

```

1  ┌────────────────── MODULE Paxos ───────────────────┐
  | This is a specification of the Paxos algorithm without explicit leaders or learners. It refines the
  | spec in Voting
6  EXTENDS Integers
7  └──────────────────┘
  | The constant parameters and the set Ballots are the same as in Voting.
11 CONSTANT Value, Acceptor, Quorum
13 ASSUME QuorumAssumption  $\triangleq \wedge \forall Q \in \text{Quorum} : Q \subseteq \text{Acceptor}$ 
14           $\wedge \forall Q1, Q2 \in \text{Quorum} : Q1 \cap Q2 \neq \{\}$ 
16 Ballot  $\triangleq \text{Nat}$ 
18 None  $\triangleq \text{CHOOSE } v : v \notin \text{Ballot}$ 
  | An unspecified value that is not a ballot number.
  |
  | This is a message-passing algorithm, so we begin by defining the set Message of all possible
  | messages. The messages are explained below with the actions that send them.
28 Message  $\triangleq$ 
29    $\cup$ 
30    $\cup$ 
31    $\cup$ 
32    $\cup$ 
33 └──────────────────┘
34 VARIABLE maxBal,
35           maxVBal, maxVal,
36           msgs
37           maxVBal[a], maxVal[a] is the vote with the largest
38           ballot number cast by a; it equals  $\langle -1, \text{None} \rangle$  if
39           a has not cast any vote.
40           The set of all messages that have been sent.
  |
  | NOTE: The algorithm is easier to understand in terms of the set msgs of all messages that have
  | ever been sent. A more accurate model would use one or more variables to represent the messages
  | actually in transit, and it would include actions representing message loss and duplication as well
  | as message receipt.
  |
  | In the current spec, there is no need to model message loss because we are mainly concerned with
  | the algorithm's safety property. The safety part of the spec says only what messages may be
  | received and does not assert that any message actually is received. Thus, there is no difference
  | between a lost message and one that is never received. The liveness property of the spec that we
  | check makes it clear what messages must be received (and hence either not lost or successfully
  | retransmitted if lost) to guarantee progress.
58 vars  $\triangleq \langle \text{maxBal}, \text{maxVBal}, \text{maxVal}, \text{msgs} \rangle$ 
  | It is convenient to define some identifier to be the tuple of all variables. I like to use the identifier
  | vars .
  |
  | The type invariant and initial predicate.
67 TypeOK  $\triangleq \wedge \text{maxBal} \in [\text{Acceptor} \rightarrow \text{Ballot} \cup \{-1\}]$ 

```

68  $\wedge \max VBal \in [Acceptor \rightarrow Ballot \cup \{-1\}]$   
69  $\wedge \max Val \in [Acceptor \rightarrow Value \cup \{None\}]$   
70  $\wedge msgs \subseteq Message$

73  $Init \triangleq \wedge \max Bal = [a \in Acceptor \mapsto -1]$   
74  $\wedge \max VBal = [a \in Acceptor \mapsto -1]$   
75  $\wedge \max Val = [a \in Acceptor \mapsto None]$   
76  $\wedge msgs = \{\}$

The actions. We begin with the subaction (an action that will be used to define the actions that make up the next-state action.

82  $Send(m) \triangleq msgs' = msgs \cup \{m\}$

In an implementation, there will be a leader process that orchestrates a ballot. The ballot  $b$  leader performs actions  $Phase1a(b)$  and  $Phase2a(b)$ . The  $Phase1a(b)$  action sends a phase 1a message (a message  $m$  with  $m.type = "1a"$ ) that begins ballot  $b$ .

91  $Phase1a(b) \triangleq \wedge Send([type \mapsto "1a", bal \mapsto b])$   
92  $\wedge UNCHANGED \langle \max Bal, \max VBal, \max Val \rangle$

Upon receipt of a ballot  $b$  phase 1a message, acceptor  $a$  can perform a  $Phase1b(a)$  action only if  $b > \max Bal[a]$ . The action sets  $\max Bal[a]$  to  $b$  and sends a phase 1b message to the leader containing the values of  $\max VBal[a]$  and  $\max Val[a]$ .

100  $Phase1b(a) \triangleq \wedge \exists m \in msgs :$   
101  $\wedge m.type = "1a"$   
102  $\wedge m.bal > \max Bal[a]$   
103  $\wedge \max Bal' = [\max Bal \text{ EXCEPT } ![a] = m.bal]$   
104  $\wedge Send([type \mapsto "1b", acc \mapsto a, bal \mapsto m.bal,$   
105  $\quad \quad \quad mbal \mapsto \max VBal[a], mval \mapsto \max Val[a]])$   
106  $\wedge UNCHANGED \langle \max VBal, \max Val \rangle$

The  $Phase2a(b, v)$  action can be performed by the ballot  $b$  leader if two conditions are satisfied: (i) it has not already performed a phase 2a action for ballot  $b$  and (ii) it has received ballot  $b$  phase 1b messages from some quorum  $Q$  from which it can deduce that the value  $v$  is safe at ballot  $b$ . These enabling conditions are the first two conjuncts in the definition of  $Phase2a(b, v)$ . This second conjunct, expressing condition (ii), is the heart of the algorithm. To understand it, observe that the existence of a phase 1b message  $m$  in  $msgs$  implies that  $m.mbal$  is the highest ballot number less than  $m.bal$  in which acceptor  $m.acc$  has or ever will cast a vote, and that  $m.mval$  is the value it voted for in that ballot if  $m.mbal \neq -1$ . It is not hard to deduce from this that the second conjunct implies that there exists a quorum  $Q$  such that  $ShowsSafeAt(Q, b, v)$  (where  $ShowsSafeAt$  is defined in module Voting).

