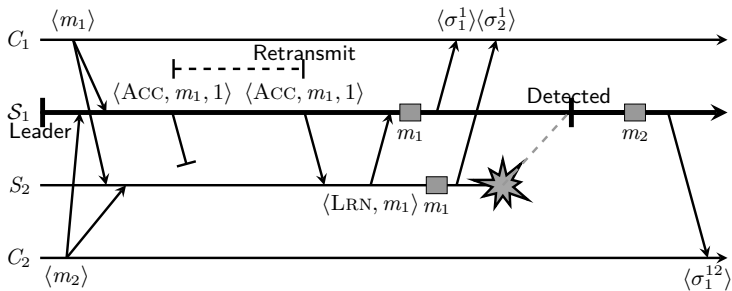
Let's Add Client Messages



Leader Remain in Control when S_2 Crash



Let's Add Client Messages Again



Recap: The Problem

- When we detect a server crash
 - Adopt the *SingleServer* protocol

Recap: The Problem

- When we detect a server crash
 - Adopt the *SingleServer* protocol
- Problem with our *RetransAccept* protocol:
 - The leader might have replied to a client and then crashed, without ensuring that S_2 saw the accept
 - S_2 takes over and may execute a different request in *SingleServer* mode

Recap: WaitForLearn Protocol

- The leader always waits for a learn message from S_2
 - Think of it as an acknowledgement

Recap: WaitForLearn Protocol

- The leader always waits for a learn message from S_2
 - Think of it as an acknowledgement
- S_2 can execute after seeing an accept from the leader
 - This is because the accept message is also an implicit learn

Recap: WaitForLearn Protocol

- The leader always waits for a learn message from S_2
 - Think of it as an acknowledgement
- S_2 can execute after seeing an accept from the leader
 - This is because the accept message is also an implicit learn
- Q: What happens if the learn message to the leader is lost?

Recap: WaitForLearn Protocol

- The leader always waits for a learn message from S_2
 - Think of it as an acknowledgement
- S_2 can execute after seeing an accept from the leader
 - This is because the accept message is also an implicit learn
- Q: What happens if the learn message to the leader is lost?
- A: The leader uses *RetransAccept*; the accept will be retransmitted.
So no need for another retransmit protocol.

Somewhat Rougher Road Ahead!

False Detection

- So far we have assumed that failure detection is accurate

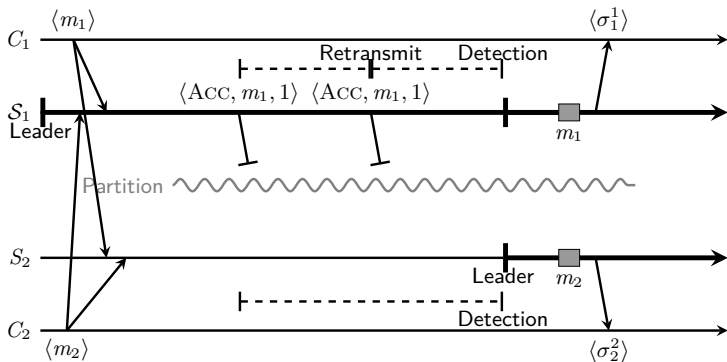
False Detection

- So far we have assumed that failure detection is accurate
- But in an asynchronous environment
 - There is always a chance of false detection
 - Because it is impossible to pick the right timeout delay

False Detection

- So far we have assumed that failure detection is accurate
- But in an asynchronous environment
 - There is always a chance of false detection
 - Because it is impossible to pick the right timeout delay
- We now consider false detection in the context of network partitions

Problem: Network Partitions



Network Partition

- Each server can switch to *SingleServer* mode (no coordination) and make progress

Network Partition

- Each server can switch to *SingleServer* mode (no coordination) and make progress
- But it will lead to inconsistencies
 - S_1 has state σ^1
 - S_2 has state σ^2

Network Partition

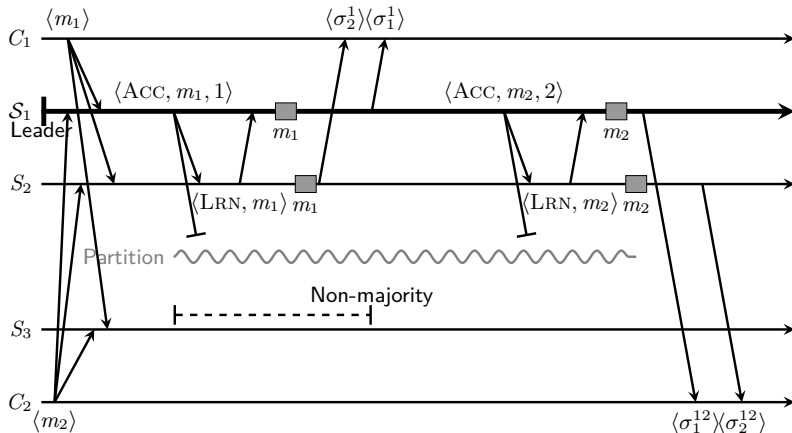
- Each server can switch to *SingleServer* mode (no coordination) and make progress
- But it will lead to inconsistencies
 - S_1 has state σ^1
 - S_2 has state σ^2
- Reconciling the state divergence
 - Involves rollback on multiple clients

Network Partition

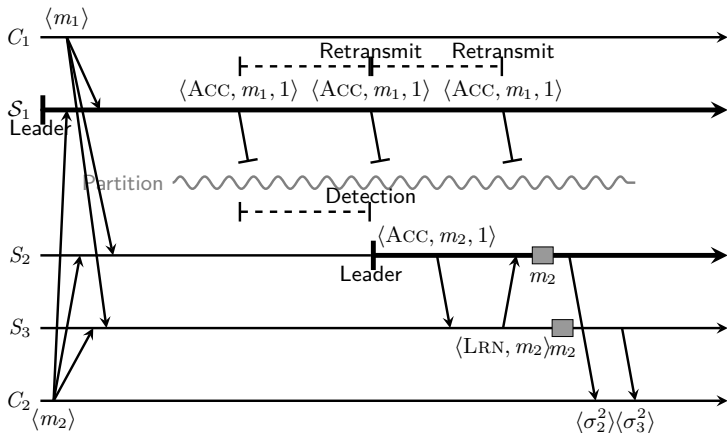
- Each server can switch to *SingleServer* mode (no coordination) and make progress
- But it will lead to inconsistencies
 - S_1 has state σ^1
 - S_2 has state σ^2
- Reconciling the state divergence
 - Involves rollback on multiple clients
 - Quickly becomes unmanageable

We Want to Avoid Relying on Clients!

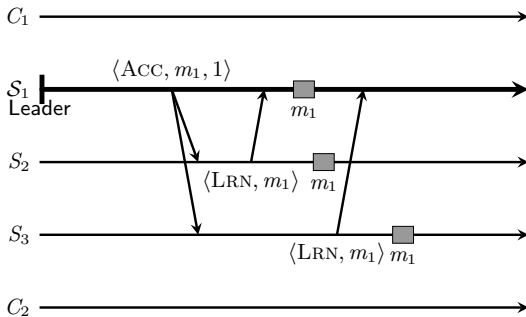
Add Another Server; Make Progress in Majority Partition



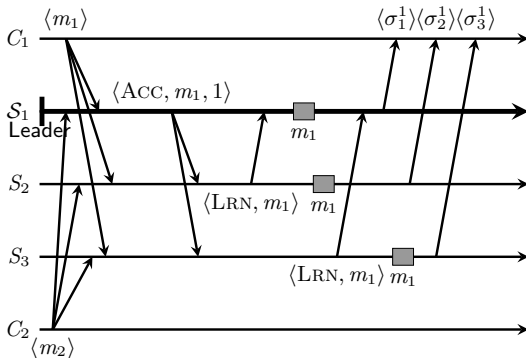
New Leader in Majority Partition



WaitForLearn Without Partition



WaitForLearn With Clients



Recap: Network Partition

- We added another server, S_3
 - To avoid rollback using clients

Recap: Network Partition

- We added another server, S_3
 - To avoid rollback using clients
- We still use the *WaitForLearn* protocol
 - To ensure that another server has seen the accept message

Recap: Network Partition

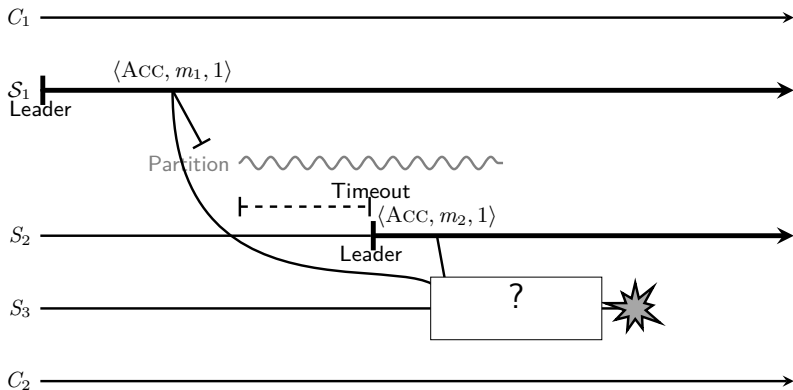
- We added another server, S_3
 - To avoid rollback using clients
- We still use the *WaitForLearn* protocol
 - To ensure that another server has seen the accept message
- Leader only needs to wait for *one* learn before executing the request
 - Allows the leader to make progress,
 - when another server has crashed or is temporarily unavailable

