──────────── MODULE *CASPaxos* ────────────

This is a high-level specification of the *CASPaxos* algorithm from the paper "" by.

*TODO*: It refines the spec in module Voting.

8 ├────────────────────────────────────────────────────

9 EXTENDS *Integers*

10 ├────────────────────────────────────────────────────

11 CONSTANTS

12    *Value*,        the set of values to be proposed and chosen from

13    *Acceptor*,     the set acceptors

14    *Quorum*        the quorum system on acceptors

16  $None \triangleq$ CHOOSE $v : v \notin Value$

18  ASSUME  $\land \forall\, Q \in Quorum : Q \subseteq Acceptor$
19          $\land \forall\, Q1, Q2 \in Quorum : Q1 \cap Q2 \neq \{\}$

21  $Ballot \triangleq Nat$

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

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

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

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

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

42  $Message \triangleq$
43          $[type : \{\text{"1a"}\}, bal : Ballot]$
44      $\cup$  $[type : \{\text{"1b"}\}, acc : Acceptor, bal : Ballot,$
45           $mbal : Ballot \cup \{-1\}, mval : Value \cup \{None\}]$
46      $\cup$  $[type : \{\text{"2a"}\}, bal : Ballot, val : Value]$
47      $\cup$  $[type : \{\text{"2b"}\}, acc : Acceptor, bal : Ballot, val : Value]$

48 ├────────────────────────────────────────────────────

$maxBal -$ Is the same as the variable of that name in the Voting algorithm.
$maxVBal$
$maxVVal -$   As in the Voting algorithm, a vote is a $\langle ballot, value \rangle$ pair.   The pair $\langle maxVBal[a], maxVVal[a] \rangle$ is the vote with the largest ballot number cast by acceptor a . It equals $\langle -1, None \rangle$ if a has not cast any vote.

58  VARIABLES

59     $maxBal$,

60     $maxVBal$,

61     $maxVVal$,

62     $msgs$,        the set of all messages that have been sent

63        $cas$             $cas[b \in Ballot]: [cmpVal : Value \cup \{None\}, swapVal : Value \cup \{None\}\rangle \setminus * \ TODO$

65   $vars \triangleq \langle maxBal, maxVBal, maxVVal, msgs, cas \rangle$

66 ⊢─────────────────────────────────────────────────────────

67   $TypeOK \triangleq \land maxBal \ \in [Acceptor \to Ballot \cup \{-1\}]$
68                    $\land maxVBal \in [Acceptor \to Ballot \cup \{-1\}]$
69                    $\land maxVVal \in [Acceptor \to Value \cup \{None\}]$
70                    $\land msgs \subseteq Message$
71                    $\land cas \in [Ballot \to [cmpVal : Value \cup \{None\},$
72                                   $swapVal : Value \cup \{None\}]]$

73 ⊢─────────────────────────────────────────────────────────

74   $Init \triangleq \land maxBal \ \ = [a \in Acceptor \mapsto -1]$
75            $\land maxVBal = [a \in Acceptor \mapsto -1]$
76            $\land maxVVal = [a \in Acceptor \mapsto None]$
77            $\land msgs = \{\}$
78            $\land cas = [b \in Ballot \mapsto [cmpVal \mapsto None, swapVal \mapsto None]]$     $TODO:$

79 ⊢─────────────────────────────────────────────────────────

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

81 ⊢─────────────────────────────────────────────────────────

The leader of ballot $b \in Ballot$ issues an $CAS(cmpVal, swapVal)$ operation by sending a $Phase1a$ message.

86   $Phase1a(b, cmpVal, swapVal) \triangleq$
87       $\land \ \ \ Send([type \mapsto \text{"1a"}, bal \mapsto b])$
88       $\land \ \ \ cas' = [cas \text{ EXCEPT } ![b] = [cmpVal \mapsto cmpVal, swapVal \mapsto swapVal]]$
89       $\land \ \ \text{UNCHANGED} \ \langle maxBal, maxVBal, maxVVal \rangle$

The acceptor $a \in Acceptor$ receives a $Phase1a$ message and sends back a $Phase1b$ message.

$TODO$: This action implements the $IncreaseMaxBal(a, b)$ action of the Voting algorithm for $b = m.bal$.

97   $Phase1b(a) \triangleq$
98      $\land \exists m \in msgs :$
99         $\land m.type = \text{"1a"}$
100        $\land m.bal > maxBal[a]$
101        $\land maxBal' = [maxBal \text{ EXCEPT } ![a] = m.bal]$
102        $\land Send([type \ \mapsto \text{"1b"}, acc \mapsto a, bal \mapsto m.bal,$
103                  $mbal \mapsto maxVBal[a], mval \mapsto maxVVal[a]])$
104      $\land \text{UNCHANGED} \ \langle maxVBal, maxVVal, cas \rangle$

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 2a 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 2a message, the phase 2a message sent by this action this ensures that these two enabling conjuncts of $VoteFor(a, b, v)$ will be true forever:

149  $Phase2a(b, v) \triangleq$
150    $\wedge \neg \exists m \in msgs : m.type =$ "2a" $\wedge m.bal = b$
151    $\wedge \exists Q \in Quorum :$
152      LET $Q1b \triangleq \{m \in msgs : \wedge m.type =$ "1b"
153                                      $\wedge m.acc \in Q$
154                                      $\wedge m.bal = b\}$
155          $Q1bv \triangleq \{m \in Q1b : m.mbal \geq 0\}$
156      IN   $\wedge \forall a \in Q : \exists m \in Q1b : m.acc = a$
157           $\wedge \vee Q1bv = \{\}$
158             $\vee \exists m \in Q1bv :$
159                 $\wedge m.mval = v$
160                 $\wedge \forall mm \in Q1bv : m.mbal \geq mm.mbal$
161    $\wedge Send([type \mapsto$ "2a"$, bal \mapsto b, val \mapsto v])$
162    $\wedge$ UNCHANGED $\langle maxBal, maxVBal, maxVVal \rangle$

176  $Phase2b(a) \triangleq$
177    $\exists m \in msgs :$
178      $\wedge m.type =$ "2a"
179      $\wedge m.bal \geq maxBal[a]$
180      $\wedge maxBal' = [maxBal$ EXCEPT $![a] = m.bal]$
181      $\wedge maxVBal' = [maxVBal$ EXCEPT $![a] = m.bal]$
182      $\wedge maxVVal' = [maxVVal$ EXCEPT $![a] = m.val]$

183      $\land Send([type \mapsto \text{``2b''}, acc \mapsto a,$
184                 $bal \mapsto m.bal, val \mapsto m.val])$

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

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

193   $Spec \triangleq Init \land \Box[Next]_{vars}$
194 ├──────────────────────────────────────────────────────────────┤

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.

240 └──────────────────────────────────────────────────────────────┘