──────── MODULE *Paxos* ────────

This is a high-level specification of the *Paxos* consensus algorithm. It refines the spec in module *Voting*, which you should read before reading this module. In the *Paxos* consensus algorithm, acceptors communicate by sending messages. There are additional processes called leaders. The specification here essentially considers there to be a separate leader for each ballot number. We can consider "leader" to be a role, where in an implementation there will be a finite number of leader processes each of which plays infinitely many of these leader roles.

Note: The algorithm described here is the *Paxos* consensus algorithm. It is the crucial component of the *Paxos* algorithm, which implements a fault-tolerant state machine using a sequence of instances of the consensus algorithm. The *Paxos* algorithm is sometimes called *MultiPaxos*, with the *Paxos* consensus algorithm being incorrectly called the *Paxos* algorithm. I'm afraid I have contributed to this confusion by being lazy and calling the module "*Paxos*" instead of "*PaxosConsensus*". But incarnations of this module have already appeared, so I'm reluctant to change its name now.

EXTENDS *Integers*

The constants and the assumptions about them are the same as for the Voting algorithm. However, the second conjunct of the assumption, which asserts that any two quorums have a non-empty intersection, is not needed for the *Paxos* consensus algorithm to implement the *Voting* algorithm. The *Voting* algorithm, and it, do not satisfy consensus without that assumption.

CONSTANTS $Value$, $Acceptor$, $Quorum$

ASSUME $\quad \wedge \forall\, Q \in Quorum : Q \subseteq Acceptor$
$\qquad\quad \wedge \forall\, Q1,\, Q2 \in Quorum : Q1 \cap Q2 \neq \{\}$

$Ballot \;\triangleq\; Nat$

$None \;\triangleq\;$ CHOOSE $v : v \notin Ballot$

  This defines *None* to be an unspecified value that is not a ballot number.

We now define *Message* toe be the set of all possible messages that can be sent in the algorithm. In TLA+, the expression

(1) $[type\; \mathbb{V}\; \text{"1a"},\, bal\; \mathbb{V}\; b]$

is a record $r$ with two components, a *type* component, $r.type$, that equals "1a" and whose *bal* component, $r.bal$, that equals $b$. The expression

(2) $[type : f\text{"1a"}g,\, bal : Ballot]$

is the set of all records $r$ with a components *type* and *bal* such that $r.type$ is an element of $f\text{"1a"}g$ and $r.bal$ is an element of *Ballot*. Since "1a" is the only element of $f\text{"1a"}g$, formula (2) is the set of all elements (1) such that $b\; 2\; Ballot$.

The function of each type of message in the set *Message* is explained below with the action that can send it.

$Message \;\triangleq\;$
$\qquad [type\; :\; \{\,\text{"1a"}\,\},\, bal\; :\; Ballot]$
$\quad \cup \quad [type\; :\; \{\,\text{"1b"}\,\},\, acc : Acceptor,\, bal : Ballot,$
$\qquad\quad mbal : Ballot \cup \{\,-1\,\},\, mval : Value \cup \{None\}]$
$\quad \cup \quad [type\; :\; \{\,\text{"2a"}\,\},\, bal : Ballot,\, val : Value]$
$\quad \cup \quad [type\; :\; \{\,\text{"2b"}\,\},\, acc : Acceptor,\, bal : Ballot,\, val : Value]$

VARIABLES $maxBal$, $maxVBal$, $maxVBal$

.            ,

, $mmaxVBal$

/.          . $maxVBal$

The ballot *b* leader can perform actions *Phase1a(b)* and *Phase2a(b)*. In the *Phase1a(b)* action, it sends to all acceptors a phase 1*a* message (a message *m* with *m.type* = "1a") that begins ballot *b*. Remember that the same process can perform the role of leader for many different ballot numbers *b*. In practice, it will stop playing the role of leader of ballot *b* when it begins a higher-numbered ballot. (Remember the definition of [*type* $\mapsto$ "1a", *bal* $\mapsto$ *b*] from the comment preceding the definition of *Message*.)