Recap: Network Partition

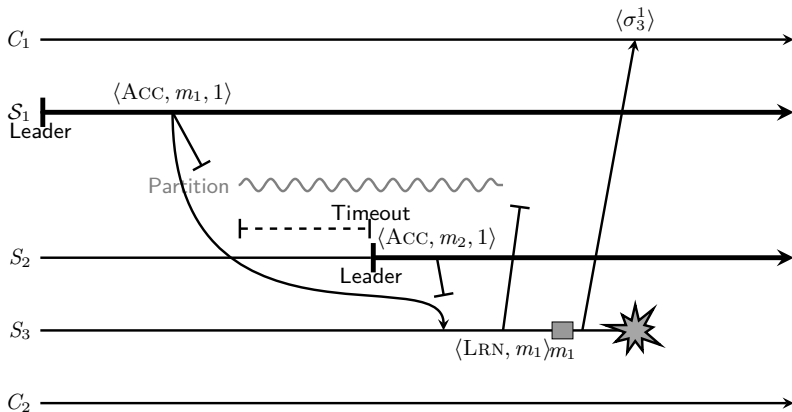
- We added another server, S_3
 - To avoid rollback using clients
- We still use the *WaitForLearn* protocol
 - To ensure that another server has seen the accept message
- Leader only needs to wait for *one* learn before executing the request
 - Allows the leader to make progress,
 - when another server has crashed or is temporarily unavailable
- But we still only tolerate one concurrent failure
 - Either a crash or a network partition

What can go Wrong: Concurrent Crash and Partition

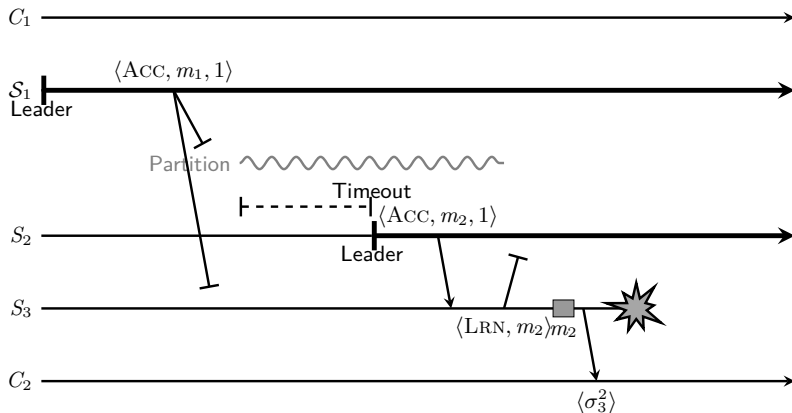
Concurrent Crash and Partition



Crash and Partition: Outcome 1 – m_1 Executed



Crash and Partition: Outcome 2 – m_2 Executed



Recap: Crash and Partition

- S_3 crashed
 - But *it could* have executed either m_1 or m_2
 - And replied to a client

Recap: Crash and Partition

- S_3 crashed
 - But *it could* have executed either m_1 or m_2
 - And replied to a client
- Other servers cannot determine which message, if any, was executed

Recap: Crash and Partition

- S_3 crashed
 - But *it could* have executed either m_1 or m_2
 - And replied to a client
- Other servers cannot determine which message, if any, was executed
 - Maybe we could talk to clients?
 - We don't want to rely on clients!

Explicit Leader Change Mechanism

- Above problem is rooted in possibility of false detection
 - Can lead to several servers thinking they are leaders
 - And sending accept messages concurrently

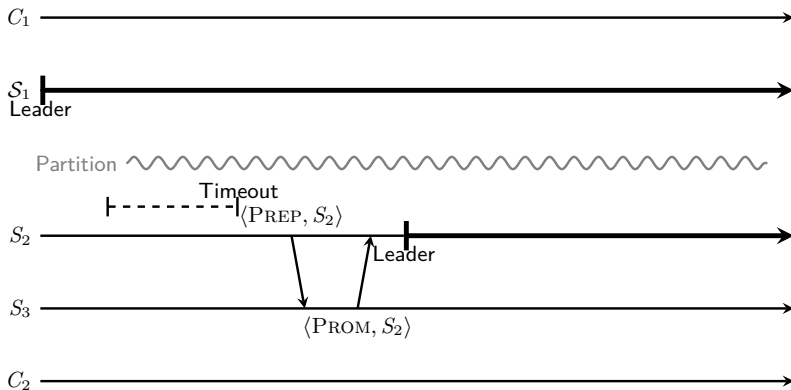
Explicit Leader Change Mechanism

- Above problem is rooted in possibility of false detection
 - Can lead to several servers thinking they are leaders
 - And sending accept messages concurrently
- It can be solved by an explicit leader takeover protocol

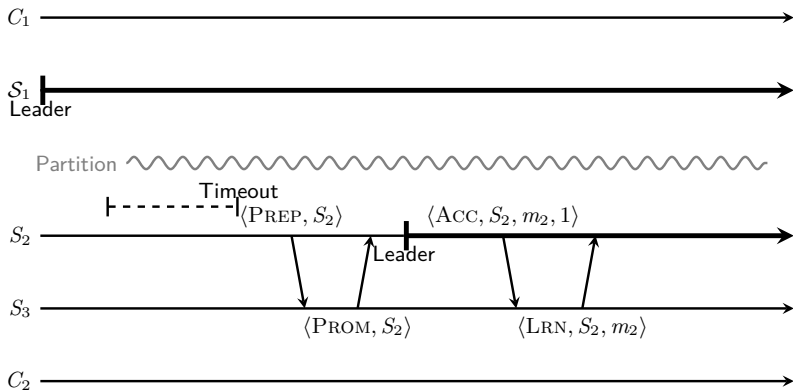
Explicit Leader Change Mechanism

- Above problem is rooted in possibility of false detection
 - Can lead to several servers thinking they are leaders
 - And sending accept messages concurrently
- It can be solved by an explicit leader takeover protocol
- We need a way to
 - Distinguish messages from different leaders
 - Change the leader

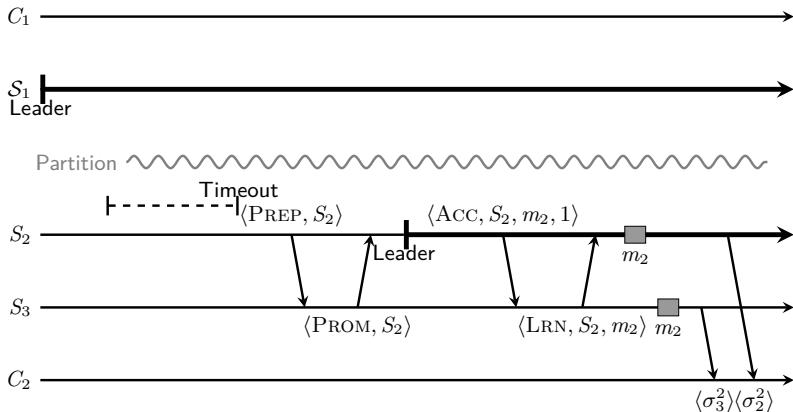
Explicit Leader Change



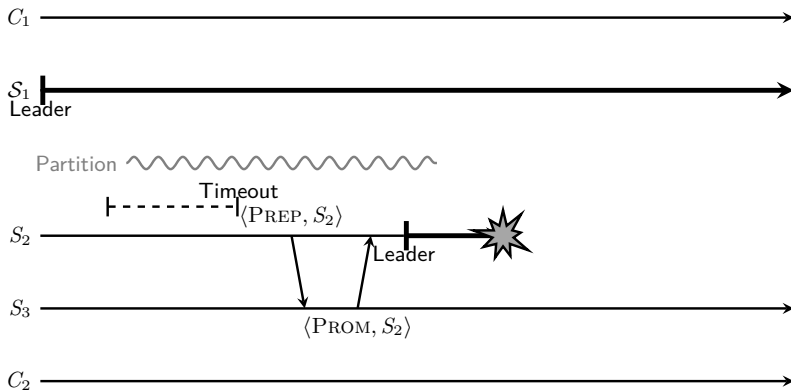
Leader Identifiers in Accept and Learn Messages



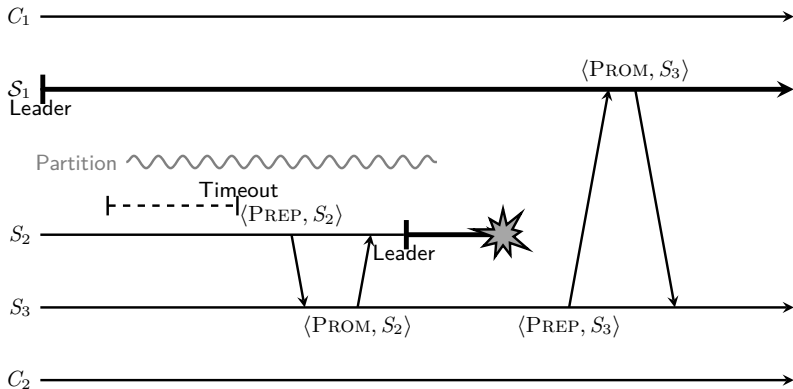
With Client Replies



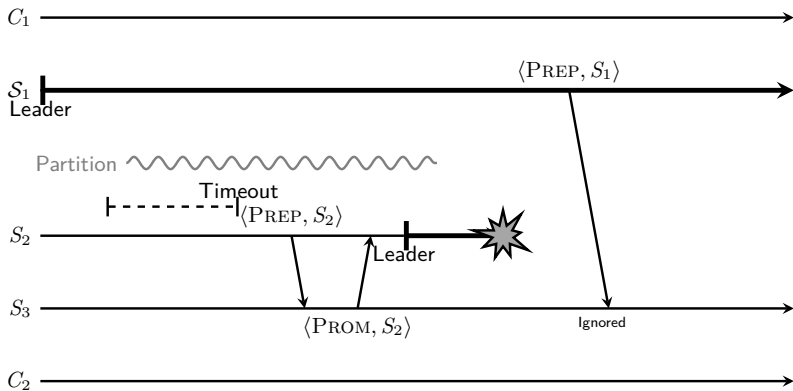
What Happens Now?



