

Paxos Made Moderately Complex

Robbert van Renesse
Cornell University
rvr@cs.cornell.edu

March 25, 2011

Abstract

For anybody who has ever tried to implement it, Paxos is by no means a simple protocol, even though it is based on relatively simple invariants. This paper provides imperative pseudo-code for the full Paxos (or Multi-Paxos) protocol without shying away from discussing various implementation details. The initial description avoids optimizations that complicate comprehension. Next we discuss liveness, and list various optimizations that make the protocol practical.

1 Introduction

Paxos [13] is a protocol for state machine replication in an asynchronous environment that admits crash failures. It is useful to consider the terms in this sentence carefully:

- A *state machine* consists of a collection of states, a collection of transitions between states, and a current state. A transition to a new current state happens in response to an issued operation and produces an output. Transitions from the current state to the same state are allowed, and are used to model read-only operations. In a *deterministic state machine*, for any state and operation, the transition enabled by the operation is unique.
- In an *asynchronous environment*, there are no bounds on timing characteristics. Clocks run arbitrarily fast, network communication takes arbitrarily long, and state machines take arbitrar-

ily long to transition in response to an operation. The term “asynchronous” as used here should not be confused with nonblocking operations on objects that are often called asynchronous as well.

- A state machine has experienced a *crash failure* if it will make no more transitions and thus its current state is fixed indefinitely. No other failures of a state machine, such as experiencing undocumented transitions, are allowed. In a “fail-stop environment” [21], crash failures can be reliably detected—not so in an asynchronous environment.
- *State Machine Replication* (SMR) [12, 22] is a technique to mask failures, and crash failures in particular. A collection of replicas of a deterministic state machine are created. The replicas are then provided with the same sequence of operations, so they end up in the same state and produce the same sequence of outputs. It is assumed that at least one replica never crashes.

Deterministic state machines are used to model server processes, such as a file server, a DNS server, and so on. A client process, using a library “stub routine,” can send a *command* to such a server over a network and await an output. A command is a triple $\langle \kappa, cid, operation \rangle$, where κ^1 is the identifier of the client that issued the command and *cid* a client-local unique command identifier (such as a sequence number). The command identifier must be included in the response from the server so the client can match

¹As in [13], we will use Greek letters to identify processes.

responses with commands. In SMR replication, that stub routine is replaced with another to provide the illusion of a single remote server that is highly available. The stub routine sends the command to all replicas and returns only the first response to the command.

The difficulty comes with multiple clients, as concurrent commands may arrive in different orders at the replicas, and thus the replicas may end up taking different transitions, producing different outputs as a result and possibly ending up in different current states. A protocol like Paxos ensures that this cannot happen: the replicated state machine behaves logically identical to a single remote state machine that never crashes [9].

While processes may crash, we assume that messaging between processes is reliable (but not necessarily FIFO):

- a message sent by a non-faulty process to a non-faulty destination process is eventually received (at least once) by the destination process;
- if a message is received by a process, it was sent by some (possibly faulty) process. That is, messages are not garbled and do not appear out of the blue.

This paper gives an *operational description* of the *multi-decree Paxos protocol*, sometimes called *multi-Paxos*. Single-decree Paxos is significantly easier to understand, and is the topic of such papers as [15, 14]. But the multi-decree Paxos protocol is the one that is used (or some variant thereof) within industrial-strength systems like Chubby [4] and ZooKeeper [10].

2 How and Why Paxos Works

Replicas and Slots

In order to tolerate f crashes, Paxos needs at least $f + 1$ replicas. When a client κ wants to execute a command $\langle \kappa, cid, op \rangle$, it broadcasts a $\langle \mathbf{request}, \langle \kappa, cid, op \rangle \rangle$ message to all replicas and waits for a $\langle \mathbf{response}, cid, result \rangle$ message from one of the replicas.

The replicas can be thought of as having a sequence of *slots* that need to be filled with commands. Each

```

process Replica(leaders, initial_state)
  var state := initial_state, slot_num := 1;
  var proposals :=  $\emptyset$ , decisions :=  $\emptyset$ ;

  function propose(p)
    if  $\nexists s : \langle s, p \rangle \in decisions$  then
       $s' := \min\{s \mid s \in \mathbb{N}^+ \wedge$ 
         $\nexists p' : \langle s, p' \rangle \in proposals \cup decisions\}$ ;
      proposals := proposals  $\cup \{\langle s', p \rangle\}$ ;
       $\forall \lambda \in leaders : send(\lambda, \langle \mathbf{propose}, s', p \rangle)$ ;
    end if
  end function

  function perform( $\langle \kappa, cid, op \rangle$ )
    if  $\exists s : s < slot\_num \wedge$ 
       $\langle s, \langle \kappa, cid, op \rangle \rangle \in decisions$  then
      slot_num := slot_num + 1;
    else
       $\langle next, result \rangle := op(state)$ ;
      atomic
        state := next;
        slot_num := slot_num + 1;
      end atomic
      send( $\kappa, \langle \mathbf{response}, cid, result \rangle$ );
    end if
  end function

  for ever
    switch receive()
      case  $\langle \mathbf{request}, p \rangle$  :
        propose(p);
      case  $\langle \mathbf{decision}, s, p \rangle$  :
        decisions := decisions  $\cup \{\langle s, p \rangle\}$ ;
        while  $\exists p' : \langle slot\_num, p' \rangle \in decisions$  do
          if  $\exists p'' : \langle slot\_num, p'' \rangle \in proposals \wedge$ 
             $p'' \neq p'$  then
            propose(p'');
          end if
          perform(p');
        end while;
    end switch
  end for
end process

```

Figure 1: Pseudo code for a replica.

slot is indexed by a *slot number*. A replica, on receipt of a $\langle \mathbf{request}, p \rangle$ message, proposes command p for the lowest unused slot. In the face of concurrently operating clients, different replicas may end up proposing different commands for the same slot. In order to avoid inconsistency, a replica awaits a decision for a slot before actually updating its state and computing a response to send back to the client.

