

Mechanised Verification of Paxos-like Consensus Protocols

Anirudh Pillai

February 3, 2018

Mechanised Verification of Paxos-like Consensus Protocols

ABSTRACT

Distributed systems are hard to reason about.

Acknowledgments

Acknowledgements

Contents

1	INTRODUCTION	2
1.1	The Problem	2
1.2	Aims and Goals	2
1.3	Project Overview	3
1.4	Report Overview	3
2	BACKGROUND	4
2.1	Distributed Systems	4
2.2	Paxos	8
2.3	Disel	12
2.4	Related Work	12
3	REQUIREMENTS AND ANALYSIS	13
3.1	Detailed Problem Statement	13
3.2	Requirements	13
3.3	Analysis	14
4	DESIGN AND IMPLEMENTATION	17
4.1	Modelling	17
4.2	Verification	25
5	CLIENT APPLICATION	26
6	CONCLUSION AND EVALUATION	27
6.1	Summary of Achievements	27
6.2	Critical Evaluation	27
6.3	Future Work	27
6.4	Final Thoughts	27
	REFERENCES	28
	APPENDIX A INTERIM REPORT	29

1

Introduction

1.1 The Problem

1.2 Aims and Goals

I have highlighted the aims and goals separately. The aims are what I want to achieve out of undertaking this project and the goals are the things that this project tries to achieve.

1.2.1 Aims

1. Learn about distributed system protocols
2. Contribute to open source software

1.2.2 Goals

1. Read about and understand the classical Paxos-like consensus algorithms.
2. Develop state transition systems for the algorithms and identify the invariants that need to be preserved during the operation of the algorithm.
3. Implement a simulation of the protocols in Python.
4. Formulate the implemented protocols in Disel by using the developed state-transition systems.

5. Mechanise the proofs of the identified protocol invariants in Disel/Coq.
6. Add additional communication channels and prove composite invariants.
7. Provide an abstract specification of the protocol, usable by third-party clients.
8. Mechanise a client application of the protocol verified out of the abstract interface.

1.3 Project Overview

1.4 Report Overview

2

Background

This chapter lays down all the previous research which the project builds on. Before going over the design decisions on the project we first need to understand this background information and look at related work to see different approaches used to solve the problem.

2.1 Distributed Systems

A distributed system is a model in which processes running on running different computers, which are connected together in a network, exchange messages to coordinate their action, often resulting in the user thinking of the entire system as one single unified computer.

A computer in the distributed system is also alternatively referred to as a processor or a node in the system. Each node in a distributed systems has its own memory.

We will now go over a few concepts of distributed systems which will help us understand the characteristics of the protocols that run on these systems. This will lay down the groundwork for us to understand the Paxos protocol on which this project is based.

2.1.1 Asynchronous Environment

An asynchronous distributed system is one where there are no guarantees about the timing and order in which events occur.

The clocks of each of the process in the system can be out of sync and may not be accurate. Therefore, there can be no guarantees about the order in which events occur.

Further, messages sent by one process to another can be delayed for an arbitrary period of time.

A protocol running in an asynchronous environment has to account for these conditions in its design and try to achieve its goal without the guarantees of timed events. An asynchronous environment is very common for a real world distributed system but it also makes reasoning about the system harder because of the aforementioned properties.

2.1.2 Fault Tolerance

A fault tolerant distributed system is one which can continue to function correctly despite the failure of some of its components. A 'failure' of a node or 'fault' in a node means any unexpected behaviour from that node eg. not responding to messages, sending corrupted messages.

Fault tolerance is one of the main reasons for using a distributed system as it increases the chances of your application continuing to functioning correctly and makes it more dependable. As Netflix mention on their blog 'Fault Tolerance is a Requirement, Not a Feature'. With their Netflix API receiving more than 1 billion requests a day, they expect that it is guaranteed that some of the components of their distributed system will fail. Using a fault tolerant distributed system they are able to ensure that a small failure in some components doesn't hinder the performance of the overall system, hence, enabling them to achieve their uptime metrics.