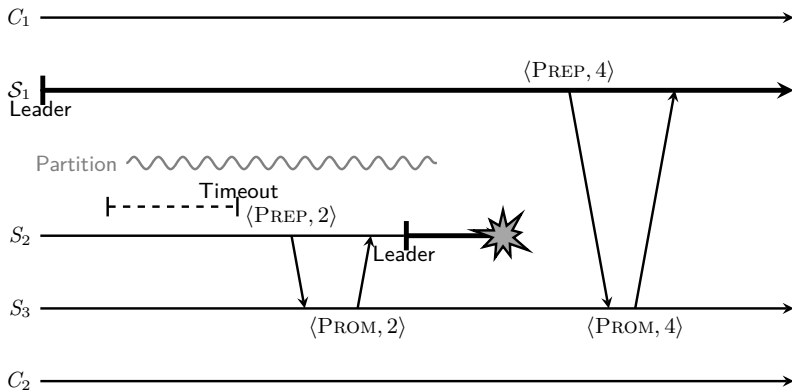
S_3 Takes Over?



S_1 Takes Over Again?



Replace Leader Identifiers With Round Numbers



Recap: Leader Change

- Added round number rnd in messages
 - To identify the leader
 - $\langle ACC, rnd, m, seqno \rangle$: Sent by leader of round rnd
 - $\langle LRN, rnd, m \rangle$: Sent to leader of round rnd

Recap: Leader Change

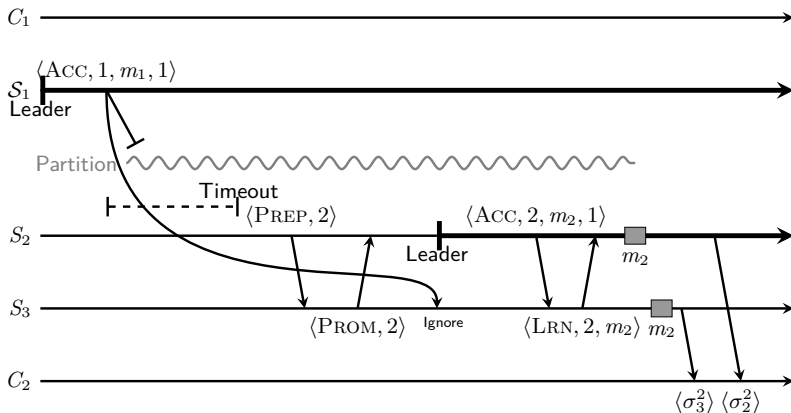
- Added round number rnd in messages
 - To identify the leader
 - $\langle ACC, rnd, m, seqno \rangle$: Sent by leader of round rnd
 - $\langle LRN, rnd, m \rangle$: Sent to leader of round rnd
 - Round numbers are assigned:
 - S_1 : 1, 4, 7, ...
 - S_2 : 2, 5, 8, ...
 - S_3 : 3, 6, 9, ...
 - Skipping rounds is possible

Recap: Leader Change

- Added round number rnd in messages
 - To identify the leader
 - $\langle ACC, rnd, m, seqno \rangle$: Sent by leader of round rnd
 - $\langle LRN, rnd, m \rangle$: Sent to leader of round rnd
 - Round numbers are assigned:
 - S_1 : 1, 4, 7, ...
 - S_2 : 2, 5, 8, ...
 - S_3 : 3, 6, 9, ...
 - Skipping rounds is possible
- Added two new messages
 - $\langle PREP, rnd \rangle$: Request to become leader for round rnd
 - $\langle PROM, rnd \rangle$: Promise not to accept messages from a lower round than rnd (i.e. an older leader)

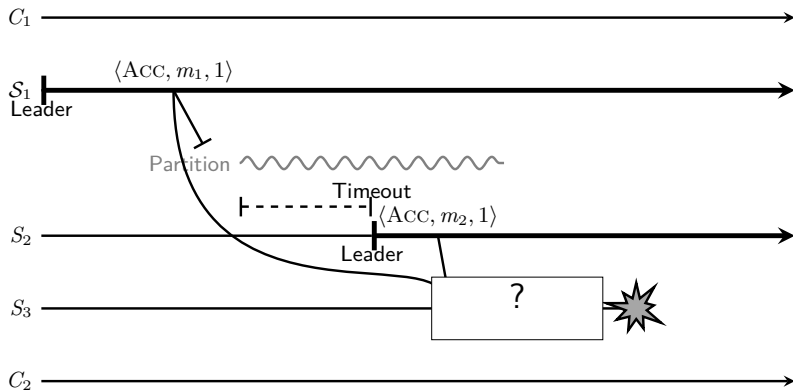
Let's Apply This Together With Accept and Learn

S_3 Ignores Accept Message From Old Leader

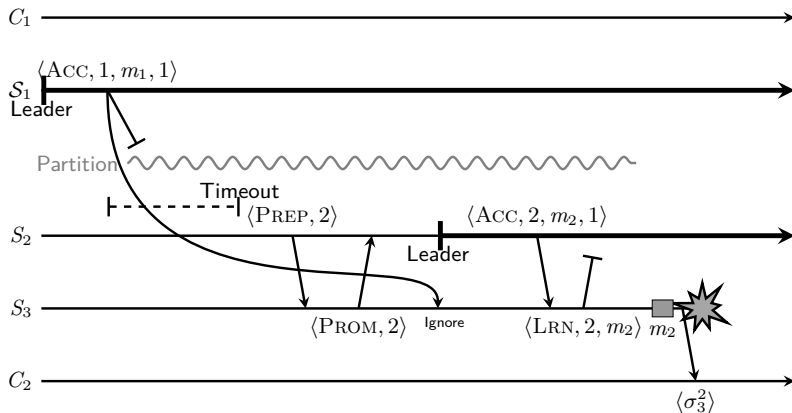


Let's Recall the Problem we are Trying to Solve

We Don't Know What S_3 Did Before Crashing



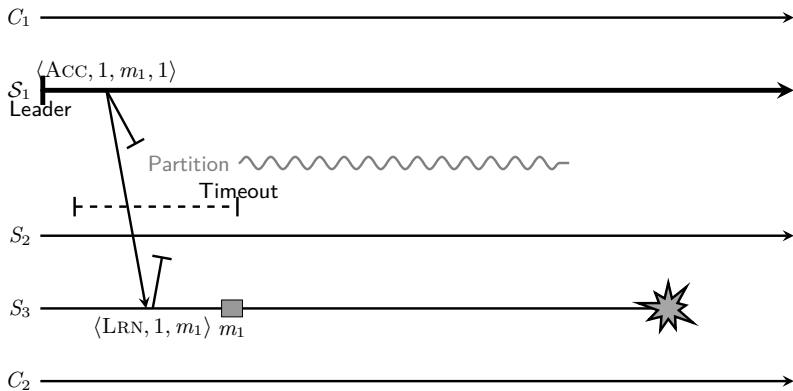
Do We Know Now?



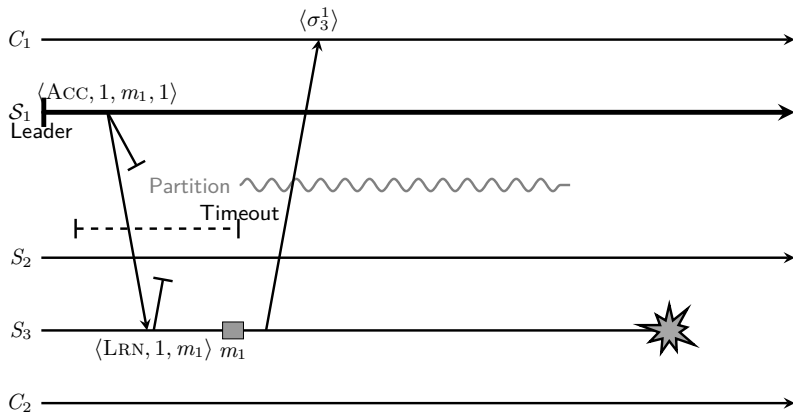
No we don't!

