

```

1  |----- MODULE PaxosTuple -----|
   | This is a specification of the Paxos algorithm without explicit leaders or learners. It refines the |
   | spec in Voting NB: This specification was written before TLAPS had any support for records |
   | and therefore uses tuples for representing messages. |
8  | EXTENDS Integers, Sets |
9  |-----|
   | The constant parameters and the set Ballots are the same as in Voting. |
13 | CONSTANT Value, Acceptor, Quorum
15 | ASSUME QuorumAssumption  $\triangleq$   $\wedge \forall Q \in \text{Quorum} : Q \subseteq \text{Acceptor}$ 
16 |                                      $\wedge \forall Q1, Q2 \in \text{Quorum} : Q1 \cap Q2 \neq \{\}$ 
18 | Ballot  $\triangleq$  Nat
20 | None  $\triangleq$  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. |
30 | Message  $\triangleq$ 
31 |     { "1a" }  $\times$  Ballot
32 |      $\cup$  { "1b" }  $\times$  Acceptor  $\times$  Ballot  $\times$  (Ballot  $\cup$  { - 1 })  $\times$  (Value  $\cup$  { None })
33 |      $\cup$  { "2a" }  $\times$  Ballot  $\times$  Value
34 |      $\cup$  { "2b" }  $\times$  Acceptor  $\times$  Ballot  $\times$  Value
35 |-----|
36 | VARIABLE maxBal,
37 |           maxVBal,  $\langle \text{maxVBal}[a], \text{maxVal}[a] \rangle$  is the vote with the largest
38 |           maxVal, ballot number cast by a; it equals  $\langle -1, \text{None} \rangle$  if
39 |           a has not cast any vote.
40 |           msgs 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. |
60 | 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. |
69 | TypeOK  $\triangleq$   $\wedge \text{maxBal} \in [\text{Acceptor} \rightarrow \text{Ballot} \cup \{ -1 \}]$ 

```

70 $\wedge \max VBal \in [Acceptor \rightarrow Ballot \cup \{-1\}]$
 71 $\wedge \max Val \in [Acceptor \rightarrow Value \cup \{None\}]$
 72 $\wedge msgs \subseteq Message$

75 $Init \triangleq \wedge \max Bal = [a \in Acceptor \mapsto -1]$
 76 $\wedge \max VBal = [a \in Acceptor \mapsto -1]$
 77 $\wedge \max Val = [a \in Acceptor \mapsto None]$
 78 $\wedge msgs = \{\}$

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

84 $Send(m) \triangleq msgs' = msgs \cup \{m\}$

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 .

93 $Phase1a(b) \triangleq \wedge Send(\langle "1a", b \rangle)$
 94 $\wedge UNCHANGED \langle \max Bal, \max VBal, \max Val \rangle$

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 $\max Val[a]$.

102 $Phase1b(a) \triangleq \wedge \exists m \in msgs :$
 103 $\wedge m[1] = "1a"$
 104 $\wedge m[2] > \max Bal[a]$
 105 $\wedge \max Bal' = [\max Bal \text{ EXCEPT } ![a] = m[2]]$
 106 $\wedge Send(\langle "1b", a, m[2],$
 107 $\max VBal[a], \max Val[a] \rangle)$
 108 $\wedge UNCHANGED \langle \max VBal, \max Val \rangle$

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).

130 $Phase2a(b, v) \triangleq$
 131 $\wedge \neg \exists m \in msgs : m[1] = "2a" \wedge m[3] = b$
 132 $\wedge \exists Q \in Quorum :$
 133 $LET Q1b \triangleq \{m \in msgs : \wedge m[1] = "1b"$

134 $\wedge m[2] \in Q$
 135 $\wedge m[3] = b\}$
 136 $Q1bv \triangleq \{m \in Q1b : m[4] \geq 0\}$
 137 IN $\wedge \forall a \in Q : \exists m \in Q1b : m[2] = a$
 138 $\wedge \vee Q1bv = \{\}$
 139 $\vee \exists m \in Q1bv :$
 140 $\wedge m[5] = v$
 141 $\wedge \forall mm \in Q1bv : m[4] \geq mm[4]$
 142 $\wedge Send(\langle "2a", b, v \rangle)$
 143 $\wedge UNCHANGED \langle maxBal, maxVBal, maxVal \rangle$

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 *maxBVal[a]* and *maxVal[a]* to record that vote and sending a phase 2*b* message announcing its vote. It also sets *maxBal[a]* to the message's ballot number

154 $Phase2b(a) \triangleq \exists m \in msgs : \wedge m[1] = "2a"$
 155 $\wedge m[2] \geq maxBal[a]$
 156 $\wedge maxBal' = [maxBal \text{ EXCEPT } ![a] = m[2]]$
 157 $\wedge maxVBal' = [maxVBal \text{ EXCEPT } ![a] = m[2]]$
 158 $\wedge maxVal' = [maxVal \text{ EXCEPT } ![a] = m[3]]$
 159 $\wedge Send(\langle "2b", a, m[2], m[3] \rangle)$

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.

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

170 $Next \triangleq \vee \exists b \in Ballot : \vee Phase1a(b)$
 171 $\vee \exists v \in Value : Phase2a(b, v)$
 172 $\vee \exists a \in Acceptor : Phase1b(a) \vee Phase2b(a)$
 174 $Spec \triangleq Init \wedge \Box[Next]_{vars}$

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.

186 $votes \triangleq [a \in Acceptor \mapsto$
 187 $\{ \langle m[3], m[4] \rangle : m \in \{ mm \in msgs : \wedge mm[1] = "2b"$
 188 $\wedge mm[2] = a \} \}]$

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*.

196 $V \triangleq \text{INSTANCE } Voting$

198 THEOREM $Spec \Rightarrow V!Spec$

199

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

204 $Inv \triangleq \wedge TypeOK$

205 $\wedge \forall a \in Acceptor : \text{IF } maxVBal[a] = -1$

206 $\text{THEN } maxVal[a] = None$

207 $\text{ELSE } \langle maxVBal[a], maxVal[a] \rangle \in votes[a]$

208 $\wedge \forall m \in msgs :$

209 $\wedge (m.type = "1b") \Rightarrow \wedge maxBal[m.acc] \geq m.bal$

210 $\wedge (m.mbal \geq 0) \Rightarrow$

211 $\langle m.mbal, m.mval \rangle \in votes[m.acc]$

212 $\wedge (m.type = "2a") \Rightarrow \wedge \exists Q \in Quorum :$

213 $V!ShowsSafeAt(Q, m.bal, m.val)$

214 $\wedge \forall mm \in msgs : \wedge mm.type = "2a"$

215 $\wedge mm.bal = m.bal$

216 $\Rightarrow mm.val = m.val$

217 $\wedge V!Inv$

218