Replicas are not necessarily identical at any time. They may propose different commands for different slots. However, replicas apply operations to the application state in the same order. Figure 1 shows pseudo-code for a replica. Each replica ρ maintains four variables:

- $\rho.state$, the (opaque) application state (all replicas are started with the same initial application state);
- $\rho.slot_num$, the replica’s current slot number (equivalent to the version of the state, and initially 1). It contains the index of the next slot for which it needs to learn a decision before it can update the application state;
- $\rho.proposals$, a set of $\langle \text{slot number, command} \rangle$ pairs for proposals that the replica has made in the past (initially empty); and
- $\rho.decisions$, another set of $\langle \text{slot number, command} \rangle$ pairs for decided slots (also initially empty).

Before giving an operational description of replicas, we present some important invariants that hold over the collected variables of replicas:

R1: There are no two different commands decided for the same slot: $\forall s, \rho_1, \rho_2, p_1, p_2 : \langle s, p_1 \rangle \in \rho_1.decisions \wedge \langle s, p_2 \rangle \in \rho_2.decisions \Rightarrow p_1 = p_2$

R2: All commands up to $slot_num$ are in the set of decisions: $\forall \rho, s : 1 \leq s < \rho.slot_num \Rightarrow (\exists p : \langle s, p \rangle \in \rho.decisions)$

R3: For all replicas ρ , $\rho.state$ is the result of applying the operations in $\langle s, p_s \rangle \in \rho.decisions$ for all s such that $1 \leq s < slot_num$, in order of slot number, to *initial_state*.

R4: For each ρ , the variable $\rho.slot_num$ cannot decrease over time.

From Invariants R1-3, it is clear that all replicas apply operations to the application state in the same order, and thus replicas with the same slot number have the same *state*. Invariant R4 ensures that a replica cannot go back in time.

Returning to Figure 1, a replica runs in an infinite loop, receiving requests in messages. Replicas receive two kinds of messages: requests from clients, and decisions. When it receives a request for command p from a client, the replica invokes $propose(p)$. This function checks if there has been a decision for p already. If so, the replica has already sent a response and the request can be ignored. If not, the replica determines the lowest unused slot number s' , and adds $\langle s', p \rangle$ to its set of proposals. It then sends a $\langle \mathbf{propose}, s', p \rangle$ message to all *leaders*. Leaders are described below.

Decisions may arrive out-of-order and multiple times. For each **decision** message, the replica adds the decision to the set of decisions. Then, in a loop, it considers which decisions are ready for execution before trying to receive more messages. If there is a decision p' corresponding to the current $slot_num$, the replica first checks to see if it has proposed a different command p'' . If so, it re-proposes p'' , which will be assigned a new slot number. Next, it invokes $perform(p')$.

The function $perform()$ is invoked with the same sequence of commands at all replicas. First, it checks to see if has already performed the command. Different replicas may end up proposing the same command for different slots, and thus the same command may be decided multiple times. The corresponding operation is evaluated only if the command is new.

Note that both *proposals* and *decisions* are “append-only” in that there is no code that removes entries from these sets. Doing so makes it easier to formulate invariants and reason about the correctness of the code. In Section 4.2 we will discuss correctness-preserving ways of removing entries that are no longer used.

It is clear that the code enforces Invariant R4. The variables *state* and *slot_num* are updated atomically in order to ensure that Invariant R3 holds, although

in practice it is not actually necessary to perform these updates atomically, as the intermediate state is not externally visible. Since *slot_num* is only advanced if the corresponding decision is in *decisions*, it is clear that Invariant R2 holds.

The real difficulty lies in enforcing Invariant R1. It requires that the set of replicas agree on the order of commands. For each slot, the Paxos protocol *chooses* a command from among a collection of commands proposed by clients. This is called *consensus*, and in Paxos the subprotocol that implements consensus is called the multi-decree *Synod* protocol.

The Synod Protocol, Ballots, and Acceptors

In the Synod protocol, there is an infinite collection of *ballots*. Ballots are not created; they just are. As we shall see later, ballots are the key to liveness in Paxos. Each ballot has a unique *leader*,² a deterministic state machine in its own right. A leader can be working on arbitrarily many ballots, although it will be predominantly working on one at a time, even as multiple slots are being decided. A leader process has a unique identifier called the *leader identifier*. A ballot has a unique identifier as well, called its *ballot number*. Ballot numbers are totally ordered, that is, for any two different ballot numbers, one is before or after the other.

In this description, we will have ballot numbers be lexicographically ordered pairs of an integer and its leader identifier (consequently, leader identifiers need to be totally ordered as well). This way, given a ballot number, it is trivial to see who the leader of the ballot is. We will use one special ballot number \perp that is ordered before any normal ballot number.

Besides replicas and leaders, there is a fixed collection of *acceptors*, deterministic state machines themselves (although *not* replicas of one another, because they get different sequences of input). Acceptors are servers, and leaders are their clients. As we shall see, acceptors are the memory of Paxos, preventing conflicting decision from being made. We will assume that at most a proper minority of acceptors can crash. Thus, in order to tolerate f crash failures, Paxos needs at least $2f + 1$ acceptors.

²We stress here that the leader of a ballot is fixed: it is not elected, and it is allowed to crash.

An acceptor is quite simple, as it is passive and only sends messages in response to requests. Its state consists of two variables. Let a *pvalue* be a triple consisting of a ballot number, a slot number, and a proposal (which is a command). If α is the identifier of an acceptor, then the acceptor's state is described by

- $\alpha.ballot_num$: a ballot number, initially \perp ;
- $\alpha.accepted$: a set of pvalues, initially empty.

Under the direction of request messages sent by leaders, the state of an acceptor can change. Let $e = \langle b, s, p \rangle$ be a pvalue consisting of a ballot number b , a slot number s , and a proposal p . When an acceptor α adds e to $\alpha.accepted$, we say that α *accepts* e . (An acceptor may accept the same pvalue multiple times.) When α sets its ballot number to b for the first time, we say that α *adopts* b .