But it is Safe to Continue
as If m_2 Had Been Executed

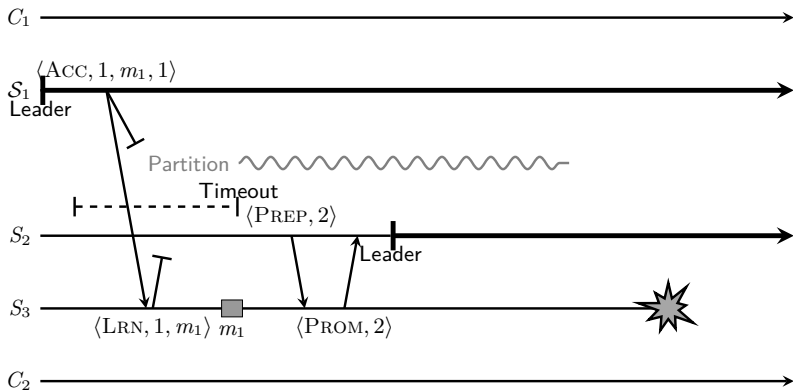
What Happens If S_3 Learn m_1 ?



What Happens If S_3 Learn m_1 ?



Does Leader Change Help?

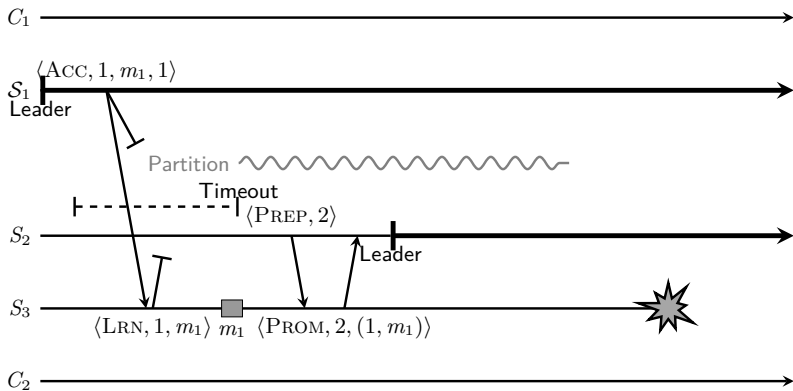


No!

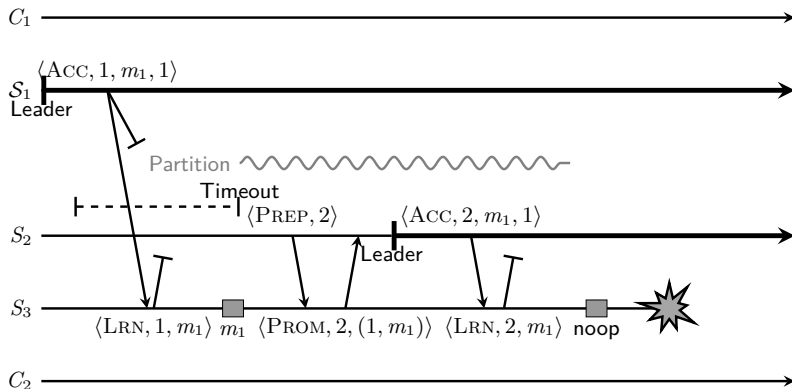
We Still don't Know What
 S_3 Did Before Crashing.

But the fix is Easy!

Tell new Leader About Accepted Messages

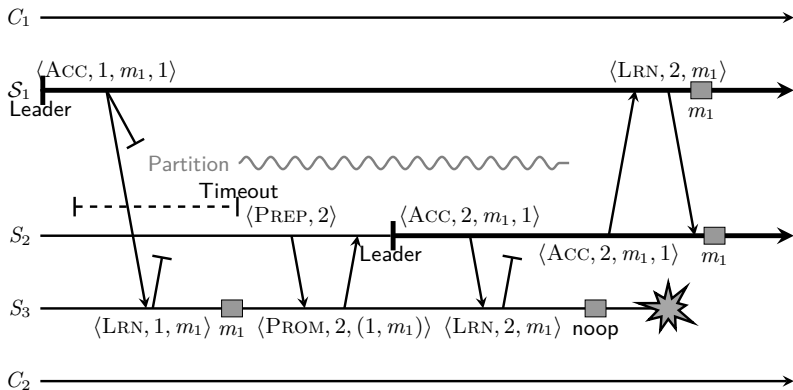


The new Leader Resends Accept for Those Messages

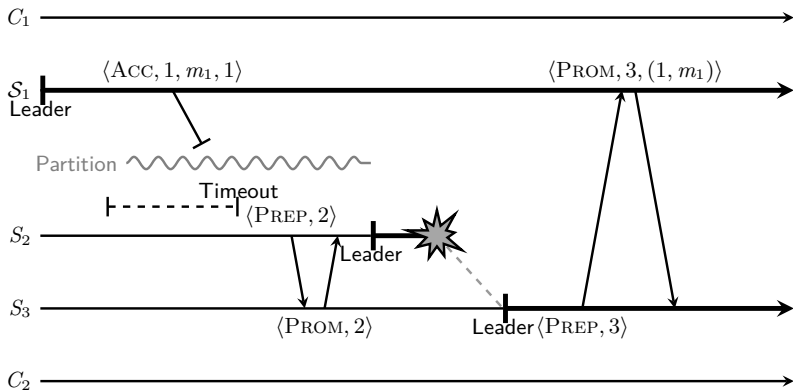


Learn was Lost and S_3 Crashed.
Leader Still can't Execute m_1 .

Leader Also Resends Accept After Merge



Promise from old Leader Includes Accepted Messages



Recap: Leader Change 2

- Added information about accept from previous leader:
 $\langle \text{PROM}, rnd, (1, m_1) \rangle$
 - Promise not to accept messages from a lower round than rnd
 - Last leader did send m_1 in round 1
 - Typical naming: $\langle \text{PROM}, rnd, (vrnd, vval) \rangle$

Recap: Leader Change 2

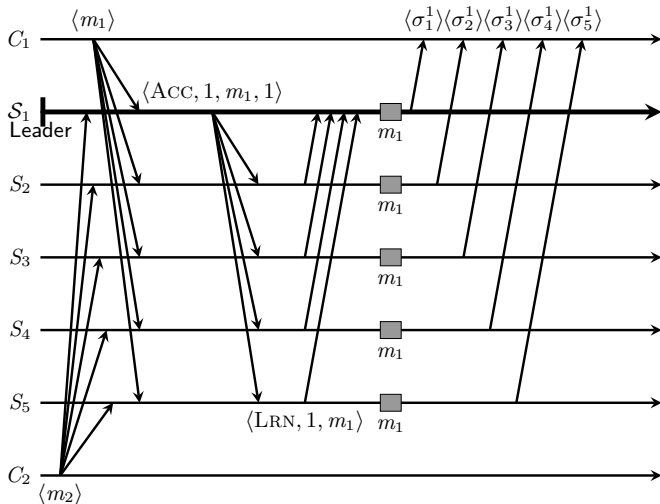
- Added information about accept from previous leader:
 $\langle \text{PROM}, rnd, (1, m_1) \rangle$
 - Promise not to accept messages from a lower round than rnd
 - Last leader did send m_1 in round 1
 - Typical naming: $\langle \text{PROM}, rnd, (vrnd, vval) \rangle$
- Leader resends accept for messages identified in the promise message
 - After receiving the promise
 - After a partition merge

What About More Than one Crash?

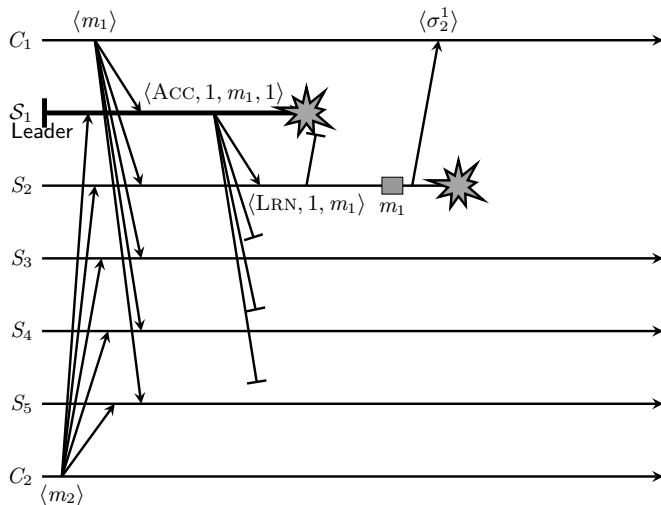
What About More Than one Crash?

- Increase the number of servers
- To limit progress to a majority partition:
 - We can only tolerate fewer than half of the servers fail
 - To tolerate f crashes, we need at least $2f + 1$

With Five Servers



With Five Servers, S_2 Cannot Execute After Accept



With Five Servers, S_2 Cannot Execute After Accept

- A combination of message loss and crashes
 - Prevent non-leader servers from executing after receiving an accept

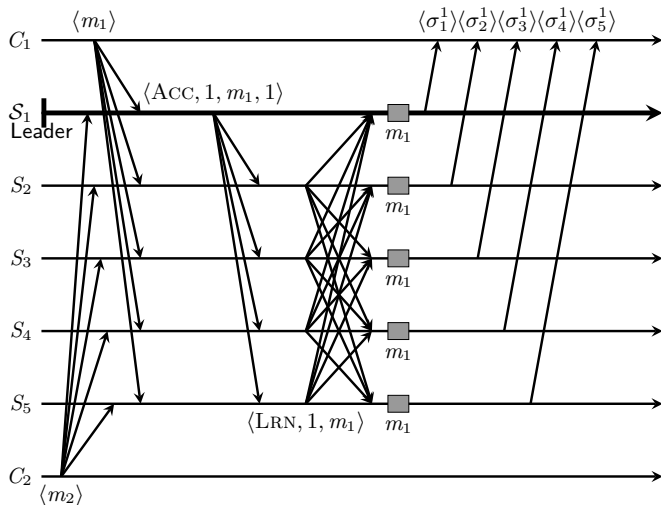
With Five Servers, S_2 Cannot Execute After Accept

- A combination of message loss and crashes
 - Prevent non-leader servers from executing after receiving an accept
 - This was not necessary for the three server case
 - The accept from the leader is an implicit learn
 - And together with its own "learn", can execute!

