

This is a specification of the *Paxos* protocol without explicit leaders or learners.

EXTENDS *TLC*, *Naturals*, *FiniteSets*, *Integers*

CONSTANTS *any*, *none*, *Replicas*, *Values*, *Ballots*, *Quorums*

VARIABLES *messages* Set of all messages sent.

VARIABLES *decision* Decided value of an acceptor.

VARIABLES *maxBallot* Maximum ballot an acceptor has seen.

VARIABLES *maxVBallot* Maximum ballot an acceptor has accepted.

VARIABLES *maxValue* Maximum value an acceptor has accepted.

Set of all possible messages.

$P1aMessage \triangleq [type : \{ "P1a" \},$
 $ballot : Ballots \setminus \{0\}]$

$P1bMessage \triangleq [type : \{ "P1b" \},$
 $ballot : Ballots,$
 $acceptor : Replicas,$
 $maxVBallot : Ballots,$
 $maxValue : Values \cup \{none\}] \quad (maxVBallot = 0) \equiv (maxValue = none)$

$P2aMessage \triangleq [type : \{ "P2a" \},$
 $ballot : Ballots,$
 $value : Values \cup \{any\}]$

$P2bMessage \triangleq [type : \{ "P2b" \},$
 $ballot : Ballots,$
 $acceptor : Replicas,$
 $value : Values]$

$Message \triangleq P1aMessage \cup P1bMessage \cup P2aMessage \cup P2bMessage$

ASSUME $PaxosAssume \triangleq$

- $\wedge \quad IsFiniteSet(Replicas)$
- $\wedge \quad any \notin Values \cup \{none\}$
- $\wedge \quad none \notin Values \cup \{any\}$
- $\wedge \quad Ballots \subseteq Nat \wedge 0 \in Ballots$
- $\wedge \quad \forall q \in Quorums : q \subseteq Replicas$
- $\wedge \quad \forall q \in Quorums : Cardinality(Replicas) \div 2 < Cardinality(q)$
- $\wedge \quad \forall q, r \in Quorums : q \cap r \neq \{\}$

$p1aMessages \triangleq \{m \in messages : m.type = "P1a"\}$ Set of all *P1a* messages sent.

$p1bMessages \triangleq \{m \in messages : m.type = "P1b"\}$ Set of all *P1b* messages sent.

$p2aMessages \triangleq \{m \in messages : m.type = "P2a"\}$ Set of all *P2a* messages sent.

$p2bMessages \triangleq \{m \in messages : m.type = "P2b"\}$ Set of all *P2b* messages sent.

$ForcedValue(M) \triangleq (CHOOSE m \in M : \forall n \in M : n.maxVBallot \leq m.maxVBallot).maxValue$

$SendMessage(m) \triangleq messages' = messages \cup \{m\}$

Phase 1a:

A proposer creates a message, which we call a "Prepare", identified with a ballot number b . Note that b is not the value to be proposed and maybe agreed on, but just a number which uniquely identifies this initial message by the proposer (to be sent to the acceptors).

While ballot number b must be greater than any ballot number used in any of the previous Prepare messages by this proposer, since the system is asynchronous and messages may be delayed and arrive out-of-order, there is no need to explicitly model this.

Then, it sends the Prepare message containing b to at least a quorum of acceptors. Note that the Prepare message only contains the ballot number b (that is, it does not have to contain the proposed value).

A Proposer should not initiate *Paxos* if it cannot communicate with at least a Quorum of Acceptors.

Some implementations may include the identity of the proposer, but that is omitted in this specification. Because while it is possible for multiple proposers to send a Prepare message with the same ballot number, only one of them can possibly receive a quorum of Promise replies. Thus, it is impossible for more than one proposer to have the same ballot number in Phase 2a.

$$\begin{aligned}
PaxosPrepare &\triangleq \\
&\wedge \text{UNCHANGED } \langle decision, maxBallot, maxVBallot, maxValue \rangle \\
&\wedge \exists b \in Ballots \setminus \{0\} : \\
&\quad SendMessage([type \mapsto \text{"P1a"}, \\
&\quad \quad ballot \mapsto b])
\end{aligned}$$

Phase 1b:

Any of the acceptors waits for a Prepare message from any of the proposers. If an acceptor receives a Prepare message, the acceptor must look at the ballot number b of the just received Prepare message.