We start by presenting some important invariants that hold over the collected variables of acceptors. Knowing these invariants are an invaluable help to understanding the Synod protocol:

- A1: an acceptor can only adopt strictly increasing ballot numbers;
- A2: an acceptor α can only add $\langle b, s, p \rangle$ to $\alpha.accepted$ (*i.e.*, accept $\langle b, s, p \rangle$) if $b = ballot_num$;
- A3: acceptor α cannot remove pvalues from $\alpha.accepted$ (we will modify this impractical restriction later);
- A4: Suppose α and α' are acceptors, with $\langle b, s, p \rangle \in \alpha.accepted$ and $\langle b, s, p' \rangle \in \alpha'.accepted$. Then $p = p'$. Informally, given a particular ballot number and slot number, there can be at most one proposal under consideration by the set of acceptors.
- A5: Suppose that for each α among a majority of acceptors, $\langle b, s, p \rangle \in \alpha.accepted$. If $b' > b$ and $\langle b', s, p' \rangle \in \alpha'.accepted$, then $p = p'$. We will consider this crucial invariant in more detail later.

Figure 2 shows pseudo-code for an acceptor. It runs in an infinite loop, receiving two kinds of request messages from leaders (note the use of pattern matching):

```

process Acceptor()
  var ballot_num :=  $\perp$ , accepted :=  $\emptyset$ ;

  for ever
    switch receive()
      case  $\langle \mathbf{p1a}, \lambda, b \rangle$  :
        if  $b > \text{ballot\_num}$  then
          ballot_num :=  $b$ ;
        end if;
        send( $\lambda, \langle \mathbf{p1b}, \text{self}(), \text{ballot\_num}, \text{accepted} \rangle$ );
      end case
      case  $\langle \mathbf{p2a}, \lambda, \langle b, s, p \rangle \rangle$  :
        if  $b \geq \text{ballot\_num}$  then
          ballot_num :=  $b$ ;
          accepted := accepted  $\cup$   $\{ \langle b, s, p \rangle \}$ ;
        end if
        send( $\lambda, \langle \mathbf{p2b}, \text{self}(), \text{ballot\_num} \rangle$ );
      end case
    end switch
  end for
end process

```

Figure 2: Pseudo code for an acceptor.

- $\langle \mathbf{p1a}, \lambda, b \rangle$: Upon receiving a “phase 1a” request message from a leader with identifier λ , for a ballot number b , an acceptor makes the following transition. First, the acceptor adopts b if and only if it exceeds its current ballot number. Then it returns to λ a “phase 1b” response message containing all pvalues accepted thus far by the acceptor.
- $\langle \mathbf{p2a}, \lambda, \langle b, s, p \rangle \rangle$: Upon receiving a “phase 2a” request message from leader λ with pvalue $\langle b, s, p \rangle$, an acceptor makes the following transition. If b exceeds the current ballot number, then the acceptor first adopts b . If the current ballot number equals b , then the acceptor accepts $\langle b, s, p \rangle$. The acceptor returns to λ a “phase 2b” response message containing its current ballot number.

It is easy to see that the code enforces Invariants A1, A2, and A3. For checking the remaining two invariants, which involve multiple acceptors, we have to study what a leader does first.

Leaders and Commanders

Leaders are responsible for selecting proposals within a ballot, and have to make sure that they do not select proposals that could conflict with decisions on other ballots. A leader may work on multiple slots at the same time. When, in ballot b , its leader tries to get a proposal p for slot number s chosen, it spawns a local *commander* thread for $\langle b, s, p \rangle$. While we present it here as a separate process, the commander is really just a thread running within the leader. As we shall see, the following invariants hold in the Synod protocol:

- C1: For any b and s , at most one commander is spawned;
- C2: Suppose that for each α among a majority of acceptors $\langle b, s, p \rangle \in \alpha.\text{accepted}$. If $b' > b$ and a commander is spawned for $\langle b', s, p' \rangle$, then $p = p'$.

Invariant C1 implies Invariant A4, because by C1 all acceptors that accept a pvalue for a particular ballot and slot number received the pvalue from the same commander. Similarly, Invariant C2 implies Invariant A5.

Figure 3(a) shows the pseudo-code for a commander. A commander sends a $\langle \mathbf{p2a}, \lambda, \langle b, s, p \rangle \rangle$ message to all acceptors, and waits for responses of the form $\langle \mathbf{p2b}, \alpha, b' \rangle$. In each such response $b' \geq b$ will hold (see the code for acceptors). There are two cases:

1. If a commander receives $\langle \mathbf{p2b}, \alpha, b \rangle$ from all acceptors in a majority of acceptors, then the commander learns that proposal p has been chosen for slot s . In this case the commander notifies the replicas and exits. In order to satisfy Invariant R1, we need to enforce that if a commander learns that p is chosen for slot s , and another commander learns that p' is chosen for the same slot s , then $p = p'$. This is a consequence of Invariant A5: if a majority of acceptors accept $\langle b, s, p \rangle$, then for any later ballot b' and the same slot number s , acceptors can only accept $\langle b', s, p \rangle$. Thus if the commander of $\langle b', s, p' \rangle$ learns that p' has been chosen for s , it is guaranteed that $p = p'$ and no inconsistency occurs, assuming—of course—that Invariant C2 holds.

```

process Commander( $\lambda$ , acceptors, replicas,  $\langle b, s, p \rangle$ )
  var waitfor := acceptors;

   $\forall \alpha \in \textit{acceptors}$  : send( $\alpha$ ,  $\langle \mathbf{p2a}, \textit{self}(), \langle b, s, p \rangle \rangle$ );
  for ever
    switch receive()
      case  $\langle \mathbf{p2b}, \alpha, b' \rangle$  :
        if  $b' = b$  then
          waitfor := waitfor -  $\{\alpha\}$ ;
          if  $|\textit{waitfor}| < |\textit{acceptors}|/2$  then
             $\forall \rho \in \textit{replicas}$  :
              send( $\rho$ ,  $\langle \mathbf{decision}, s, p \rangle$ );
            exit();
          end if;
        else
          send( $\lambda$ ,  $\langle \mathbf{preempted}, b' \rangle$ );
          exit();
        end if;
      end case
    end switch
  end for
end process

```

(a)

```

process Scout( $\lambda$ , acceptors,  $b$ )
  var waitfor := acceptors, pvalues :=  $\emptyset$ ;

   $\forall \alpha \in \textit{acceptors}$  : send( $\alpha$ ,  $\langle \mathbf{p1a}, \textit{self}(), b \rangle$ );
  for ever
    switch receive()
      case  $\langle \mathbf{p1b}, \alpha, b', r \rangle$  :
        if  $b' = b$  then
          pvalues := pvalues  $\cup$   $r$ ;
          waitfor := waitfor -  $\{\alpha\}$ ;
          if  $|\textit{waitfor}| < |\textit{acceptors}|/2$  then
            send( $\lambda$ ,  $\langle \mathbf{adopted}, b, \textit{pvalues} \rangle$ );
            exit();
          end if;
        else
          send( $\lambda$ ,  $\langle \mathbf{preempted}, b' \rangle$ );
          exit();
        end if;
      end case
    end switch
  end for
end process

```

(b)

Figure 3: (a) Pseudo code for a commander. Here λ is the identifier of its leader, *acceptors* the set of acceptor identifiers, *replicas* the set of replicas, and $\langle b, s, p \rangle$ the pvalue the commander is responsible for. (b) Pseudo code for a scout. Here λ is the identifier of its leader, *acceptors* the identifiers of the acceptors, and b the desired ballot number.