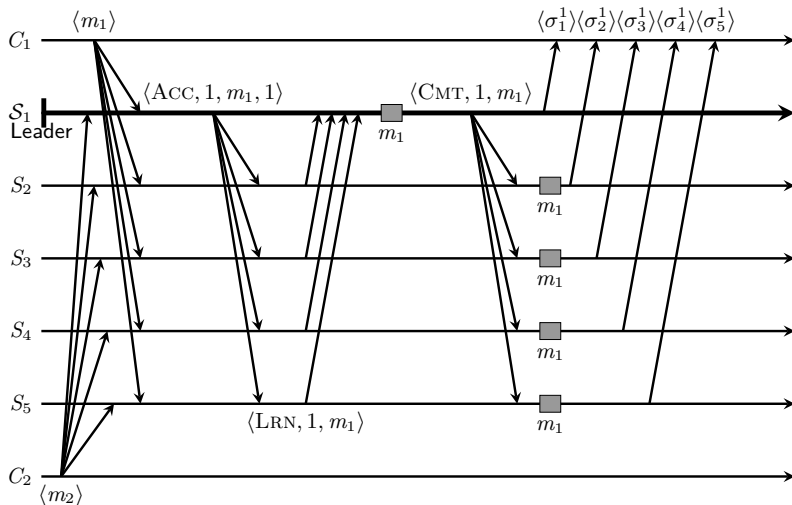
With Five Servers, S_2 Cannot Execute After Accept

- A combination of message loss and crashes
 - Prevent non-leader servers from executing after receiving an accept
 - This was not necessary for the three server case
 - The accept from the leader is an implicit learn
 - And together with its own "learn", can execute!
- There are two solutions:
 - Wait for all-to-all learn
 - Wait for commit from leader

All-to-All Learn Before Execute

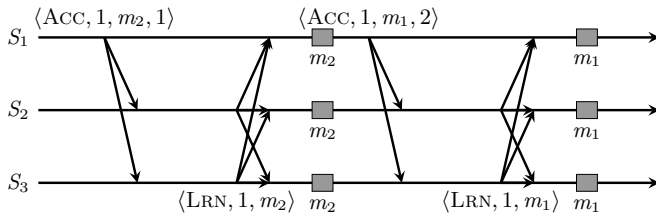


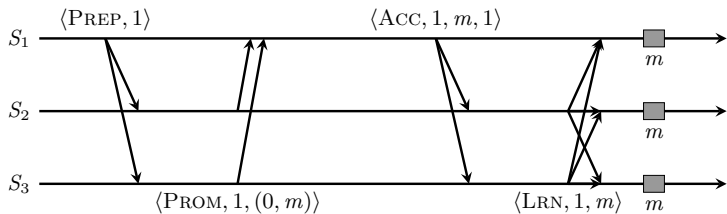
Await Commit Before Execute

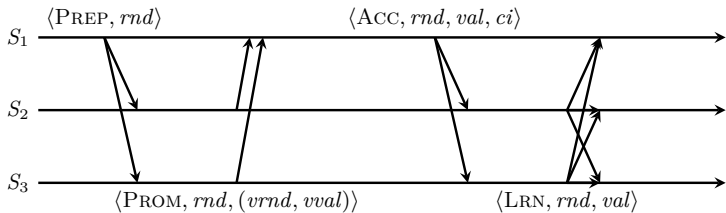


Wrapping it up!

Multi-Paxos







- Proposer = Leader
 - Sends prepare and accept messages
 - Receive promise messages
- Acceptor
 - Receive accept messages
 - Sends learn messages
- Learner
 - Receive learn messages

Agreement in Distributed Systems: Possible Solutions

- ▶ Two-Phase Commit (2PC)
- ▶ Paxos



Two-Phase Commit (2PC)

The Two-Phase Commit (2PC) Problem

- ▶ The problem first was encountered in **database systems**.
- ▶ Suppose a database system is updating some complicated data structures that include parts residing on **more than one machine**.

The Two-Phase Commit (2PC) Problem

- ▶ The problem first was encountered in **database systems**.
- ▶ Suppose a database system is updating some complicated data structures that include parts residing on **more than one machine**.
- ▶ System model:
 - **Concurrent processes** and uncertainty of **timing**, **order of events** and inputs (**asynchronous systems**).
 - **Failure** and recovery of machines/processors, of communication channels.

Intuitive Example (1/3)

- ▶ You want to organize outing with 3 friends at 6pm Tuesday.
 - Go out only if all friends can make it.



Intuitive Example (2/3)

- What do you do?

Intuitive Example (2/3)

► What do you do?

- Call each of them and ask if can do 6pm on Tuesday (voting phase)



Intuitive Example (2/3)

► What do you do?

- Call each of them and ask if can do 6pm on Tuesday (voting phase)
- If all can do Tuesday, call each friend back to ACK (commit)



Intuitive Example (2/3)

► What do you do?

- Call each of them and ask if can do 6pm on Tuesday (voting phase)
- If all can do Tuesday, call each friend back to ACK (commit)
- If one cannot do Tuesday, call other three to cancel (abort)



Intuitive Example (3/3)

- ▶ Critical details

Intuitive Example (3/3)

► Critical details

- While you were calling everyone to ask, people who have **promised** they can do 6pm Tuesday **must reserve that slot**.

Intuitive Example (3/3)

► Critical details

- While you were calling everyone to ask, people who have **promised** they can do 6pm Tuesday **must reserve that slot**.
- You need to **remember the decision** and **tell anyone** whom you have not been able to reach during commit/abort phase.

Intuitive Example (3/3)

► Critical details

- While you were calling everyone to ask, people who have **promised** they can do 6pm Tuesday **must reserve that slot**.
- You need to **remember the decision** and **tell anyone** whom you have not been able to reach during commit/abort phase.

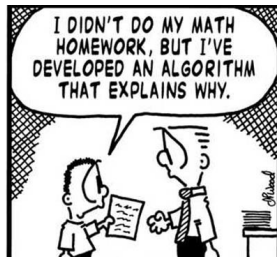
► That is exactly how 2PC works.

The 2PC Players

- ▶ **Coordinator** (Transaction Manager)
 - Begins transaction.
 - Responsible for commit/abort.
- ▶ **Participants** (Resource Managers)
 - The servers with the data used in the distributed transaction.

The 2PC Algorithm

- ▶ Phase 1 - prepare phase
- ▶ Phase 2 - commit phase



The 2PC Algorithm - Prepare Phase

- ▶ Coordinator asks each participant `canCommit`.

The 2PC Algorithm - Prepare Phase

- ▶ Coordinator asks each participant `canCommit`.
- ▶ Participants must prepare to `commit` using permanent `storage` before answering `yes`.

The 2PC Algorithm - Prepare Phase

- ▶ Coordinator asks each participant `canCommit`.
- ▶ Participants must prepare to `commit` using permanent `storage` before answering `yes`.
 - Lock the objects.
 - Participants are **not allowed** to cause an `abort` after it replies `yes` to `canCommit`.

The 2PC Algorithm - Prepare Phase

- ▶ Coordinator asks each participant `canCommit`.
- ▶ Participants must prepare to `commit` using permanent `storage` before answering `yes`.
 - Lock the objects.
 - Participants are **not allowed** to cause an `abort` after it replies `yes` to `canCommit`.
- ▶ Outcome of the transaction is `uncertain` until `doCommit` or `doAbort`.
 - Other participants might still cause an abort.

The 2PC Algorithm - Commit Phase

- ▶ The coordinator collects all votes.
 - If unanimous yes, causes commit.
 - If any participant voted no, causes abort.

The 2PC Algorithm - Commit Phase

- ▶ The **coordinator** collects **all votes**.
 - If **unanimous yes**, causes **commit**.
 - If **any participant voted no**, causes **abort**.
- ▶ The fate of the transaction is decided **atomically** at the **coordinator**, once all **participants** vote.
 - Coordinator records fate using **permanent storage**.
 - Then **broadcasts doCommit** or **doAbort** to participants.

2PC Sequence of Events

Coordinator

Participant

"prepared"

canCommit?

"prepared"
(persistently)

Yes

"committed"
(persistently)

doCommit

"uncertain"
(objects still
locked)

haveCommitted

"committed"

"done"

Recovery in 2PC

- ▶ Recovery after **timeouts**.
- ▶ Recovery after **crashes and reboot**.

Recovery in 2PC

- ▶ Recovery after **timeouts**.
- ▶ Recovery after **crashes and reboot**.
- ▶ Note: you **cannot differentiate** between the above in a realistic **asynchronous network**.