$Phase1a(b) \triangleq \land Send([type \mapsto \text{``1a''}, bal \mapsto b])$
$\qquad\qquad\qquad\;\land \text{UNCHANGED } \langle maxBal, maxVBal, maxVal \rangle$

Note that there is no enabling condition to prevent sending the phase 1*a* message a second time. Since messages are never removed from msg, performing the action a second time leaves msg and all the other spec variables unchanged, so it's a stuttering step. Since stuttering steps are always allowed, there's no reason to try to prevent them.

Upon receipt of a ballot *b* phase 1*a* message, acceptor a can perform a *Phase1b(a)* action only if *b* > *maxBal[a]*. The action sets *maxBal[a]* to *b* and sends a phase 1*b* message to the leader containing the values of *maxVBal[a]* and *maxVal[a]*. This action implements the *IncreaseMaxBal(a, b)* action of the *Voting* algorithm for *b* = *m.bal*.

$Phase1b(a) \triangleq$
$\quad \land \exists m \in msgs :$
$\qquad \land m.type = \text{``1a''}$
$\qquad \land m.bal > maxBal[a]$
$\qquad \land maxBal' = [maxBal \text{ EXCEPT } ![a] = m.bal]$
$\qquad \land Send([type \mapsto \text{``1b''}, acc \mapsto a, bal \mapsto m.bal,$
$\qquad\qquad\qquad\quad mbal \mapsto maxVBal[a], mval \mapsto maxVal[a]])$
$\quad \land \text{UNCHANGED } \langle maxVBal, maxVal \rangle$

In the *Phase2a(b, v)* action, the ballot *b* leader sends a type "2a" message asking the acceptors to vote for *v* in ballot number *b*. The enabling conditions of the action–its first two conjuncts–ensure that three of the four enabling conditions of action *VoteFor(a, b, v)* in module *Voting* will be true when acceptor a receives that message. Those three enabling conditions are the second through fourth conjuncts of that action.

The first conjunct of *Phase2a(b, v)* asserts that at most one phase 2*a* message is ever sent for ballot *b*. Since an acceptor will vote for a value in ballot *b* only when it receives the appropriate phase 2*a* message, the phase 2*a* message sent by this action this ensures that these two enabling conjuncts of *VoteFor(a, b, v)* will be true forever:

$\quad \land \forall vt \in votes[a] : vt[1] \neq b$
$\quad \land \forall c \in Acceptor \setminus fag : \forall vt \in votes[c] : (vt[1] = b) \Rightarrow (vt[2] = v)$

The second conjunct of the *Phase2a(b, v)* action is the heart of the *Paxos* consensus algorithm. It's a bit complicated, but I've tried a number of times to write it in *English*, and it's much easier to understand when written in mathematics. The LET /IN construct locally defines *Q1* to be the set of phase 1*b* messages sent in ballot number *b* by acceptors in quorum *Q*; and it defines *Q1bv* to be the subset of those messages indicating that the sender had voted in some ballot (which must have been numbered less than *b*). You should study the IN clause to convince yourself that it equals *ShowsSafeAt(Q, b, v)*, defined in module *Voting*, using the values of *maxBal[a]*, *maxVBal[a]*, and *maxVal[a]a* sent in its phase 1*b* message to describe what votes it had cast when it sent that message. Moreover, since *a* will no longer cast any votes in ballots numbered less than *b*, the IN clause implies that *ShowsSafeAt(Q, b, v)* is still true and will remain true forever. Hence, this conjunct of *Phase2a(b, v)* checks the enabling condition

$\quad \land \exists Q \in Quorum : ShowsSafeAt(Q, b, v)$

3

of module *Voting*'s *VoteFor*(*a*, *b*, *v*) action.

The type "2a" message sent by this action therefore tells every acceptor *a* that, when it receives the message, all the enabling conditions of *VoteFor*(*a*, *b*, *v*) but the first, *maxBal*[*a*]   *b*, are satisfied.