2. If a commander receives $\langle \mathbf{p2b}, \alpha', b' \rangle$ from some acceptor α' , with $b' \neq b$, then it learns that a ballot b' (which must be larger than b) is active. This means that ballot b may no longer be able to make progress, as there may no longer exist a majority of acceptors that can accept $\langle b, s, p \rangle$. In this case, the commander notifies its leader about the existence of b' , and exits.

Under the assumption that at most a minority of acceptors can fail and messages are delivered reliably, the commander will eventually do one or the other.

Scouts, Passive and Active Modes

The leader must enforce Invariants C1 and C2. Invariant C1 is trivial to enforce by not spawning more than one commander per ballot number and slot number. In order to enforce Invariant C2, the leader of a ballot runs what is often called a *view change* protocol before spawning commanders for that ballot.³ The leader spawns a *scout* thread to run the view change protocol for some ballot b . A leader starts at most one of these for any ballot b , and only for its own ballots.

Figure 3(b) shows the pseudo-code for a scout. The code is similar to that of a commander, except that

³The term “view change” is used in the Viewstamped Replication protocol [19], largely identical to the Paxos protocol.

it sends and receives phase 1 instead of phase 2 messages. A scout completes successfully when it has collected $\langle \mathbf{p1b}, \alpha, b, r_\alpha \rangle$ messages from all acceptors in a majority (again, guaranteed to complete eventually), and returns an $\langle \mathbf{adopted}, b, \bigcup r_\alpha \rangle$ message to its leader λ . As we will see later, the leader needs $\bigcup r_\alpha$, the union of all pvalues accepted by this majority of acceptors, in order to enforce Invariant C2.

Figure 4 shows the main code of a leader. Leader λ maintains three state variables:

- $\lambda.ballot_num$: a monotonically increasing ballot number, initially $(0, \lambda)$;
- $\lambda.active$: a boolean flag, initially **false**; and
- $\lambda.proposals$: a map of slot numbers to proposals in the form of a set of $\langle slot\ number, proposal \rangle$ pairs, initially empty. At any time, there is at most one entry per slot number in the set.

The leader starts by spawning a scout for its initial ballot number, and then enters into a loop awaiting messages. There are three types of messages that cause transitions:

- $\langle \mathbf{propose}, s, p \rangle$: A replica proposes command p for slot number s ;
- $\langle \mathbf{adopted}, ballot_num, pvals \rangle$: Sent by a scout, this message signifies that the current ballot number $ballot_num$ has been adopted by a majority of acceptors. (If an **adopted** message arrives for an old ballot number, it is ignored.) The set $pvals$ contains all pvalues accepted by these acceptors prior to $ballot_num$.
- $\langle \mathbf{preempted}, \langle r', \lambda' \rangle \rangle$: Sent by either a scout or a commander, it means that some acceptor has adopted $\langle r', \lambda' \rangle$. If $\langle r', \lambda' \rangle > ballot_num$, it may no longer be possible to use $ballot_num$.

A leader goes between *passive* and *active* modes. When passive, the leader is waiting for an $\langle \mathbf{adopted}, ballot_num, pvals \rangle$ message from the last scout that it spawned. When this message arrives, the leader becomes active and spawns commanders for each of the slots for which it has a proposal, but must select proposals that satisfy Invariant C2. We will now consider how the leader goes about this.

The leader knows that a majority of acceptors, say \mathcal{A} , have adopted $ballot_num$ and thus no longer accept pvalues for ballot numbers less than $ballot_num$ (because of Invariants A1 and A2). There are two cases to consider:

1. If, for some slot s , there is no pvalue in $pvals$, then, prior to $ballot_num$, it is not possible that any pvalue has been chosen or will be chosen for slot s . After all, suppose that some pvalue $\langle b, s, p \rangle$ were chosen, with $b < ballot_num$. This would require a majority of acceptors \mathcal{A}' to accept $\langle b, s, p \rangle$, but we have responses from a majority \mathcal{A} that have adopted $ballot_num$ and have not accepted, nor can accept, pvalues with a ballot number smaller than $ballot_num$ (Invariants A1 and A2). Clearly, there is a contradiction, because $\mathcal{A} \cap \mathcal{A}'$ is non-empty. Thus any proposal for slot s will satisfy Invariant C2.
2. Otherwise, let $\langle b, s, p \rangle$ be the pvalue with the maximum ballot number for slot s . Because of Invariant A4, this pvalue is unique—there cannot be two different proposals for the same ballot number and slot number. Also note that $b < ballot_num$ (because acceptors only report pvalues they accepted before $ballot_num$). Like the leader of $ballot_num$, the leader of b must have picked p carefully to ensure that Invariant C2 holds, and thus if a pvalue is chosen before or at b , its proposal must be p . Since all acceptors in \mathcal{A} have adopted $ballot_num$, no pvalues between b and $ballot_num$ can be chosen (Invariants A1 and A2). Thus, by using p as a proposal, λ enforces Invariant C2.

This inductive argument is the crux for the correctness of the Synod protocol. It demonstrates that Invariant C2 holds, which in turn implies Invariant A5, which in turn implies Invariant R1 that ensures that all replicas apply the same operations in the same order.