Handling Timeout

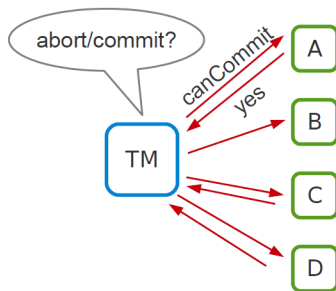
- ▶ To avoid processes blocking for ever.

Handling Timeout

- ▶ To avoid processes blocking for ever.
- ▶ Two scenarios:
 - Coordinator waits for votes from participants.
 - Participant is waiting for the final decision.

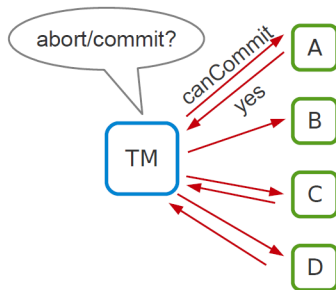
Handling Timeout at Coordinator

- ▶ If B voted no, can coordinator unilaterally abort?
- ▶ If B voted yes, can coordinator unilaterally abort/commit?



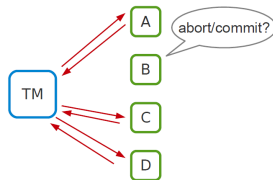
Handling Timeout at Coordinator

- ▶ If B voted no, can coordinator unilaterally abort?
- ▶ If B voted yes, can coordinator unilaterally abort/commit?
 - Coordinator waits for votes from participants.
 - Participant is waiting for the final decision.
 - Coordinator timeout abort and send **doAbort** to participants.



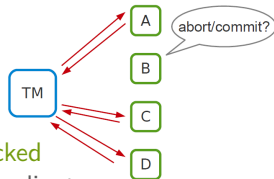
Handling Timeout at Participants

- If B times out on TM and has voted **yes**, then execute **termination protocol**.



Handling Timeout at Participants

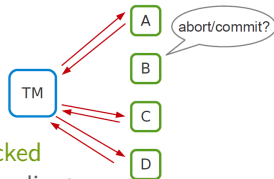
- ▶ If B times out on TM and has voted **yes**, then execute **termination protocol**.



- ▶ **Simple protocol**: participant is **remained blocked** until it can establish communication with coordinator.

Handling Timeout at Participants

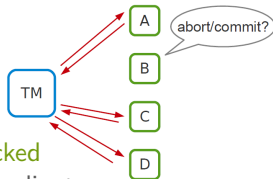
- ▶ If B times out on TM and has voted **yes**, then execute **termination protocol**.



- ▶ **Simple protocol**: participant is **remained blocked** until it can establish communication with coordinator.
- ▶ **Cooperative protocol**: participant sends a **decision-request** message to **other participants**.

Handling Timeout at Participants

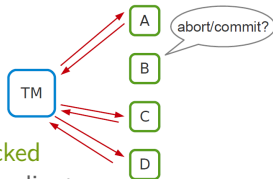
- ▶ If B times out on TM and has voted **yes**, then execute **termination protocol**.



- ▶ **Simple protocol**: participant is **remained blocked** until it can establish communication with coordinator.
- ▶ **Cooperative protocol**: participant sends a **decision-request** message to **other participants**.
 - B sends status message to A

Handling Timeout at Participants

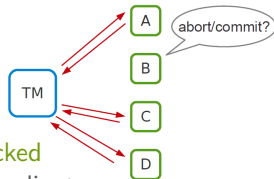
- ▶ If B times out on TM and has voted **yes**, then execute **termination protocol**.



- ▶ **Simple protocol**: participant is **remained blocked** until it can establish communication with coordinator.
- ▶ **Cooperative protocol**: participant sends a **decision-request** message to **other participants**.
 - B sends status message to A
 - If A has received commit/abort from TM, ...

Handling Timeout at Participants

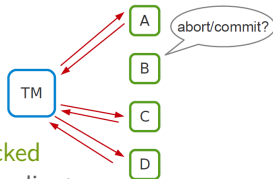
- ▶ If B times out on TM and has voted **yes**, then execute **termination protocol**.



- ▶ **Simple protocol**: participant is **remained blocked** until it can establish communication with coordinator.
- ▶ **Cooperative protocol**: participant sends a **decision-request** message to **other participants**.
 - B sends status message to A
 - If A has received commit/abort from TM, ...
 - If A has not responded to TM, ...

Handling Timeout at Participants

- ▶ If B times out on TM and has voted **yes**, then execute **termination protocol**.



- ▶ **Simple protocol**: participant is **remained blocked** until it can establish communication with coordinator.
- ▶ **Cooperative protocol**: participant sends a **decision-request** message to **other participants**.
 - B sends status message to A
 - If A has received commit/abort from TM, ...
 - If A has not responded to TM, ...
 - If A has responded with no/yes ...

Handling Crash and Recovery (1/2)

- ▶ All nodes must **log protocol progress**.
 - **Participants**: prepared, uncertain, committed/aborted
 - **Coordinator**: prepared, committed/aborted, done

Handling Crash and Recovery (1/2)

- ▶ All nodes must **log protocol progress**.
 - **Participants**: prepared, uncertain, committed/aborted
 - **Coordinator**: prepared, committed/aborted, done
- ▶ Nodes **cannot back out** if **commit is decided**.

Handling Crash and Recovery (2/2)

- ▶ **Coordinator** crashes:
 - If it finds **no commit** on disk, it **aborts**.
 - If it finds **commit**, it **commits**.

Handling Crash and Recovery (2/2)

► Coordinator crashes:

- If it finds **no commit** on disk, it **aborts**.
- If it finds **commit**, it **commits**.

► Participant crashes:

- If it finds **no yes** on disk, it **aborts**.
- If it finds **yes**, runs **termination protocol** to decide.

Fault-Tolerance Limitations of 2PC

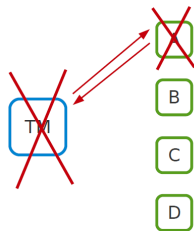
- ▶ Even with recovery enabled, 2PC is not really **fault-tolerant (or live)**, because it can be **blocked** even when one (or a few) machines **fail**.
- ▶ Blocking means that it does **not make progress during the failures**.

Fault-Tolerance Limitations of 2PC

- ▶ Even with recovery enabled, 2PC is not really **fault-tolerant (or live)**, because it can be **blocked** even when one (or a few) machines **fail**.
- ▶ Blocking means that it does **not make progress during the failures**.
- ▶ **Any scenarios?**

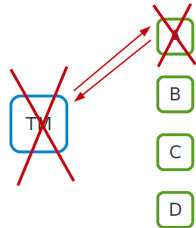
2PC Blocking Scenario

- ▶ TM sends `doCommit` to A, A gets it and commits, and then both TM and A die.



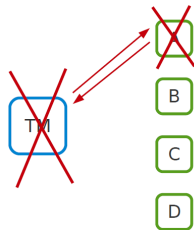
2PC Blocking Scenario

- ▶ TM sends **doCommit** to A, A gets it and commits, and then **both TM and A die**.
- ▶ B, C, D have already also replied **yes**, have **locked** their mutexes, and now need to wait for TM or A to reappear.
 - They **cannot recover** the decision with certainty until TM or A are online.



2PC Blocking Scenario

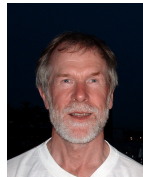
- ▶ TM sends **doCommit** to A, A gets it and commits, and then **both TM and A die**.
- ▶ B, C, D have already also replied **yes**, have **locked** their mutexes, and now need to wait for TM or A to reappear.
 - They **cannot recover** the decision with certainty until TM or A are online.
- ▶ This is why 2PC is called a **blocking protocol**:
2PC is safe, but not live.



Impossibility of Distributed Consensus with One Faulty Process

► Fischer-Lynch-Paterson (FLP)

- M.J. Fischer, N.A. Lynch, and M.S. Paterson, Impossibility of distributed consensus with one faulty process, Journal of the ACM, 1985.



FLP Impossibility Result

- ▶ It is **impossible** for a set of processors in an **asynchronous** system to **agree** on a binary value, even if only a **single** process is subject to an unannounced **failure**.
- ▶ The core of the problem is **asynchrony**.

FLP Impossibility Result

- ▶ What FLP says: you cannot guarantee both **safety** and **progress** when there is even a **single fault** at an inopportune moment.

FLP Impossibility Result

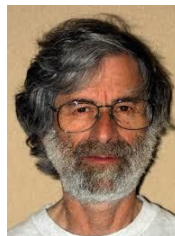
- ▶ What FLP says: you cannot guarantee both **safety** and **progress** when there is even a **single fault** at an inopportune moment.
- ▶ What FLP does not say: in practice, **how close can you get to the ideal** (always safe and live)?

FLP Impossibility Result

- ▶ What FLP says: you cannot guarantee both **safety** and **progress** when there is even a **single fault** at an inopportune moment.
- ▶ What FLP does not say: in practice, **how close can you get to the ideal** (always safe and live)?
- ▶ So, **Paxos** ...

Paxos

- ▶ The only known **completely-safe** and **largely-live agreement protocol**.
- ▶ L. Lamport, The part-time parliament, ACM Transactions on Computer Systems, 1998.



The Paxos Players

▶ Proposers

- Suggests values for consideration by acceptors.

▶ Acceptors

- Considers the values proposed by proposers.
- Renders an accept/reject decision.