$Phase2a(b, v) \triangleq$
   $\land \neg \exists m \in msgs$  $: m.type = $ "2a" $\land m.bal = b$
   $\land \exists Q \in Quorum :$
       LET $Q1b \triangleq \{m \in msgs$  $: \land m.type = $ "1b"
                           $\land m.acc \in Q$
                           $\land m.bal = b\}$
           $Q1bv \triangleq \{m \in Q1b : m.mbal \geq 0\}$
      IN    $\land \forall a \in Q : \exists m \in Q1b : m.acc = a$
           $\land \lor Q1bv = \{\}$
              $\lor \exists m$  $\in Q1bv :$
                  $\land m.mval = v$
                  $\land \forall mm \in Q1bv : m.mbal \geq mm.mbal$
   $\land Send([type \mapsto $ "2a"$, bal \mapsto b, val \mapsto v])$
   $\land$ UNCHANGED $\langle maxBal, maxVBal, maxVal \rangle$

The *Phase*2*b*(*a*) action describes what acceptor *a* does when it receives a phase 2a message *m*, which is sent by the leader of ballot *m.bal* asking acceptors to vote for *m.val* in that ballot. Acceptor *a* acts on that request, voting for *m.val* in ballot number *m.bal*, iff *m.bal*   *maxBal*[*a*], which means that *a* has not participated in any ballot numbered greater than *m.bal*. Thus, this enabling condition of the *Phase*2*b*(*a*) action together with the receipt of the phase 2*a* message *m* implies that the *VoteFor*(*a*, *m.bal*, *m.val*) action of module *Voting* is enabled and can be executed. The *Phase*2*b*(*a*) message updates *maxBal*[*a*], *maxVBal*[*a*], and *maxVal*[*a*] so their values mean what they were claimed to mean in the comments preceding the variable declarations.

$Phase2b(a) \triangleq$
  $\exists m \in msgs :$
     $\land m.type = $ "2a"
     $\land m.bal \geq maxBal[a]$
     $\land maxBal' = [maxBal$ EXCEPT $![a] = m.bal]$
     $\land maxVBal' = [maxVBal$ EXCEPT $![a] = m.bal]$
     $\land maxVal' = [maxVal$ EXCEPT $![a] = m.val]$
     $\land Send([type \mapsto $ "2b"$, acc \mapsto a,$
            $bal \mapsto m.bal, val \mapsto m.val])$

The definitions of *Next* and *Spec* are what we expect them to be.

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

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

We define *votes* to be the function such that *votes*[*a*] is the set of pairs *hb*, *vi* such that acceptor *a* has voted for *v* in ballot number *b* by sending executing the *Phase*2*b*(*a*) action to send the appropriate type "2b" message. The *Paxos* consensus algorithm implements the *Voting* algorithm by implementing the variable *votes* of module *Voting* with the expression *votes* of the current module, and implementing the variable *maxBal* of module *Voting* with the variable of the same name of the current module.

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

The following INSTANCE statement omits the redundant clause

  WITH *votes*   *votes*, *maxBal*   *maxBal*,
    *Value*   *Value*, *Acceptor*   *Acceptor*, *Quorum*   *Quorum*

$V \;\triangleq\;$ INSTANCE *Voting*

The inductive invariant *Inv* explains why the *Paxos* consensus algorithm implements the *Voting* algorithm. It is defined after the INSTANCE statement because it uses the operator *V !ShowsSafeAt* imported by that statement.

$Inv \;\triangleq\;$
  $\land\; TypeOK$
  $\land\; \forall\, a \in Acceptor : maxBal[a] \geq maxVBal[a]$
  $\land\; \forall\, a \in Acceptor :$ IF $maxVBal[a] \;=\; -1$
                            THEN $maxVal[a] \;=\; None$
                            ELSE $\langle maxVBal[a],\, maxVal[a]\rangle \in votes[a]$
  $\land\; \forall\, m \in msgs :$
    $\land\; (m.type \;=\; \text{"1b"}) \Rightarrow \;\land\; maxBal[m.acc] \geq m.bal$
                                 $\land\; (m.mbal \geq 0) \Rightarrow$
                                     $\langle m.mbal,\, m.mval\rangle \in votes[m.acc]$
    $\land\; (m.type \;=\; \text{"2a"}) \Rightarrow \;\land\; \exists\, Q \in Quorum :$
                                        $V !ShowsSafeAt(Q,\, m.bal,\, m.val)$
                               $\land\; \forall\, mm \in msgs : \;\land\; mm.type \;=\; \text{"2a"}$
                                                $\land\; mm.bal \;\; = \; m.bal$
                                                $\Rightarrow mm.val \;=\; m.val$
    $\land\; (m.type \;=\; \text{"2b"}) \Rightarrow \;\land\; maxVBal[m.acc] \geq m.bal$
                                 $\land\; \exists\, mm \in msgs : \;\land\; mm.type \;=\; \text{"2a"}$
                                                $\land\; mm.bal \;\; = \; m.bal$
                                                $\land\; mm.val \;\; = \; m.val$