Fault tolerant distributed system protocols are protocols which achieve their goals despite the failure of some of the components of the distributed system they run on. The protocol accounts for the failures and generally specifies the maximum number of failures and the types of failures it can handle before it stops functioning correctly.

2.1.3 State Machine Replication

For a client server model, the easiest way to implement it is to use one single server which handles all the client request. Obviously this isn't the most robust solution as if the single server fails, so does your service. To overcome the problem you use a collection of servers each of which is a replica of the original single server and ensure that each of these 'replicas' fails independantly, without effecting the other replicas. This adds more fault tolerance.

State Machine Replication is method for creating a fault tolerant distributed system by replicating servers and using protocols to coordinate the interactions of these replicated servers with the client.

A State Machine M can be defined as $M = \langle q_0, Q, I, O, \delta, \gamma \rangle$ where

q_0 is the starting state

Q is the set of all possible states. I is set of all valid inputs O is the set of all valide outputs δ is the state transition function, $\delta : I \times Q \rightarrow Q$ γ is the output function, $\gamma : I \times Q \rightarrow O$

The state machine begins in the start state and transitions to other states and produces outputs when it receives the inputs. The transition and output are found using the transition and output functions. A deterministic state machine is one whose state transition and output functions are injective, i.e. multiple copies of the machine when given the same input, pass through the same order of states and produce the same output in the same order.

The method of modelling a distributed system protocol as state transition system is very common and is a critical component of this project as we will see soon when we need to encode our protocol in Disel.

State machine replication involves modelling our single server, from the client server model, and using multiple copies (replicas) of the same deterministic state machine and providing all of them with the input from the client. As long as one of the replicas does not crash, while resolving the request, we can successfully return a response to the client.

2.1.4 Consensus Protocols

For handling faults in your distributed system you need to have replication. This leads to the problem of making all these replicas agree with each other to keep them consistent. Consensus protocols try to solve this problem.

Consensus protocols are the family of distributed systems protocols which aim to make a distributed network of processes agree on one result.

These protocols are of interest because of their numerous real world applications. Let us take the example of a distributed database, which is a critical part of almost all large scale real world applications. This distributed database will run over a network of computers and everytime you use the database you aren't guaranteed to be served by the same computer.

Suppose you add a file to the database. This action is performed by the processor that was serving you 'add' request. Later when you want to retrieve the file from the database you might be served by a different computer that did not perform the 'add' request. In-order for the new computer to know that the file exists in the database, you will need to use a consensus protocol which helps all the computers in the network (which handle user requests) agree upon the result that the file has been added to the database.

Popular consensus protocols include PageRank used by Google and the Blockchain consensus protocol. George and Ilya verified a subset of the protocol in Coq in their Toychain paper.

2.2 Paxos

Having understood the the main concepts behind distributed system protocols, we can now finally get to the protocol at the heart of this project. Paxos is a family of asynchronous, fault tolerant, consensus protocol which achieves consensus in a network of unreliable processes as long as a majority of them don't fail.

Paxos is used for state machine replication. Once you have multiple replicas servicing client requests, how do you makes sure that all of these replicas agree on what action to take? The solution is simply to use a consensus protocol like Paxos to make all replicas agree on something.

Paxos has many variants but the one we will focus on is the one we actually prove in Disel, single decree Paxos, also know as simple paxos. Simple Paxos is an algorithm that helps a distributed network of processors to achieve consensus. Consensus is achieved when the network of processor agree on a common value.

For simple paxos, we assume the following assumptions hold about the processors and the environment, in order for the protocol to function correctly.

- Processors communicate between each other by exchanging asynchronous messages between each other.
- Processors run at an arbitrary speed and may fail or restart. Handling this relates to the fault tolerant nature of paxos. Also, we assume that Byzantine faults don't occur. This means that all processors actually work together to try to achieve consensus on a value. There are variants of paxos which can also handle Byzantine failure buy not simple paxos. (This can be linked to the 'PBFT' paper which states that any algorithm handling Byzantine faults must have three phases. Simple paxos only has two phases.)

As for fault tolerance of paxos, in order to handle a failure of upto f processors, we need to have a minimum $2f + 1$ processors participating in the algorithm. This means paxos

functions correctly as long as a majority of the processors in the network do not fail. We will see shortly why just a majority needs to function correctly.