▶ Learners

- Learns the chosen value.

The Paxos Players

▶ Proposers

- Suggests values for consideration by acceptors.

▶ Acceptors

- Considers the values proposed by proposers.
- Renders an accept/reject decision.

▶ Learners

- Learns the chosen value.

- ▶ A node can act as more than one roles (**usually 3**).

Single Proposal, Single Acceptor

- ▶ Use just **one acceptor**
 - Collects proposers' **proposals**.
 - Decides the value and tells everyone else.



Single Proposal, Single Acceptor

- ▶ Use just **one acceptor**
 - Collects proposers' **proposals**.
 - Decides the value and tells everyone else.



- ▶ Sounds familiar?

Single Proposal, Single Acceptor

- ▶ Use just **one acceptor**
 - Collects proposers' **proposals**.
 - Decides the value and tells everyone else.



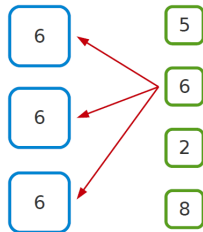
- ▶ Sounds familiar?
 - **two-phase commit (2PC)**
 - **acceptor fails = protocol blocks**

Single Proposal, Multiple Acceptors

- ▶ One acceptor is **not fault-tolerant** enough.
- ▶ Let's have **multiple acceptors**.

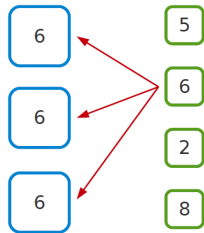
Single Proposal, Multiple Acceptors

- ▶ One acceptor is **not fault-tolerant** enough.
- ▶ Let's have **multiple acceptors**.



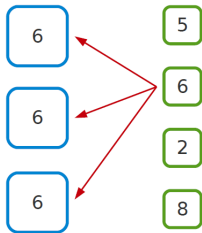
Single Proposal, Multiple Acceptors

- ▶ One acceptor is **not fault-tolerant** enough.
- ▶ Let's have **multiple acceptors**.
- ▶ From there, must reach a decision. **How?**



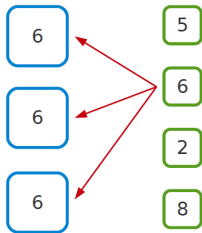
Single Proposal, Multiple Acceptors

- ▶ One acceptor is **not fault-tolerant** enough.
- ▶ Let's have **multiple acceptors**.
- ▶ From there, must reach a decision. **How?**
- ▶ **Decision** = value accepted by the **majority**.



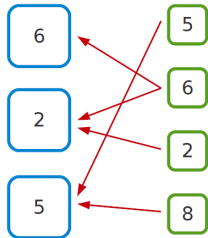
Single Proposal, Multiple Acceptors

- ▶ One acceptor is **not fault-tolerant** enough.
- ▶ Let's have **multiple acceptors**.
- ▶ From there, must reach a decision. **How?**
- ▶ **Decision = value accepted by the majority.**
- ▶ **P1: an acceptor must accept first proposal it receives.**



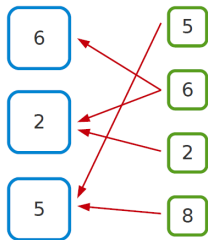
Multiple Proposals, Multiple Acceptors

- ▶ If there are **multiple proposals**, **no proposal** may get the **majority**.
 - 3 proposals may each get 1/3 of the acceptors.



Multiple Proposals, Multiple Acceptors

- ▶ If there are **multiple proposals**, **no proposal** may get the **majority**.
 - 3 proposals may each get 1/3 of the acceptors.



- ▶ **Solution:** acceptors can **accept multiple proposals**, distinguished by a **unique proposal number**.

Handling Multiple Proposals

- ▶ All chosen proposals must have the same value.

Handling Multiple Proposals

- ▶ All chosen proposals must have the same value.
- ▶ P2: If a proposal with value v is chosen, then every higher-numbered proposal that is chosen also has value v .

Handling Multiple Proposals

- ▶ All chosen proposals must have the same value.
- ▶ P2: If a proposal with value v is chosen, then every higher-numbered proposal that is chosen also has value v .
 - P2a: ... accepted ...

Handling Multiple Proposals

- ▶ All chosen proposals must have the same value.
- ▶ P2: If a proposal with value v is chosen, then every higher-numbered proposal that is chosen also has value v .
 - P2a: ... accepted ...
 - P2b: ... proposed ...

The Paxos Algorithm

- ▶ Phase 1a - prepare phase
- ▶ Phase 1b - promise phase
- ▶ Phase 2a - accept phase
- ▶ Phase 2b - accepted phase



Paxos Algorithm

Phase 1a: “Prepare”

Select proposal number* N and send a ***prepare(N)*** request to a majority of acceptors.

proposer

Phase 1b: “Promise”

If $N > \text{number of any previous promises or acceptances}$,
* promise to never accept any future proposal less than N ,
- send a ***promise(N, U)*** response
(where U is the highest-numbered proposal accepted so far (if any))

Phase 2a: “Accept!”

If proposer received promise responses from a majority,
- send an ***accept(N, V)*** request to those acceptors
(where V is the value of the highest-numbered proposal among the ***promise*** responses, or any value if no ***promise*** contained a proposal)

acceptor

Phase 2b: “Accepted”

If $N \geq \text{number of any previous promise}$,
* accept the proposal
- send an ***accepted*** notification to the learner

Paxos Algorithm - Prepare Phase

- ▶ A proposer selects a proposal number n and sends a prepare request with number n to majority of acceptors.

Paxos Algorithm - Promise Phase

- ▶ If an **acceptor** receives a **prepare** request with number n greater than that of any prepare request it saw

Paxos Algorithm - Promise Phase

- ▶ If an **acceptor** receives a **prepare** request with number n greater than that of any prepare request it saw
 - It responds **yes** to that request with a **promise** not to accept any more proposals numbered **less than** n .

Paxos Algorithm - Promise Phase

- ▶ If an **acceptor** receives a **prepare** request with number n greater than that of any prepare request it saw
 - It responds **yes** to that request with a **promise** not to accept any more proposals numbered **less than** n .
 - It includes the **highest-numbered proposal** (if any) that it has accepted.

Paxos Algorithm - Accept Phase

- ▶ If the proposer receives a response **yes** to its **prepare** requests from a **majority** of **acceptors**

Paxos Algorithm - Accept Phase

- ▶ If the **proposer** receives a response **yes** to its **prepare** requests from a **majority** of **acceptors**
 - It sends an **accept** request to each of those **acceptors** for a proposal numbered n with a values v , which is the value of the **highest-numbered proposal** among the responses.

Paxos Algorithm - Accepted Phase

- ▶ If an **acceptor** receives an **accept** request for a proposal numbered n

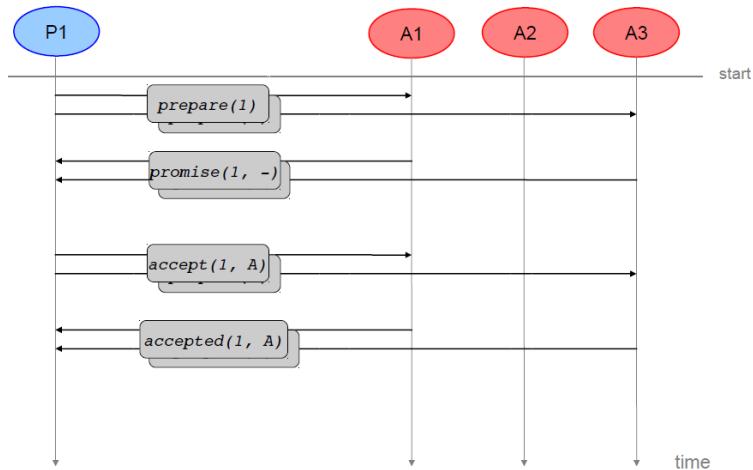
Paxos Algorithm - Accepted Phase

- ▶ 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 .

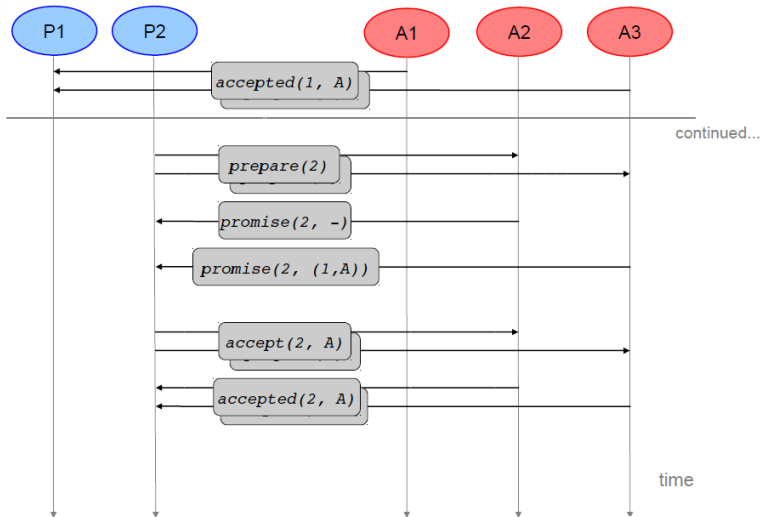
Definition of Chosen

- ▶ A value is **chosen** at proposal number n , iff **majority** of acceptors **accept** that value in phase 2 of the proposal number.