The following two theorems assert that *Inv* is an invariant of the *Paxos* cconsensus algorithm and that this algorithm implements the *Voting* algorithm with the declared variables and constants of that algorithm implemented by the correspondingly-named expressions in the current module.

THEOREM $Invariance \;\triangleq\; Spec \Rightarrow \Box Inv$

THEOREM $Implementation \;\triangleq\; Spec \Rightarrow V !Spec$

The ASSUME statement of this module trivial implies trivially implies the instantiated version of the ASSUME statement of module *Voting*. (Because the INSTANCE statement substitutes the constants of the current module for the constants of the same name in module *Voting*, the imported assumption is the same as the assumption of the current module.) Hence, this theorem imported from module *Voting* is true in the current module

THEOREM $V!Implementation \triangleq V!Spec \supset V!C!Spec$

Theorems *Implementation* and $V!Implementation$ imply

THEOREM $Spec \supset V!C!Spec$

This theorem asserts that the *Paxos* consensus algorithm implements the Consensus specification by substituting for the variable *chosen* of the *Consensus* specification the value $V!chosen$ of the current module. The expression $V!chosen$ is obtained by substituting the expression *votes* of the current module for the variable *votes* of module *Voting* in the expression *chosen* of module *Voting*.

In other words, as we should expect, "implements" is a transitive relation–under a suitable understanding of what transitivity means in this situation.

This current module is distributed with two models, *TinyModel* and *SmallModel*. *SmallModel* is the same as the model by that name for the Voting specification. *TinyModel* is the same except it defines *Ballot* to contain only two elements. Run *TLC* on these models. You should find that it takes a couple of seconds to run *TinyModel* and two or three minutes to run *SmallModel*.

Next, try the same thing you did with the *Voting* algorithm: Modify the models so the assumption that any pair of quorums has an element in common is no longer true. (Again, it's best to modify clones of the models.) This time, running *TLC* will not find an error. The correctness of theorems *Invariance* and *Implementation* does not depend on that assumption. The *Paxos* consensus algorithm still correctly implements the *Voting* algorithm; but the *Voting* algorithm is incorrect if the assumption does not hold.

Now go back to the original *SmallModel*, in which the quorum assumption holds. The sets *Acceptor* and *Value* are symmetry sets for the spec. (See the "Model Values and Symmetry" help page to find out what that means.) Try editing the values substituted for *Acceptor* and/or *Value* by selecting the "Symmetry set" option and comparing the number of reachable states *TLC* found and the time it took. (Remember to use cloned models.)

When you have other things to do while *TLC* is running, try increasing the size of the model a very little bit at a time and see how the running time increases. You'll find that it increases exponentially with the numbers of acceptors, values, and ballots.

Fortunately, exhaustively checking a small model is very effective at finding errors. Since the *Paxos* consensus algorithm has been proved correct, and that proof has been read by many people, I'm sure that the basic algorithm is correct. Checking this spec on *SmallModel* makes me quite confident that there are no "coding errors" in this TLA+ specification of the algorithm.

For checking safety properties, *TLC* can obtain close to linear speedup using dozens of processors. After designing a new distributed algorithm, you will have plenty of time to run *TLC* while the algorithm is being implemented and the implementation tested. Use that time to run it for as long as you can on the largest *machine(s)* that you can. Testing the implementation is unlikely to find subtle errors in the algorithm.