

```

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, TLC
7  ┌────────────────────────────────────────────────────────────────────────┐
8  CONSTANT Value, Acceptor, Quorum

10 ASSUME QuorumAssumption  $\triangleq$ 
11      $\wedge \forall Q \in \text{Quorum} : Q \subseteq \text{Acceptor}$ 
12      $\wedge \forall Q1, Q2 \in \text{Quorum} : Q1 \cap Q2 \neq \{\}$ 

14 Ballot  $\triangleq$  Nat
15 None  $\triangleq$  CHOOSE  $v : v \notin \text{Ballot}$ 

17 Message  $\triangleq$ 
18     [type : {"1a"}, bal : Ballot]
19      $\cup$  [type : {"1b"}, acc : Acceptor, bal : Ballot,
20         mbal : Ballot  $\cup$  { - 1 }, mval : Value  $\cup$  { None }]
21      $\cup$  [type : {"2a"}, bal : Ballot, val : Value]
22      $\cup$  [type : {"2b"}, acc : Acceptor, bal : Ballot, val : Value]
23 ┌────────────────────────────────────────────────────────────────────────┐
24 VARIABLE maxBal,      maxBal[a]: the largest ballot number a has seen
25           maxVBal,     $\langle \text{maxVBal}[a], \text{maxVal}[a] \rangle$  is the vote with the largest
26           maxVal,      ballot number cast by a; it is  $\langle -1, \text{None} \rangle$  if a has not cast any vote.
27           msgs         The set of all messages that have been sent.

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

31 vars  $\triangleq$   $\langle \text{maxBal}, \text{maxVBal}, \text{maxVal}, \text{msgs} \rangle$ 

    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.

49 TypeOK  $\triangleq$ 
50      $\wedge \text{maxBal} \in [\text{Acceptor} \rightarrow \text{Ballot} \cup \{ -1 \}]$ 
51      $\wedge \text{maxVBal} \in [\text{Acceptor} \rightarrow \text{Ballot} \cup \{ -1 \}]$ 
52      $\wedge \text{maxVal} \in [\text{Acceptor} \rightarrow \text{Value} \cup \{ \text{None} \}]$ 
53      $\wedge \text{msgs} \subseteq \text{Message}$ 
54 ┌────────────────────────────────────────────────────────────────────────┐
55 Init  $\triangleq$ 
56      $\wedge \text{maxBal} = [a \in \text{Acceptor} \mapsto -1]$ 
57      $\wedge \text{maxVBal} = [a \in \text{Acceptor} \mapsto -1]$ 

```

58 $\wedge \text{maxVal} = [a \in \text{Acceptor} \mapsto \text{None}]$
59 $\wedge \text{msgs} = \{\}$

60 |
In an implementation, there will be a leader process that orchestrates a ballot. The ballot b leader performs actions $\text{Phase1a}(b)$ and $\text{Phase2a}(b)$. The $\text{Phase1a}(b)$ action sends a phase 1a message that begins ballot b .

66 $\text{Phase1a}(b) \triangleq$
67 $\wedge \text{Send}([type \mapsto \text{"1a"}, bal \mapsto b])$
68 $\wedge \text{UNCHANGED } \langle \text{maxBal}, \text{maxVbal}, \text{maxVal} \rangle$

Upon receipt of a ballot b phase 1a message, acceptor a can perform a $\text{Phase1b}(a)$ action only if $b > \text{maxBal}[a]$. The action sets $\text{maxBal}[a]$ to b and sends a phase 1b message to the leader containing the values of $\text{maxVbal}[a]$ and $\text{maxVal}[a]$.

74 $\text{Phase1b}(a) \triangleq$
75 $\wedge \exists m \in \text{msgs} :$
76 $\wedge m.type = \text{"1a"}$
77 $\wedge m.bal > \text{maxBal}[a]$
78 $\wedge \text{maxBal}' = [\text{maxBal} \text{ EXCEPT } ![a] = m.bal]$ make promise
79 $\wedge \text{Send}([type \mapsto \text{"1b"}, acc \mapsto a, bal \mapsto m.bal,$
80 $\quad m.bal \mapsto \text{maxVbal}[a], mval \mapsto \text{maxVal}[a]])$
81 $\wedge \text{UNCHANGED } \langle \text{maxVbal}, \text{maxVal} \rangle$

83 $\text{NoBackInTime} \triangleq$
84 $\forall m \in \text{msgs} : m.type = \text{"1b"} \Rightarrow m.mbal < m.bal$

The $\text{Phase2a}(b, v)$ 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 $\text{Phase2a}(b, v)$. The 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 $\text{ShowsSafeAt}(Q, b, v)$ (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 $\text{maxBal}[a]$ greater than b (thereby promising not to vote in ballot b).