Back to the code, after the leader receives $\langle \mathbf{adopted}, ballot_num, pvals \rangle$, it determines for each slot the proposal corresponding to the maximum ballot number in $pvals$ by invoking the function $pmax$. Formally, the function $pmax(pvals)$ is defined as follows:

```

process Leader(acceptors, replicas)
  var ballot_num = (0, self()), active = false, proposals =  $\emptyset$ ;

  spawn(Scout(self()), acceptors, ballot_num);
  for ever
    switch receive()
      case  $\langle$ propose, s, p $\rangle$  :
        if  $\nexists p' : \langle s, p' \rangle \in$  proposals then
          proposals := proposals  $\cup$   $\{\langle s, p \rangle\}$ ;
          if active then
            spawn(Commander(self()), acceptors, replicas,  $\langle$ ballot_num, s, p $\rangle$ );
          end if
        end if
      end case
      case  $\langle$ adopted, ballot_num, pvals $\rangle$  :
        proposals := proposals  $\oplus$  pmax(pvals);
         $\forall \langle s, p \rangle \in$  proposals : spawn(Commander(self()), acceptors, replicas,  $\langle$ ballot_num, s, p $\rangle$ );
        active := true;
      end case
      case  $\langle$ preempted,  $\langle r', \lambda' \rangle$  $\rangle$  :
        if  $\langle r', \lambda' \rangle >$  ballot_num then
          active := false;
          ballot_num :=  $\langle r' + 1, \text{self}() \rangle$ ;
          spawn(Scout(self()), acceptors, ballot_num);
        end if
      end case
    end switch
  end for
end process

```

Figure 4: Pseudo code skeleton for a leader. Here *acceptors* is the set of acceptor identifiers, and *replicas* the set of replica identifiers.

$$\begin{aligned}
\textit{pmax}(\textit{pvals}) \equiv & \{ \langle s, p \rangle \mid \exists b : \langle b, s, p \rangle \in \textit{pvals} \wedge \\
& \forall b', p' : \langle b', s, p' \rangle \in \textit{pvals} \Rightarrow b' \leq b \}
\end{aligned}$$

The update operator \oplus applies to two maps of slot numbers to proposals (sets of $\langle \textit{slot number}, \textit{proposal} \rangle$ pairs). $x \oplus y$ returns the elements of y as well as the elements of x that are not in y . Formally:

$$\begin{aligned}
x \oplus y \equiv & \{ \langle s, p \rangle \mid \langle s, p \rangle \in y \vee \\
& (\langle s, p \rangle \in x \wedge \nexists p' : \langle s, p' \rangle \in y) \}
\end{aligned}$$

Thus the line $\textit{proposals} := \textit{proposals} \oplus \textit{pmax}(\textit{pvals})$; updates the map of slot numbers to proposals, replacing for each slot number the proposal corresponding to the maximum pvalues in *pvals*. Now the leader can start commanders for each proposal while satisfying Invariant C2.

If a new proposal arrives while the leader is active, the leader checks to see if it already has a proposal for the same slot (and has thus spawned a commander for that slot) in its set *proposers*. If not, the new proposal will satisfy Invariant C2, and thus the leader adds the proposal to *proposers* and spawns a commander.

If either a scout or a commander notifies that an acceptor has adopted a ballot number b , with $b > \text{ballot_num}$, then it sends the leader a **preempted** message. The leader becomes passive and spawns a new scout with a ballot number that is higher than b .

Figure 5 shows an example of a leader spawning a scout to become active, and a client sending a request to a replica, which in turns sends a proposal to an active leader.

3 When Paxos Works

It would clearly be desirable that, if a client broadcasts a new command to all replicas, that it eventually receives at least one response. This is often referred to as *liveness*. It requires that if one or more commands have been proposed for a particular slot, that some command is eventually decided for that slot. Unfortunately, the Synod protocol as described does not guarantee this, even in the absence of any failure whatsoever.

Consider the following scenario, with two leaders with identifiers λ and λ' such that $\lambda < \lambda'$. Both start at the same time, respectively proposing commands p and p' for slot number 1. Suppose there are three acceptors, α_1 , α_2 , and α_3 . In ballot $\langle 0, \lambda \rangle$, leader λ is successful getting α_1 and α_2 to adopt the ballot, and α_1 to accept pvalue $\langle \langle 0, \lambda \rangle, 1, p \rangle$.

Now leader λ' gets α_2 and α_3 to adopt $\langle 0, \lambda' \rangle$ (which is after $\langle 0, \lambda \rangle$ because $\lambda < \lambda'$). Note that neither α_2 or α_3 accepted any pvalues. Leader λ' then gets α_3 to accept $\langle \langle 0, \lambda' \rangle, 1, p' \rangle$.

Going on, leader λ gets α_1 and α_2 to adopt $\langle 1, \lambda \rangle$. The maximum pvalue accepted by α_1 and α_2 is $\langle \langle 0, \lambda \rangle, 1, p \rangle$, and thus λ must propose p . Suppose λ gets α_1 to accept $\langle \langle 1, \lambda \rangle, 1, p \rangle$. Subsequently, leader λ' gets α_2 and α_3 to adopt $\langle 1, \lambda' \rangle$, and gets α_3 to accept $\langle \langle 1, \lambda' \rangle, 1, p' \rangle$.

As is now clear, this can be indefinitely continued, with no ballot ever succeeding in choosing a pvalue. This is true even if $p = p'$, that is, the leaders propose the same command. The well-known “FLP impossibility result” [7] demonstrates that in an asynchronous environment that admits crash failures, no consensus protocol can guarantee termination, and the Synod protocol is no exception. The argu-

ment does not apply directly if transitions have non-deterministic actions—for example changing state in a randomized manner. However, it can be demonstrated that such protocols cannot guarantee a decision either.

If we could somehow guarantee that some leader would be able to work long enough to get a majority of acceptors to adopt a high ballot and also accept a pvalue, then Paxos would be guaranteed to choose a proposed command. A possible approach could be as follows: when a leader λ discovers (through a **preempted** message) that there is a higher ballot with leader λ' active, rather than starting a new scout with an even higher ballot number, it starts monitoring the leader of b by pinging it on a regular basis. As long as λ' responds timely to pings, leader λ waits patiently. Only if λ' stops responding will λ select a higher ballot number and start a scout.