The action sends a phase 2a message that tells any acceptor  $a$  that it can vote for  $v$  in ballot  $b$ , unless it has already set  $\max Bal[a]$  greater than  $b$  (thereby promising not to vote in ballot  $b$ ).

128  $Phase2a(b, v) \triangleq$   
129  $\wedge \neg \exists m \in msgs : m.type = "2a" \wedge m.bal = b$   
130  $\wedge \exists Q \in Quorum :$   
131  $\quad LET \ Q1b \triangleq \{m \in msgs : \wedge m.type = "1b"$

132  $\wedge m.acc \in Q$   
 133  $\wedge m.bal = b\}$   
 134  $Q1bv \triangleq \{m \in Q1b : m.mbal \geq 0\}$   
 135 IN  $\wedge \forall a \in Q : \exists m \in Q1b : m.acc = a$   
 136  $\wedge \vee Q1bv = \{\}$   
 137  $\vee \exists m \in Q1bv :$   
 138  $\wedge m.mval = v$   
 139  $\wedge \forall mm \in Q1bv : m.mbal \geq mm.mbal$   
 140  $\wedge Send([type \mapsto "2a", bal \mapsto b, val \mapsto v])$   
 141  $\wedge UNCHANGED \langle maxBal, maxVbal, maxVal \rangle$

The  $Phase2b(a)$  action is performed by acceptor  $a$  upon receipt of a phase  $2a$  message. Acceptor  $a$  can perform this action only if the message is for a ballot number greater than or equal to  $maxBal[a]$ . In that case, the acceptor votes as directed by the phase  $2a$  message, setting  $maxVbal[a]$  and  $maxVal[a]$  to record that vote and sending a phase  $2b$  message announcing its vote. It also sets  $maxBal[a]$  to the message's ballot number

152  $Phase2b(a) \triangleq \exists m \in msgs : \wedge m.type = "2a"$   
 153  $\wedge m.bal \geq maxBal[a]$   
 154  $\wedge maxBal' = [maxBal \text{ EXCEPT } ![a] = m.bal]$   
 155  $\wedge maxVbal' = [maxVbal \text{ EXCEPT } ![a] = m.bal]$   
 156  $\wedge maxVal' = [maxVal \text{ EXCEPT } ![a] = m.val]$   
 157  $\wedge Send([type \mapsto "2b", acc \mapsto a,$   
 158  $bal \mapsto m.bal, val \mapsto m.val])$

In an implementation, there will be learner processes that learn from the phase  $2b$  messages if a value has been chosen. The learners are omitted from this abstract specification of the algorithm.

Below are defined the next-state action and the complete spec.

169  $Next \triangleq \vee \exists b \in Ballot : \vee Phase1a(b)$   
 170  $\vee \exists v \in Value : Phase2a(b, v)$   
 171  $\vee \exists a \in Acceptor : Phase1b(a) \vee Phase2b(a)$

173  $Spec \triangleq Init \wedge \Box [Next]_{vars}$

174 |

We now define the refinement mapping under which this algorithm implements the specification in module *Voting*.

As we observed, votes are registered by sending phase  $2b$  messages. So the array *votes* describing the votes cast by the acceptors is defined as follows.

185  $votes \triangleq [a \in Acceptor \mapsto$   
 186  $\{ \langle m.bal, m.val \rangle : m \in \{ mm \in msgs : \wedge mm.type = "2b" \}$   
 187  $\wedge mm.acc = a \} \}]$

We now instantiate module *Voting*, substituting the constants *Value*, *Acceptor*, and *Quorum* declared in this module for the corresponding constants of that module *Voting*, and substituting the variable *maxBal* and the defined state function *votes* for the correspondingly-named variables of module *Voting*.

195  $V \triangleq \text{INSTANCE } Voting$

197 THEOREM  $Spec \Rightarrow V!Spec$

198

Here is a first attempt at an inductive invariant used to prove this theorem.

203  $Inv \triangleq \wedge TypeOK$

204  $\wedge \forall a \in Acceptor : \text{IF } maxVBal[a] = -1$

205  $\text{THEN } maxVal[a] = None$

206  $\text{ELSE } \langle maxVBal[a], maxVal[a] \rangle \in votes[a]$

207  $\wedge \forall m \in msgs :$

208  $\wedge (m.type = "1b") \Rightarrow \wedge maxBal[m.acc] \geq m.bal$

209  $\wedge (m.mbal \geq 0) \Rightarrow$

210  $\langle m.mbal, m.mval \rangle \in votes[m.acc]$

211  $\wedge (m.type = "2a") \Rightarrow \wedge \exists Q \in Quorum :$

212  $V!ShowsSafeAt(Q, m.bal, m.val)$

213  $\wedge \forall mm \in msgs : \wedge mm.type = "2a"$

214  $\wedge mm.bal = m.bal$

215  $\Rightarrow mm.val = m.val$

216  $\wedge V!Inv$

217