A processor participating in simple paxos, may have one or more of these three different roles - proposer, acceptor or learner.

- **Proposer** - A process acting as a proposer listens for client request and proposes a value which the network of processes tries to agree upon.
- **Acceptor** - acceptors receive proposed values from the proposers and then respond to them stating whether they are in a position to accept the value or not. For a proposed value to be accepted, a majority of all the existing acceptors have to accept the proposed value.
- **Learner** - The learner has to be informed when an acceptor accepts a value. The learner can then figure out when consensus has been achieved by calculating when a majority of acceptors have accepted the same proposal. Once the acceptors agree on a value, the learner may act on the value eg. Send request to client informing them about the agreed value.

2.2.1 Choosing a Value

For passing around the value to be chosen from one processor to the other, a processor must send a 'proposal' to the other processor. You can think of a proposal as just a tuple $\langle n, v \rangle$. n is just a natural number associated with a proposal which makes it easy to keep track of all the different proposals.

A *quorum* of acceptors is a subset of the set of all acceptors with length greater than $N/2$ where N is the length of the set of acceptors. A *quorum* is just a set denoting a majority of all the available acceptors.

Consensus is achieved when a proposal is accepted by a majority of acceptors.

THE ALGORITHM

Simple paxos runs in rounds until consensus is achieved (a successful round has occurred, where a majority of acceptors have accepted a proposal). A successful round of the algorithm has two phases, each of which can be subdivided into parts a, b.

- **Phase 1a: Prepare Request.** A proposer sends a proposal $\langle n, v \rangle$ to each acceptor in any randomly chosen *quorum* of acceptors. This first message that the proposer sends out is called a *prepare request*. As it the proposer tries to 'prepare' the acceptors to 'accept' a value in the future.
- **Phase 1b: Promise Response.** An acceptor on receiving a prepare request, responds with a *promise response*, if and only if the acceptor has not already sent a promise response with a proposal containing a proposal number n' where $n' > n$.

A promise response for proposal $\langle n, v \rangle$ is basically a guarantee (a 'promise') that this acceptor will not respond to any messages with proposals that have a proposal number n' where $n' < n$.

Thus, if an incoming prepare request has proposal number less than what the acceptor has already promised earlier, then the acceptor can ignore this prepare request by not responding to it. Although, for speeding up the protocol, the acceptor can send out a *nack response* which tells the proposer to stop trying to achieve consensus with this proposal.

If the acceptor has not sent any promise response before, then the body of the promise response can be empty, otherwise the acceptor must include the last proposal that it promised (before the current one) in the body of the message.

- **Phase 2a: Accept Request.** If the proposer successfully receives promise responses from a majority of acceptors, then it can send out an *accept request*. A accept request is a message containing a proposal which tells an acceptor to accept this proposal if it can.

The proposer creates a new proposal, $\langle n, v' \rangle$ where n is the same as in the proposal which the proposer sent in its prepare request. But, v' is the value from the highest

numbered proposal, selected from all the proposals that the proposer receives in the promise responses. If none of the promise responses received by the proposer contain a proposal, the proposer is free to set v' to any value it likes. The proposer then sends this accept request with proposal $\langle n, v' \rangle$ to another quorum of acceptors.

- **Phase 2b: Accepted Response.** Any acceptor that receives the accept request with proposal $\langle n, v \rangle$ responds with an accepted response if and only if it hasn't already promised not to respond to any proposals with proposal number n' where $n' > n$.

2.2.2 Informing learner

When consensus is achieved, a learner must be informed that a majority of acceptors have agreed on a value. There are various ways to do this.

1. Whenever an acceptor accepts a value, it should send the accepted proposal to all the learners. The learner will then know when a majority of acceptors have accepted the same value.
2. We can have a distinguished learner which informs other learners about the chosen value. The acceptors only need to inform this particular learner when they accept a value. This reduces number of messages sent but the distinguished learner becomes the single point of failure and also requires an additional round of sending messages where the distinguished learner informs other learners that a value has been chosen.
3. We can use a set of distinguished learners. The acceptors inform these distinguished learners who then inform the other learners. This increases reliability but also increases the number of messages exchanged.