This concept is called *failure detection*, and theoreticians have been interested in the weakest properties failure detection should have in order to support a consensus algorithm that is guaranteed to terminate [6]. In a purely asynchronous environment it is impossible to determine through pinging or any other method whether a particular leader has crashed or is simply slow. However, under fairly weak assumptions about timing, we can design a version of Paxos that is guaranteed to choose a proposal. In particular, we will assume that both the following are bounded:

- the clock drift of a process, that is, the rate of its clock is within some factor of the rate of real-time;
- the time between when a non-faulty process initiates sending a message, and the message having been received and handled by a non-faulty destination process.

We do not need to assume that we know what those bounds are—only that such bounds exist. From a practical point of view, this seems entirely reasonable. Modern clocks are certainly within a factor of 2 of real-time. A message between two non-faulty processes is likely delivered within a year, say. Even if the network was partitioned at the time the sender started sending the message, by the time a year has

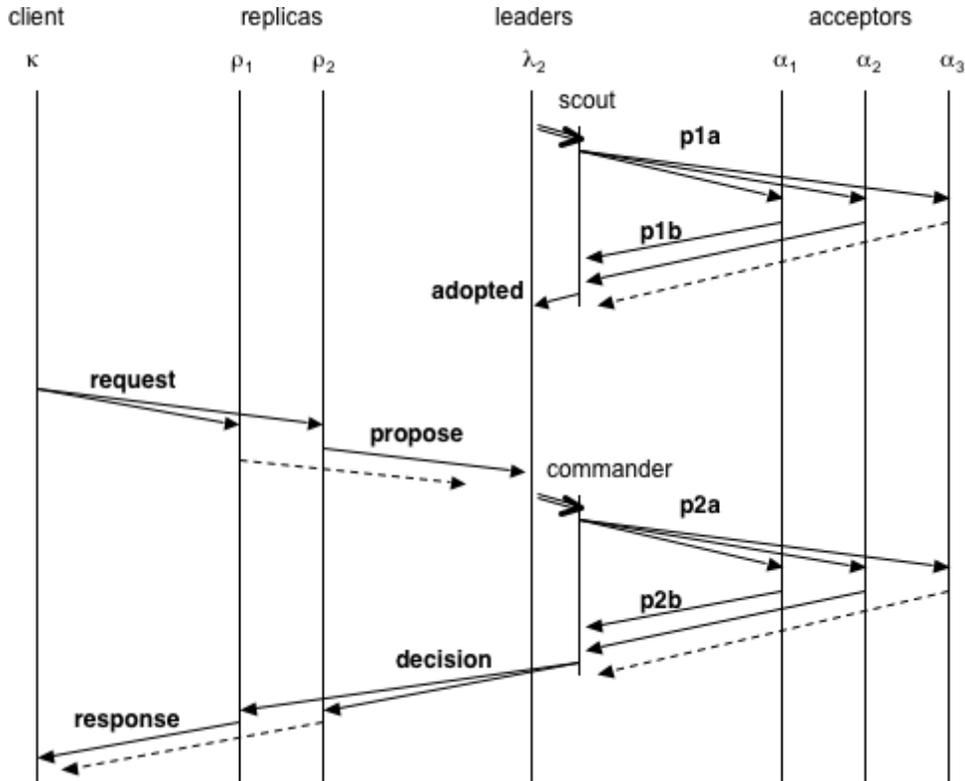


Figure 5: The time diagram shows a client, two replicas, a leader (with a scout and a commander), and three acceptors, with time progressing downward. Arrows represent messages. Dashed arrows are messages that end up being ignored. The leader first runs a scout to become active. Later, when a replica proposes a command (in response to a client’s request), the leader runs a commander, which notifies the replicas upon learning a decision.

expired the message is highly likely to have been delivered and processed.

These assumptions can be exploited as follows: we use a scheme similar to the one described above, based on pinging and timeouts, but the value of the timeout interval depends on the ballot number: the higher the competing ballot number, the longer a leader waits before trying to preempt it with a higher ballot number. Eventually the timeout at each of the leaders becomes so high that some correct leader will always be able to get its proposals chosen.

For good performance, one would like the timeout period to be long enough so that a leader can be successful, but short enough so that a faulty leader is pre-

empted as quickly as possible. This can be achieved with a TCP-like AIMD (Additive Increase, Multiplicative Decrease) approach for choosing timeouts. The leader associates an initial timeout with each ballot. If a ballot gets preempted, the next ballot uses a timeout that is multiplied by some factor larger than one. With each chosen proposal, this initial timeout is decreased linearly. Eventually the timeout will become too short, and the ballot replaced with another even if its leader is non-faulty, but this does not affect correctness.

(As an aside: some people call this process *leader election* or *weak leader election*. This is, however,

confusing, as each ballot has a fixed leader that is not elected.)

For further improved liveness, crashes should be avoided. The Paxos protocol can tolerate a minority of its acceptors failing, and all but one of its replicas failing. If more than that fail, consistency is still guaranteed, but liveness will be violated. For this reason, one may want to keep the state of acceptors and replicas on disk. A process that suffers from a power failure but can recover from disk is not theoretically considered crashed—it is simply slow for a while. However, a process that suffers a permanent disk failure would be considered crashed.

4 Paxos Made Pragmatic

We have described a simple version of the Paxos protocol with the intention to make it understandable, but the described protocol is not practical. The state of the various components, as well as the contents of `p1b` messages, grows much too quickly. Also, we have not made it clear where the various components should run. This section is a list of various optimizations and design decisions.

4.1 State Reduction

First note that although a leader gets a set of all accepted pvalues from a majority of acceptors, it only needs to know if this set is empty or not, and if not, what the maximum pvalue is. Thus, a large step toward practicality is that acceptors only maintain the most recently accepted pvalue for each slot (\perp if no pvalue has been accepted) and return only these pvalues in a `p1b` message to the scout. This gives the leader all information needed to enforce Invariant C2.

The `p1b` message can be further reduced in size if a leader keeps track of which slots have been completed. A leader includes on the `p1a` request the first slot for which it does not know the decision. Acceptors do not need to respond with pvalues for smaller slot numbers.

Also note that the set *requests* maintained by a replica only needs to contain those requests for slot numbers higher than *slot_num*.