There are two cases:

1. If b is higher than every previous proposal number received, from any of the proposers, by the acceptor, then the acceptor must return a message, which we call a "Promise", to the proposer, to ignore all future proposals having a ballot less than b . If the acceptor accepted a proposal at some point in the past, it must include the previous proposal number and the corresponding accepted value in its response to the proposer.
2. Otherwise the acceptor can ignore the received proposal. It does not have to answer in this case for *Paxos* to work.

$$\begin{aligned}
PaxosPromise &\triangleq \\
&\wedge \text{UNCHANGED } \langle decision, maxVBallot, maxValue \rangle \\
&\wedge \exists a \in Replicas, m \in p1aMessages : \\
&\quad \wedge maxBallot[a] < m.ballot \\
&\quad \wedge maxBallot' = [maxBallot \text{ EXCEPT } ![a] = m.ballot] \\
&\quad \wedge SendMessage([type \mapsto \text{"P1b"}, \\
&\quad \quad ballot \mapsto m.ballot, \\
&\quad \quad acceptor \mapsto a, \\
&\quad \quad maxVBallot \mapsto maxVBallot[a], \\
&\quad \quad maxValue \mapsto maxValue[a]])
\end{aligned}$$

Phase 2a:

If a proposer receives Promises from a quorum of acceptors, it needs to set a value v to its proposal. If any acceptors had previously accepted any proposal, then they'll have sent their values to the proposer, who now must set the value of its proposal, v , to the value associated with the highest proposal ballot reported by the acceptors, let's call it z . If none of the acceptors had accepted a proposal up to this point, then the proposer may choose the value it originally wanted to propose, say x .

The proposer sends an Accept message, (b, v) , to a quorum of acceptors with the chosen value for its proposal, v , and the ballot number b (which is the same as the number contained in the Prepare message previously sent to the acceptors). So, the Accept message is either $(b, v = z)$ or, in case none of the Acceptors previously accepted a value, $(b, v = x)$.

This Accept message should be interpreted as a "request", as in "Accept this proposal, please!".

$$\begin{aligned}
PaxosAccept &\triangleq \\
&\wedge \text{UNCHANGED } \langle decision, maxBallot, maxVBallot, maxValue \rangle \\
&\wedge \exists b \in Ballots, q \in Quorums, v \in Values : \\
&\quad \wedge \forall m \in p2aMessages : \neg(m.ballot = b) \\
&\quad \wedge \text{LET } M \triangleq \{m \in p1bMessages : m.ballot = b \wedge m.acceptor \in q\} \\
&\quad \text{IN } \quad \wedge \forall a \in q : \exists m \in M : m.acceptor = a \\
&\quad \quad \wedge \forall m \in M : m.maxValue = none \\
&\quad \quad \vee v = ForcedValue(M) \\
&\quad \wedge SendMessage([type \mapsto "P2a", \\
&\quad \quad \quad ballot \mapsto b, \\
&\quad \quad \quad value \mapsto v])
\end{aligned}$$

Phase 2b:

If an acceptor receives an Accept message, (b, v) , from a proposer, it must accept it if and only if it has not already promised (in Phase 1b of the *Paxos* protocol) to only consider proposals having a ballot greater than b .

If the acceptor has not already promised (in Phase 1b) to only consider proposals having a ballot greater than b , it should register the value v (of the just received Accept message) as the accepted value (of the protocol), and send an Accepted message to the proposer and every acceptor.