2.3 Disel

2.3.1 Inductive Invariant

2.4 Related Work

3

Requirements and Analysis

3.1 Detailed Problem Statement

3.2 Requirements

1. Adapt Paxos for encoding in Disel and devise the state-transition system for this protocol.
2. Develop an inductive invariant for the adapted protocol that ensures the protocol functions correctly by imposing requirements on the global state of the system
3. Implement a simulation of the adapted protocol with the developed state transition system.
4. Mechanise the proof of the adapted protocol in Disel/Coq. Thereby, providing a library of reusable verified distributed components.
5. Mechanise a client application of the protocol verified out of the abstract interface.

3.3 Analysis

3.3.1 Requirement 1: Adapted Protocol and State Transition System

In order to mechanise the proof of Paxos in Disel, we had to first adapt the protocol in order to simplify the proof. There are many variants of Paxos and we had to study the we decided to focus on single decree paxos, the variant that was first proposed by Leslie Lamport

We studied the protocol in detail and also decided to focus on proving the part of the protocol that deals with achieving consensus, which is also the main function of the protocol. For this reason we did not include the learner in our client application or adapted protocol, nor did we focus on the part where the chosen value is learnt by all the nodes. We instead let our inductive invariant handle the case to detect when consensus had been achieved as the inductive invariant can impose requirements on the global system state.

Additionally, we also had to create a state transition system for the nodes in the protocol as Disel relies on this to impose pre and post conditions on the states of a node. In Paxos each node can have different roles but we had to split up each role into different states depending on the current data held by the node and the current function of the node in the protocol. We decided on the states each node could be in and how and when it transitions between them. This helped us come up with precondition and postcondition for the state of each node when it transitions on receiving or sending a message. We tried to minimise the number of transitions and the data held in each node's state in order to simplify the proof in Disel.

We look at these in more detail in the next chapter when we talk about modelling the protocol.

3.3.2 Requirement 2: Inductive Invariant

We also had to come up with an inductive invariant for the protocol such that if the inductive invariant holds in some state then it holds in every state reachable from that state. The inductive invariant was critical as it helped ensure that the protocol functions correctly by imposing requirements on the global state of the system. For proving the correctness of Paxos we found that our invariant had to capture when consensus is achieved on a value and also that once consensus is achieved on a particular value, further rounds of the protocol don't change this value. We then also came up with a proof for how this inductive invariant holds in our adapted protocol.

3.3.3 Requirement 3: Simulation

I also needed to implement a simulation of our adapted protocol. The simulator must be based on a state-transition system like Disel and should be able to simulate different nodes in the distributed system. The main reason for implementing this simulation is to be sure that our adapted protocol, designed in Requirement 1, will actually be provable in Disel. The simulation will enable us to detect errors and fix them much faster than having to fixing them in the middle of the Disel proof which is a much more time consuming to implement. The simulator will be implemented with the same state transitions we decided upon in Requirement 1, thus, the correct working of the simulator will give us confidence that our state transition system for Paxos will work correctly in Disel.

I decided to use Python to implement the simulation as that was the language I was strongest in. I studied how the simulation for Multi Paxos was implemented in the Paxos made moderately complex paper, which helped me learn how to create separate process for each node and also how to communicate by exchanging messages between the processes.

3.3.4 Requirement 4: Proof

Having designed the state transition diagram and the inductive invariant, the task of mechanising the proof in Disel becomes much easier. For implementing the proofs, I needed to learn more about Coq and SSReflect. I had to study examples of protocols proved in Disel, like the proof of the Two Phase Commit protocol.

3.3.5 Requirement 5: Client Application

After studying the Disel paper and looking at similar examples, I implemented the core of the adapted protocol in Disel. I also implemented a client application in Disel. The pre and post conditions from the state transition system helped me to implement the client application in such a way to adhere with the main protocol. Using the extraction feature in Disel and the shims runtime, I successfully extracted a working program of the client application in OCaml.