Finally, note that the leader maintains proposals for all slots and starts new commanders upon turning active for each slot for which it has a proposal, even if those slots already have been decided. Clearly a lot of storage and work could be avoided if the leader kept track of which slots have been decided already. Leaders can learn this from their co-located commanders, or alternatively from replicas in case leaders and replicas are co-located (see Section 4.3).

4.2 Garbage Collection

When *all* replicas have learned that some slot has been decided, then there is no longer a good reason for an acceptor to maintain the corresponding pvalues in its *accepted* set. To enable this garbage collection, replicas could respond to leaders when they have performed a command, and upon a leader learning that all replicas have done so it could notify the acceptors to release the corresponding state.

The state of an acceptor would have to be extended with a new variable that contains a slot number: all pvalues lower than that slot number have been garbage collected. This slot number must be included in `p1b` messages so that a leader does not mistakenly conclude that the acceptors have not accepted any pvalues for those commands.

This garbage collection technique does not work if one of the replicas is faulty or slow, leaving the acceptors no option but to maintain state for all slots. A solution is to use $2f + 1$ or more replicas instead of just $f + 1$. Acceptor state is garbage collected when more than f replicas have performed a command. If because of this a replica is not able to learn a particular command, then it can always obtain a snapshot of the state from another replica that has learned the operation and continue on.

Another, but more complicated, solution is to make the set of replicas adaptive, by having the replicas themselves keep track of which $f + 1$ replicas are currently active. A special command is used to change the set of replicas in case one or more are suspected of having failed.

4.3 Co-location

So far we have treated leaders (along with their respective scouts and commanders) as if they were running on separate machines. In practice, however, leaders are typically co-located with replicas. That is, each machine that runs a replica also runs a leader.

A client sends its proposals to replicas. If co-located, the replica can send a proposal for a particular slot to its local leader, say λ , rather than broadcasting the request to all leaders. If λ is passive, monitoring another leader λ' , it may forward the request to λ' . If λ is active, it will start a commander.

An alternative not often considered is to have clients and leaders be co-located instead of replicas and leaders. Thus, each client runs a local leader. By doing so, one obtains a protocol that is much like Quorum Replication [23, 2]. While traditional quorum replication protocols can only support read and write operations, this Paxos version could support arbitrary (deterministic) operations. In practice, however, services probably want to be in control over their own liveness, and choose to co-locate leaders with replicas.

Replicas are also often co-located with acceptors. As shown in Section 4.2, one may need as many replicas as acceptors in any case. When leaders are co-located with acceptors, one has to be careful that they use separate ballot number variables.

4.4 Read-only Commands

The Paxos protocol does not treat read-only commands any different from other commands, and this leads to more overhead than necessary. One would however be naive in thinking that a client that wants to do a read-only command could simply query one of the replicas—doing so would easily violate consistency as the selected replica may not be up-to-date.

Therefore, read-only commands are typically sent to the leader just like update commands. One simple optimization is for a leader to send a chosen read-only command to only a single replica instead of to all replicas.⁴ After all, the state of none of the replicas needs to change, but one of the replicas has to

⁴For this to work, a leader needs to be able to recognize read-only commands.

compute the result and send it to the client. (One has to consider the case that the selected replica is faulty and does not send a result to the client. Again, the end-to-end argument applies, and the client may have to query the service to see if its command has been performed and, if so, what the result was.)

A read-only command does not actually require that acceptors accept a pvalue. A leader does have to run a scout and wait for an `adopted` message to ensure that its ballot is current. The leader then “attaches” the command to the highest slot number for which it knows the decision or for which it is proposing a command that contains an update command. Once decided, the leader can send all read-only commands attached to the slot number to one of the replicas, which can perform the read-only commands after it has performed the update command for the corresponding slot number.

However, while this avoids unnecessary accepts, running a view change for each read-only command is an expensive proposition. The better solution involves so-called *leases* [8, 13]. Leases require an additional assumption on timing, which is that there is a *known* bound on clock drift. For simplicity, we will assume that there is no clock drift whatsoever.

Before a leader sends a `p1a` request to the acceptors, it records the time. The leader includes in the `p1a` request a *lease period*. For example, the lease period could be “10 seconds.” An acceptor that adopts the ballot number promises not to adopt another (higher) ballot number until the lease period expires (measured on its local clock from the time the acceptor received the `p1a` request). If a majority of acceptors accept the ballot, the leader can be certain that from the recorded time, until the lease period expires on its own clock, no other leader can preempt its ballot, and thus it is impossible that other leaders introduce update commands.

Knowing that its ballot is current, a leader can treat read-only commands as above, attaching them to the highest slot number that is outstanding or chosen. The leasing technique can be integrated with adaptive timeout technique described in Section 3.

5 Exercises

This paper is accompanied by a Java package that contains a Java implementation for each of the pseudo-codes presented in this paper. Below find a list of suggestions for exercises using this code.

1. Implement the state reduction techniques for acceptors and `p1b` messages described in Section 4.1.
2. In the current implementation, ballot numbers are pairs of round numbers and leader process identifiers. If the set of leaders is fixed and ordered, then we can simplify ballot numbers. For example, if the leaders are $\{\lambda_1, \dots, \lambda_n\}$, then the ballot numbers for leader λ_i could be $i, i+n, i+2n, \dots$. Ballot number \perp could be represented as 0. Modify the Java code accordingly.
3. Implement a simple replicated bank application. The bank service maintains a set of client records, a set of accounts (a client can have zero or more accounts), and operations such as deposit, withdraw, transfer, inquiry.
4. In the Java implementation, all processes run as threads within the same Java machine, and communicate over message queue. Allow processes to run in different machines and have them communicate over TCP connections. Hint: do not consider TCP connections as reliable. If they break, have them periodically try to re-connect until successful.
5. Implement the failure detection scheme of Section 3 so that most of the time only one leader is active.
6. Improve the security of the Java implementation by securing connections. For example, one can use the TCP MD5 option for internal (non-client) connections. As a bonus, also implement SSL connections for clients.
7. Co-locate leaders and replicas as suggested in Section 4.3, and garbage collect unnecessary leader state, that is, leaders can forget about proposals for command numbers that have already been decided. Upon becoming active,

leaders do not have to start commanders for such slots either.