104 $\text{P2C}(b, v) \triangleq$
105 $\exists Q \in \text{Quorum} :$
106 $\text{LET } Q2bv \triangleq \{m \in \text{msgs} : m.type = \text{"2b"} \wedge m.acc \in Q \wedge m.bal < b\}$
107 $\text{IN } \vee Q2bv = \{\}$
108 $\vee \exists m \in Q2bv :$
109 $\wedge m.val = v$
110 $\wedge \forall mm \in Q2bv : m.bal \geq mm.bal$
112 $\text{Phase2a}(b, v) \triangleq$
113 $\wedge \neg \exists m \in \text{msgs} : m.type = \text{"2a"} \wedge m.bal = b$
114 $\wedge \exists Q \in \text{Quorum} :$
115 $\text{LET } Q1b \triangleq \{m \in \text{msgs} : m.type = \text{"1b"} \wedge m.acc \in Q \wedge m.bal = b\}$

```

116       $Q1bv \triangleq \{m \in Q1b : m.mbal \geq 0\}$ 
117      IN  $\wedge \forall a \in Q : \exists m \in Q1b : m.acc = a$ 
118       $\wedge \vee Q1bv = \{\}$ 
119       $\vee \exists m \in Q1bv :$ 
120       $\wedge m.mval = v$ 
121       $\wedge \forall mm \in Q1bv : m.mbal \geq mm.mbal$ 
122       $\wedge Send([type \mapsto "2a", bal \mapsto b, val \mapsto v])$ 
123       $\wedge Assert(P2C(b, v), "P2C Fails!")$ 
124       $\wedge UNCHANGED \langle maxBal, maxVbal, maxVal \rangle$ 

```

The *Phase2b(a)* action is performed by acceptor *a* upon receipt of a phase 2*a* 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 2*a* message, setting *maxVbal[a]* and *maxVal[a]* to record that vote and sending a phase 2*b* message announcing its vote.

Note: It also sets *maxBal[a]* to the message's ballot number. Otherwise,

- (1) *NoBackInTime* for *Phase1b* does not hold.
- (2) "Non-Increasing Error" assertion in *Phase2b(a)* fails.
- (3) *P2C* assertion for *Phase2a* does not hold ???

```

138  Phase2b(a)  $\triangleq$ 
139       $\wedge \exists m \in msgs :$ 
140       $\wedge m.type = "2a"$ 
141       $\wedge m.bal \geq maxBal[a]$ 
142       $\wedge maxBal' = [maxBal \text{ EXCEPT } ![a] = m.bal]$ 
143       $\wedge maxVbal' = [maxVbal \text{ EXCEPT } ![a] = m.bal]$ 
144       $\wedge Assert(maxVbal'[a] \geq maxVbal[a], "Non-Increasing Error!")$ 
145       $\wedge maxVal' = [maxVal \text{ EXCEPT } ![a] = m.val]$ 
146       $\wedge Send([type \mapsto "2b", acc \mapsto a, bal \mapsto m.bal, val \mapsto m.val])$ 
147       $\wedge UNCHANGED \langle \rangle$ 

```

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

```

153 |-----|
154  Next  $\triangleq$ 
155       $\vee \exists b \in Ballot :$ 
156       $\vee Phase1a(b)$ 
157       $\vee \exists v \in Value : Phase2a(b, v)$ 
158       $\vee \exists a \in Acceptor : Phase1b(a) \vee Phase2b(a)$ 

```

```

160  Spec  $\triangleq Init \wedge \Box [Next]_{vars}$ 
161 |-----|

```

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 2*b* messages. So the array *votes* describing the votes cast by the acceptors is defined as follows.

```

172  votes  $\triangleq [a \in Acceptor \mapsto$ 
173       $\{ \langle m.bal, m.val \rangle : m \in \{ mm \in msgs : \wedge mm.type = "2b" \}$ 

```

174 $\wedge mm.acc = a\}}]$

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*.

181 $V \triangleq \text{INSTANCE } Voting$

183 THEOREM $Spec \Rightarrow V!Spec$

184 $\left| \begin{array}{l} \text{Here is a first attempt at an inductive invariant used to prove this theorem.} \\ 189 \text{ } Inv \triangleq \wedge TypeOK \\ 190 \quad \wedge \forall a \in Acceptor : \text{IF } maxVBal[a] = -1 \\ 191 \quad \quad \quad \text{THEN } maxVal[a] = None \\ 192 \quad \quad \quad \text{ELSE } \langle maxVBal[a], maxVal[a] \rangle \in votes[a] \\ 193 \quad \wedge \forall m \in msgs : \\ 194 \quad \quad \wedge (m.type = "1b") \Rightarrow \wedge maxBal[m.acc] \geq m.bal \\ 195 \quad \quad \quad \wedge (m.mbal \geq 0) \Rightarrow \\ 196 \quad \quad \quad \langle m.mbal, m.mval \rangle \in votes[m.acc] \\ 197 \quad \quad \wedge (m.type = "2a") \Rightarrow \wedge \exists Q \in Quorum : \\ 198 \quad \quad \quad V!ShowsSafeAt(Q, m.bal, m.val) \\ 199 \quad \quad \quad \wedge \forall mm \in msgs : \wedge mm.type = "2a" \\ 200 \quad \quad \quad \quad \wedge mm.bal = m.bal \\ 201 \quad \quad \quad \quad \Rightarrow mm.val = m.val \\ 202 \quad \wedge V!Inv \end{array} \right|$

203