4

Design and Implementation

4.1 Modelling

4.1.1 The Protocol

4.1.2 State Transitions

Having adapted the protocol, we then had to create a state transition system for the nodes in your protocol in order to encode it in Disel. For creating the states, we need to look at what the function of each node is in the protocol at a particular moment and what type of data it holds at that time.

A node should only be able to transition from one state to another when it either receives or sends a message. Therefore, the data held by the node in each state should be enough for it to be able to create the message it wants to send or to be able to correctly process the message it receives.

We tried to minimise the number of states and transitions between them, in order to simplify the proof in Disel. This was important because each state transition has to be shown to hold with the invariant so reducing the state transitions, reduces the number of proofs.

We decided that a node can either be initialised as an acceptor or a proposer. The state transitions of each node will depend upon this initial state, so below we will separately look

at the state transition systems for the proposer and the acceptor.

The main difference between the state transition for the acceptor and the proposer is that the acceptor sends and receives messages from a single proposer while a proposer has to send and receive messages from all the acceptors.

PROPOSER

The proposer starts off in the `PInit` state where it is initialised with a `proposal` (a custom defined data type which is tuple of two natural numbers), $\langle p, v \rangle$. The natural number p is the proposal number and v is the value that the proposer tries to achieve consensus on. This means that the first prepare request this proposer sends will this proposal $\langle p, v \rangle$.

The proposer then moves to the `PSentPrep` state when it starts to send prepare requests to the acceptors. In this state, the proposer still holds the proposal but additionally now also stores a list of natural numbers, `sent_to`. This list stores the natural number identifiers of the acceptors, this proposer has sent requests to. Whenever it sends a prepare request to an acceptor, it adds the identifier of the acceptor to this list. The proposer remains in the `PSentPrep` state and keeps sending prepare requests until the contents of `sent_to` become equal to the global list `acceptors` which holds the identifiers of all the acceptors in the system. This means the proposer stays in this state until it has sent a prepare request to every single acceptor in the system.

Once the proposer has sent the last prepare request, it then transitions to the `PWait-PrepResp` state. In this state the proposer again holds a proposal and another list `promises` which is defined as below to be a list of tuples each containing a `nid` (a natural number identifier for a node), a boolean and a `proposal`.

Definition `promises := seq (nid * bool * proposal)`.

The proposer stays in this states and keeps receiving messages from the acceptor until one of the following two things happen:

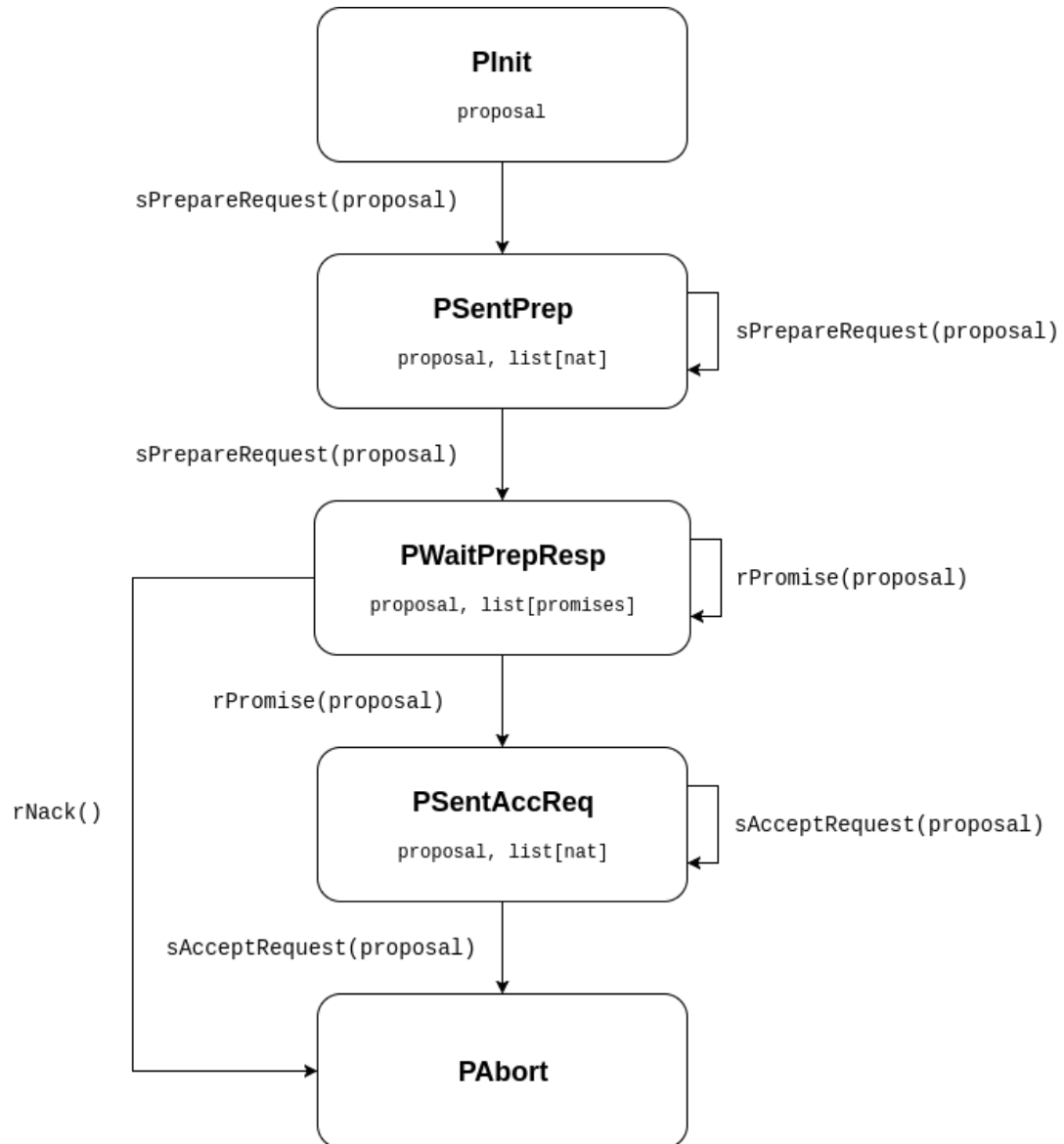


Figure 4.1: Proposer State Transition Diagram

1. It receives a nack response from the acceptor. This indicates that the acceptor might already have promised a proposal with a proposal number greater than p . This leads to the proposer to transition into the `PAbort` state. In this state the proposer basically gives up trying to achieve consensus using the proposal number p that it was initialised with and completely stops sending and receiving messages. Hence, the proposer doesn't need to hold any data in this state.
2. It receives a promise response from every single acceptor. When this happens, the proposer transitions to the `PSentAccReq` state.

When the proposer reaches the `PSentAccReq`, it means it has received a promise from every single acceptor and it can now start sending accept requests to each of the acceptors in the system. In the `PSentAccReq` the proposer again stores a list `sent_to` to keep track of every single acceptor it has already sent the accept request to. It also stores another `proposal` which has the same proposal number p that the proposer was initialised with but the value v is the value from the highest numbered proposal it received in a promise response. In the verification section, we will look at how it determines this value by looping over the `promises` list from the `PWaitPrepResp` state. The sending of accept requests works similar to sending prepare requests in the `PSentPrep` state. Finally, when the proposer finishes sending the accept requests to all the acceptors, it transitions to the `PAbort` state where it stops sending and receiving messages.

ACCEPTOR

The Acceptor starts off in the `AInit` state. It doesn't hold any data in this state as it is not sending any messages. It keeps listening for messages and on receiving a prepare request message, it transitions to `APromised` state.

In the `APromised` state, the acceptor holds a `proposal`. This is the highest numbered proposal that it has received so far in a prepare request message. In this state, on receiving