8. In order to increase fault tolerance, the state of acceptors and leaders can be kept on stable storage (disk). This would allow such processes to recover from crashes. Implement this. Take into consideration that a process may crash while it is saving its state.
9. Acceptors can garbage collect pvalues for decided commands that have been learned by all replicas. Implement this.
10. Implement the leasing scheme to optimize read-only operations as suggested in Section 4.4.

6 Conclusion

In this paper we presented Paxos as a collection of five kinds of processes, each with a simple operational specification. We started with an impractical but relatively easy to understand description, and then showed how various aspects can be improved to render a practical protocol.

The paper is the next in a long line of papers that describe the Paxos protocol or the experience of implementing it. A partial list follows. The View-stamped Replication protocol by Oki and Liskov [19], largely identical to Paxos, was published in 1988. Leslie Lamport's Part-time Parliament paper [13] was first written in 1989, but not published until 1998. Butler Lampson wrote a paper explaining the Part-time Parliament paper using pseudo-code in 1996 [15], and in 2001 gives an invariant-based description of various variants of Paxos [16]. In 2000, De Prisco, Lampson, and Lynch present an implementation of Paxos in the General Timed Automaton formalism [20]. Lamport wrote "Paxos made Simple" in 2001, giving a simple invariant-based explanation of the protocol [14]. In their 2003 presentation, Boichat et al. [3] give a formal account of various variants of Paxos along with extensive pseudo-code. Chandra et al. describe Google's challenges in implementing Paxos in 2007 [5]. Also in 2007, Li et al. give a novel simplified presentation of Paxos using a write-once register [17]. In his 2007 report "Paxos

made Practical” [18], David Mazières gives details of how to build replicated services using Paxos. Kirch and Amir describe their 2008 Paxos implementation experiences in [11]. Alvaro et al., in 2009 [1], describe implementing Paxos in Overlog, a declarative language.

References

- [1] P. Alvaro, T. Condie, N. Conway, J.M. Hellerstein, and R.C. Sears. I Do Declare: Consensus in a logic language. In *Proceedings of the SOSP Workshop on Networking Meets Databases (NetDB)*, 2009.
- [2] H. Attiya, A. Bar Noy, and D. Dolev. Sharing memory robustly in message passing systems. *Journal of the ACM*, 42(1):121–132, 1995.
- [3] R. Boichat, P. Dutta, S. Frolund, and R. Guerraoui. Deconstructing Paxos. *ACM SIGACT News*, 34(1), March 2003.
- [4] M. Burrows. The Chubby Lock Service for loosely-coupled distributed systems. In *7th Symposium on Operating System Design and Implementation*, Seattle, WA, November 2006.
- [5] T.D. Chandra, R. Griesemer, and J. Redstone. Paxos made live: an engineering perspective. In *Proc. of the 26th ACM Symp. on Principles of Distributed Computing*, pages 398–407, Portland, OR, May 2007. ACM.
- [6] T.D. Chandra and S. Toueg. Unreliable failure detectors for asynchronous systems. In *Proc. of the 11th ACM Symp. on Principles of Distributed Computing*, pages 325–340, Montreal, Quebec, August 1991. ACM SIGOPS-SIGACT.
- [7] M.J. Fischer, N.A. Lynch, and M.S. Patterson. Impossibility of distributed consensus with one faulty process. *J. ACM*, 32(2):374–382, April 1985.
- [8] C. Gray and D. Cheriton. Leases: an efficient fault-tolerant mechanism for distributed file cache consistency. In *Proc. of the Twelfth ACM Symp. on Operating Systems Principles*, pages 202–210, Litchfield Park, AZ, November 1989.
- [9] M. Herlihy and J. Wing. Linearizability: a correctness condition for concurrent objects. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 12(3):463 – 492, 1990.
- [10] F. Junqueira, P. Hunt, M. Konar, and B. Reed. The ZooKeeper Coordination Service (poster). In *Symposium on Operating Systems Principles (SOSP)*, 2009.
- [11] J. Kirsch and Y. Amir. Paxos for system builders. Technical Report CNDS-2008-2, Johns Hopkins University, 2008.
- [12] L. Lamport. Time, clocks, and the ordering of events in a distributed system. *CACM*, 21(7):558–565, July 1978.
- [13] L. Lamport. The part-time parliament. *Trans. on Computer Systems*, 16(2):133–169, 1998.
- [14] L. Lamport. Paxos made simple. *ACM SIGACT News (Distributed Computing Column)*, 32(4):51–58, 2001.
- [15] B. Lampson. How to build a highly available system using consensus. In O. Babaoglu and K. Marzullo, editors, *Distributed Algorithms*, volume 115 of *Lecture Notes on Computer Science*, pages 1–17. Springer-Verlag, 1996.
- [16] B.W. Lampson. The ABCD’s of Paxos. In *Proc. of the 20th ACM Symp. on Principles of Distributed Computing*, page 13, Newport, RI, 2001. ACM Press.
- [17] H.C. Li, A. Clement, A. S. Aiyer, and L. Alvisi. The Paxos register. In *Proceedings of the 26th IEEE International Symposium on Reliable Distributed Systems (SRDS 07)*, 2007.
- [18] D. Mazières. Paxos made practical. Technical Report on the web at scs.stanford.edu/~dm/home/papers/paxos.pdf, Stanford University, 2007.

- [19] B.M. Oki and B.H. Liskov. Viewstamped replication: A general primary-copy method to support highly-available distributed systems. In *Proc. of the 7th ACM Symp. on Principles of Distributed Computing*, pages 8–17, Toronto, Ontario, August 1988. ACM SIGOPS-SIGACT.
- [20] R. De Prisco, B. Lampson, and N. Lynch. Revisiting the Paxos algorithm. *Theoretical Computer Science*, 243(1-2):35–91, July 2000.
- [21] R.D. Schlichting and F.B. Schneider. Fail-stop processors: an approach to designing fault-tolerant computing systems. *Trans. on Computer Systems*, 1(3):222–238, August 1983.
- [22] F.B. Schneider. Implementing fault-tolerant services using the state machine approach: A tutorial. *ACM Computing Surveys*, 22(4):299–319, December 1990.
- [23] R.H. Thomas. A solution to the concurrency control problem for multiple copy databases. In *Proc. of COMPCON 78 Spring*, pages 88–93, Washington, D.C., February 1978. IEEE Computer Society.