Else, it can ignore the Accept message or request.

$$\begin{aligned}
PaxosAccepted &\triangleq \\
&\wedge \text{UNCHANGED } \langle decision \rangle \\
&\wedge \exists a \in Replicas, m \in p2aMessages : \\
&\quad \wedge m.value \in Values \\
&\quad \wedge maxBallot[a] \leq m.ballot \\
&\quad \wedge maxBallot' = [maxBallot \text{ EXCEPT } ![a] = m.ballot] \\
&\quad \wedge maxVBallot' = [maxVBallot \text{ EXCEPT } ![a] = m.ballot] \\
&\quad \wedge maxValue' = [maxValue \text{ EXCEPT } ![a] = m.value] \\
&\quad \wedge SendMessage([type \mapsto "P2b", \\
&\quad \quad \quad ballot \mapsto m.ballot, \\
&\quad \quad \quad acceptor \mapsto a, \\
&\quad \quad \quad value \mapsto m.value])
\end{aligned}$$

Consensus is achieved when a majority of acceptors accept the same ballot number (rather than the same value). Because each ballot is unique to a proposer and only one value may be proposed per ballot, all acceptors that accept the same ballot thereby accept the same value.

There is no need to model the variable decision for every acceptor. In this specification, the variable decision represents the decision of any acceptor, can its value may potentially be changed multiple times. Instead, we use the consistency safety property to proof that the decision for every acceptor is the same.

$$\begin{aligned}
PaxosDecide &\triangleq \\
&\wedge \text{UNCHANGED } \langle messages, maxBallot, maxVBallot, maxValue \rangle \\
&\wedge \exists b \in Ballots, q \in Quorums : \\
&\quad \text{LET } M \triangleq \{m \in p2bMessages : m.ballot = b \wedge m.acceptor \in q\} \\
&\quad \text{IN } \wedge \forall a \in q : \exists m \in M : m.acceptor = a \\
&\quad \wedge \exists m \in M : decision' = m.value \\
\\
PaxosTypeOK &\triangleq \wedge messages \subseteq Message \\
&\wedge decision \in Values \cup \{none\} \\
&\wedge maxBallot \in [Replicas \rightarrow Ballots] \\
&\wedge maxVBallot \in [Replicas \rightarrow Ballots] \\
&\wedge maxValue \in [Replicas \rightarrow Values \cup \{none\}] \\
\\
PaxosInit &\triangleq \wedge messages = \{\} \\
&\wedge decision = none \\
&\wedge maxBallot = [r \in Replicas \mapsto 0] \\
&\wedge maxVBallot = [r \in Replicas \mapsto 0] \\
&\wedge maxValue = [r \in Replicas \mapsto none] \\
\\
PaxosNext &\triangleq \vee PaxosPrepare \\
&\vee PaxosPromise \\
&\vee PaxosAccept \\
&\vee PaxosAccepted \\
&\vee PaxosDecide \\
\\
PaxosSpec &\triangleq \wedge PaxosInit \\
&\wedge \Box [PaxosNext]_{\langle messages, decision, maxBallot, maxVBallot, maxValue \rangle} \\
&\wedge SF_{\langle messages, decision, maxBallot, maxVBallot, maxValue \rangle} (PaxosDecide)
\end{aligned}$$

Non-triviality safety property: Only proposed values can be learnt.

$$\begin{aligned}
PaxosNontriviality &\triangleq \\
&\wedge \vee decision = none \\
&\quad \vee \exists m \in p2aMessages : m.value = decision \\
&\wedge \forall m \in p1bMessages : \wedge m.maxValue \in Values \vee 0 = m.maxVBallot \\
&\quad \wedge m.maxValue = none \vee 0 < m.maxVBallot
\end{aligned}$$

Consistency safety property: At most 1 value can be learnt.

$$PaxosConsistency \triangleq \Box [decision = none]_{\langle decision \rangle}$$

From *Wikipedia*:

Note that *Paxos* is not guaranteed to terminate, and thus does not have the liveness property. This is supported by the *Fischer Lynch Paterson* impossibility result (*FLP*) which states that a consistency protocol can only have two of safety, liveness, and fault tolerance.

As *Paxos*'s point is to ensure fault tolerance and it guarantees safety, it cannot also guarantee liveness.

$$PaxosLiveness \triangleq \text{FALSE}$$

Define symmetry for faster computations.

$$PaxosSymmetry \triangleq \text{Permutations}(\text{Values}) \cup \text{Permutations}(\text{Replicas})$$