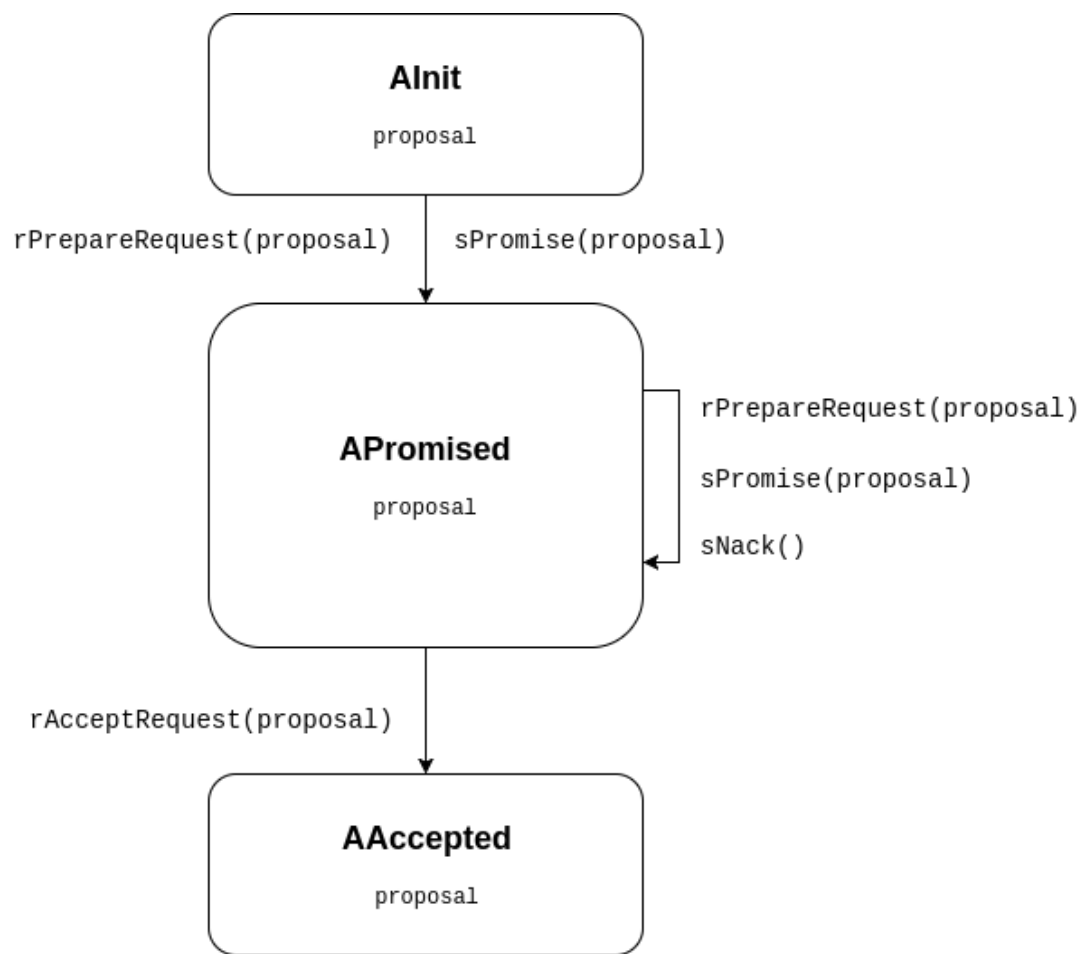


Figure 4.2: Acceptor State Transition Diagram

a prepare request, if the proposal number of the proposal in the prepare request is greater than the proposal number of the proposal it currently holds, it updates its current state to hold the new proposal but still remains in the `APromised` state. If the proposal number of the proposal in the prepare request is not greater, the acceptor sends a nack response to the proposer who sent the prepare request and does not update its state.

In the `APromised`, on receiving an accept request, if and only if the value of the proposal number proposal in the accept request is greater than the proposal number of the proposal that it currently holds, it transitions to the `AAccepted` state where it now holds the new proposal with the greater proposal number. If the proposal number of the new proposal number is not greater, the acceptor remains in the same state.

In the `AAccepted` state, the acceptor stops listening for and responding to messages. This is similar to the `PAbort` state for the proposer.

4.1.3 Inductive Invariant

A critical part of proving our protocol in Diesel was designing an inductive invariant for our adapted protocol. The inductive invariant helps ensure the correctness of our adapted protocol enabling us to imposing requirements on the global state of the system. For proving the correctness of paxos we found that our invariant had to capture when consensus is achieved on a value and also that once consensus is achieved on a particular value, further rounds of the protocol don't change this value.