Paxos Example



Paxos Example

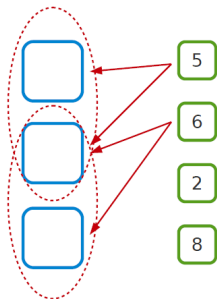


Paxos - Safety (1/3)

- ▶ If a value v is chosen at proposal number n , any value that is sent out in phase 2 of any later proposal numbers must be also v .

Paxos - Safety (2/3)

- ▶ Decision = Majority (any two majorities share at least one element)
- ▶ Therefore after the first round in which there is a decision, any subsequent round involves at least one acceptor that has accepted v .

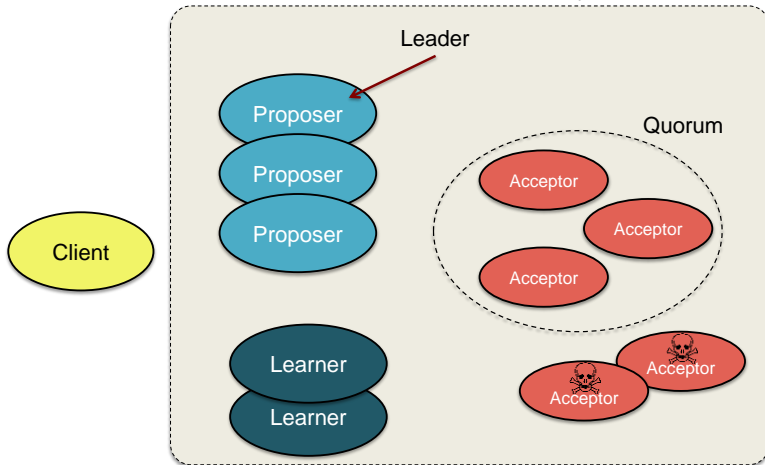


Paxos - Safety (3/3)

- ▶ Now suppose our claim is not true, and let m is the first proposal number that is later than n and in 2nd phase, the value sent out is $w \neq v$.
- ▶ This is not possible, because if the proposer P was able to start 2nd phase for w , it means it got a majority to accept round for m (for $m > n$). So, either:
 - v would not have been the value decided, or
 - v would have been proposed by P
- ▶ Therefore, once a majority accepts v , that never changes.

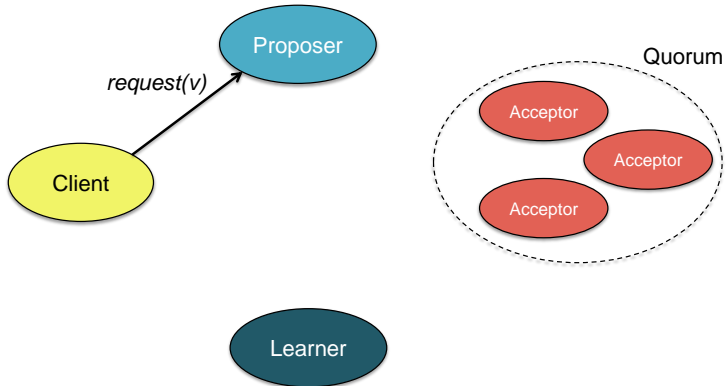
Paxos in action

Paxos nodes: one machine may serve several roles



Paxos in action: Phase 0

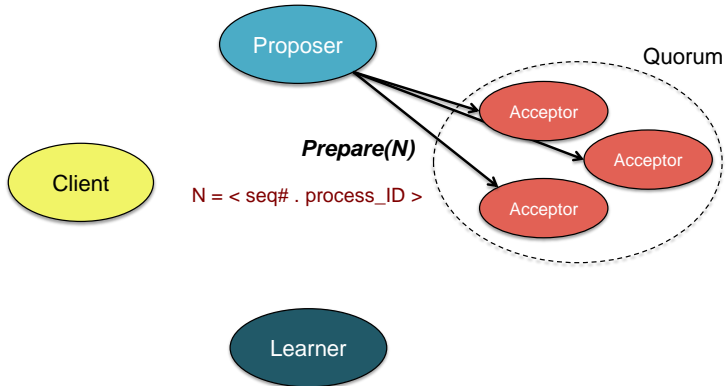
Client sends a request to a proposer



Paxos in action: Phase 1a

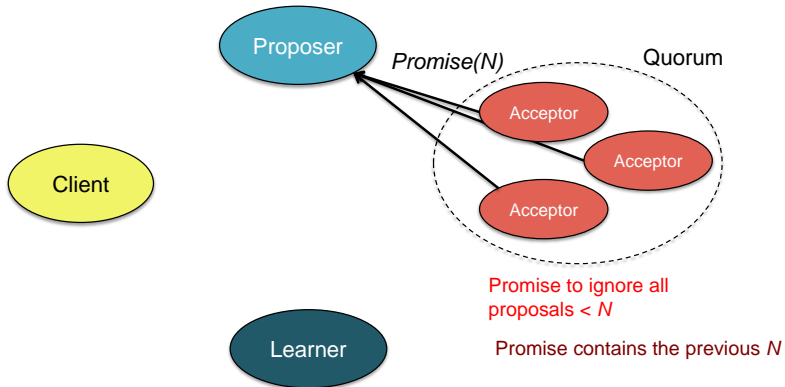
Proposer: creates a *proposal #N* (*N* acts like a Lamport time stamp), where *N* is greater than any previous proposal number used by this proposer

Send to Quorum of Acceptors (however many you can reach – but a majority)



Paxos in action: Phase 1b

Acceptor: if proposer's ID > any previous proposal
promise to ignore all requests with IDs < N
reply with info about highest past proposal: { N, value }



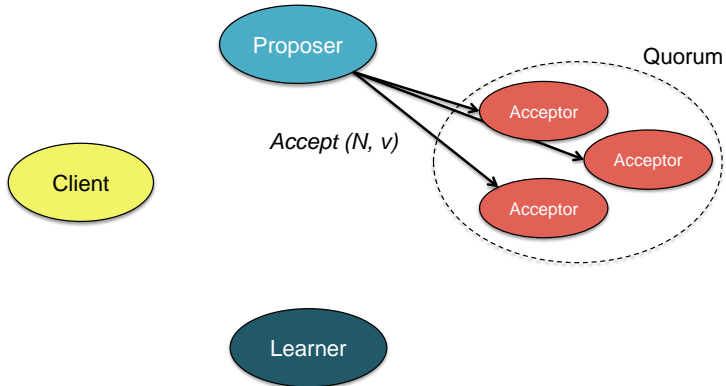
Paxos in action: Phase 2a

Proposer: if proposer receives promises from the quorum (majority):

Attach a value v to the proposal (the event).

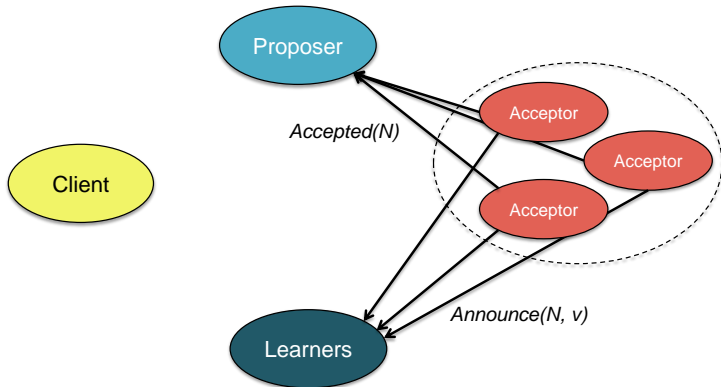
Send **Accept** to quorum with the **chosen** value

If promise was for another $\{N, v\}$, proposer **MUST** accept that



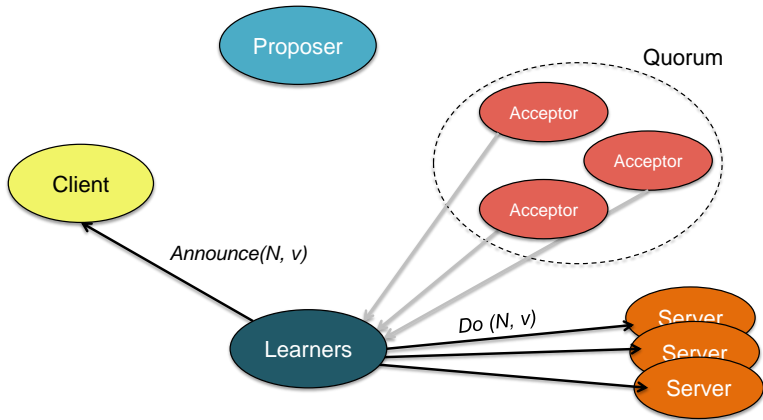
Paxos in action: Phase 2b

Acceptor: if the promise still holds, then announce the value v
Send **Accepted** message to Proposer and every Learner
else ignore the message (or send *NACK*)



Paxos in action: Phase 3

Learner: Respond to client and/or take action on the request



Paxos: Keep trying

- A proposal N may fail because
 - The acceptor may have made a new promise to ignore all proposals less than some value $M > N$
 - A proposer does not receive a quorum of responses: either *promise* (phase 1b) or *accept* (phase 2b)
- Algorithm then has to be restarted with a higher proposal #