Inductive invariant means a property which when it holds for a state s , it will hold for any state s' reachable from s .

The crux of Paxos' correctness lies in the prepare phase where, before sending the accept request, the Proposer must first set the value of the proposal, that it wants to propose, to be the value of the highest numbered proposal it receives as a promise. This ensures that when consensus has been achieved on a value 'v', further rounds of the protocol also ensure that

consensus will only be achieved on 'v'.

We established two invariants **I1** and **I2** which together form an inductive invariant for our protocol that also proves its safety.

- **I1** simply tries to say that there can only be one unique value associated with a particular proposal number for any proposal that has been accepted.
- **I2** states that once consensus has been achieved on a value v, every higher number proposal accepted by an acceptor also has the value v.

The mathematical representations for the invariants is given by.

- **I1** - $\forall a_i, a_j \in A, \langle p_i, v_i \rangle \in a_i.\text{accepted}, \langle p_j, v_j \rangle \in a_j.\text{accepted} \rightarrow v_i = v_j$.
- **I2** - $\forall \langle p_i, v_i \rangle, \forall a_j \in A, \exists \langle p_j, v_j \rangle \in a_j.\text{accepted}, p_j > p_i \rightarrow v_i = v_j$

I1 is preserved because if there are n proposers, they are initialised with a unique proposal numbers and throughout the running of our adapted protocol, the proposer always uses this unique proposal number for any value that it proposes. Hence, two different proposers never propose a proposal with the same proposal number. Additionally, each proposer only sends one round of accept requests with the same proposal. So as each proposer proposes only one value with a unique proposal number, we can deduce that each accepted proposal will have a unique value associated with a particular proposal number.

Once consensus has been achieved on a value, further runs of the algorithm don't change the value on which consensus has been achieved. We need to show that once consensus has been achieved on a proposal with value v then every other proposal, with a higher proposal number, on which consensus is achieved will also have proposal value set to v.

In order for consensus to be achieved on a new proposal, the new proposal first needs to be accepted by an acceptor.

\Rightarrow If consensus has been achieved on a proposal $\langle p_1, v_1 \rangle$ then every other proposal $\langle p_2, v_2 \rangle$ accepted by any acceptor, where $p_2 > p_1$, has $v_2 = v_1$.

Further, acceptors can only accept a proposal which has been proposed by a proposer. So we can reduce the requirement as follows.

\Rightarrow If consensus has been achieved on a proposal $\langle p_1, v_1 \rangle$ then every accept request $\langle p_2, v_2 \rangle$ sent by the proposer with $p_2 > p_1$, has $v_2 = v_1$.

In order to prove the above, let's assume that consensus has been achieved on a proposal $\langle p_1, v_1 \rangle$.
(4.1)

After that let's say that the system achieves consensus on $\langle p_2, v_2 \rangle$ where $p_2 > p_1$ and there does not exist p_x such that consensus has been achieved on a proposal with proposal number p_x where $p_1 < p_x < p_2$.

So from our assumption (4.1), there must be a majority of acceptors such that they have accepted the proposal $\langle p_2, v_2 \rangle$. So we need to show that $v_2 = v_1$. This is ensured in Paxos because of Phase 1 where the proposer must first get promises from a majority.

So any majority the proposer for p_2, v_2 gets in Phase 1, will have at least one acceptor a which has accepted $\langle p, v \rangle$. Paxos also ensures that before sending the accept request for p_2, v_2 , the proposer must select the value of the highest numbered proposals which it receives in its promises.

So the Acceptor a will send $\langle p_1, v_1 \rangle$ in its promise message to the proposer. As $\langle p_2, v_2 \rangle$ is the only proposal number which has proposal number greater than p , the proposer must set $v_2 = v_1$ in its accept request message $\langle p_2, v_2 \rangle$ as v_1 is the value of the highest numbered proposal that it receives as a promise response. Thus, meeting our above requirement.

4.2 Verification

5

Client Application

6

Conclusion and Evaluation

6.1 Summary of Achievements

6.2 Critical Evaluation

6.3 Future Work

6.4 Final Thoughts

References



Interim Report