─────────── MODULE *PaxosCommit* ───────────

This module specifies the *Paxos Commit* algorithm. We specify only safety properties, not liveness properties. We simplify the specification in the following ways. \begin{*itemize*} \item As in the specification of module $TwoPhase$, and for the same reasons, we let the variable $msgs$ be the set of all messages that have ever been sent. If a message is sent to a set of recipients, only one copy of the message appears in $msgs$.

\item We do not explicitly model the receipt of messages. If an operation can be performed when a process has received a certain set of messages, then the operation is represented by an action that is enabled when those messages are in the set $msgs$ of sent messages. (We are specifying only safety properties, which assert what events can occur, and the operation can occur if the messages that enable it have been sent.)

\item We do not model leader selection. We define actions that the current leader may perform, but do not specify who performs them. \end{*itemize*}

As in the specification of Two-Phase commit in module $TwoPhase$, we have *RMs* spontaneously issue Prepared messages and we ignore $Prepare$ messages.

EXTENDS *Integers*

$Maximum(S) \triangleq$

If $S$ is a set of numbers, then this define $Maximum(S)$ to be the maximum of those numbers, or $-1$ if $S$ is empty.

LET $Max[T \in \text{SUBSET } S] \triangleq$
      IF $T = \{\}$ THEN $-1$
            ELSE LET $n \quad \triangleq$ CHOOSE $n \in T$ : TRUE
                         $rmax \triangleq Max[T \setminus \{n\}]$
                IN   IF $n \geq rmax$ THEN $n$ ELSE $rmax$
IN   $Max[S]$

CONSTANT $RM$,           The set of resource managers.
         $Acceptor$,      The set of acceptors.
         $Majority$,      The set of majorities of acceptors
         $Ballot$        The set of ballot numbers

ASSUME    We assume these properties of the declared constants.
  $\wedge\, Ballot \subseteq Nat$
  $\wedge\, 0 \in Ballot$
  $\wedge\, Majority \subseteq \text{SUBSET } Acceptor$
  $\wedge\, \forall\, MS1,\, MS2 \in Majority : MS1 \cap MS2 \neq \{\}$

All we assume about the set $Majority$ of majorities is that any two majorities have non-empty intersection.

$Message \triangleq$

The set of all possible messages. There are messages of type $\text{“Commit”}$ and $\text{“Abort”}$ to announce the decision, as well as messages for each phase of each instance of $ins$ of the *Paxos* consensus algorithm. The $acc$ field indicates the sender of a message from an acceptor to the leader; messages from a leader are broadcast to all acceptors.

$[type : \{\text{“phase1a”}\},\, ins : RM,\, bal : Ballot \setminus \{0\}]$
    $\cup$

$[type : \{\text{"phase1b"}\}, ins : RM, mbal : Ballot, bal : Ballot \cup \{-1\},$
$\quad val : \{\text{"prepared"}, \text{"aborted"}, \text{"none"}\}, acc : Acceptor]$
$\qquad \cup$
$[type : \{\text{"phase2a"}\}, ins : RM, bal : Ballot, val : \{\text{"prepared"}, \text{"aborted"}\}]$
$\qquad \cup$
$[type : \{\text{"phase2b"}\}, acc : Acceptor, ins : RM, bal : Ballot,$
$\quad val : \{\text{"prepared"}, \text{"aborted"}\}]$
$\qquad \cup$
$[type : \{\text{"Commit"}, \text{"Abort"}\}]$

---

VARIABLES

| | |
|---|---|
| $rmState,$ | $rmState[rm]$ is the state of resource manager $rm$. |
| $aState,$ | $aState[ins][ac]$ is the state of acceptor $ac$ for instance $ins$ of the *Paxos* algorithm |
| $msgs$ | The set of all messages ever sent. |

$PCTypeOK \triangleq$

The type-correctness invariant. Each acceptor maintains the values $mbal$, $bal$, and $val$ for each instance of the *Paxos* consensus algorithm.

$\quad \wedge rmState \in [RM \to \{\text{"working"}, \text{"prepared"}, \text{"committed"}, \text{"aborted"}\}]$
$\quad \wedge aState \quad \in [RM \to [Acceptor \to [mbal : Ballot,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad bal \quad : Ballot \cup \{-1\},$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad val \quad : \{\text{"prepared"}, \text{"aborted"}, \text{"none"}\}]]]$
$\quad \wedge msgs \in \text{SUBSET } Message$

$PCInit \triangleq$     The initial predicate.
$\quad \wedge rmState = [rm \in RM \mapsto \text{"working"}]$
$\quad \wedge aState \quad = [ins \in RM \mapsto$
$\qquad\qquad\qquad\qquad [ac \in Acceptor$
$\qquad\qquad\qquad\qquad\qquad \mapsto [mbal \mapsto 0, bal \mapsto -1, val \mapsto \text{"none"}]]]$
$\quad \wedge msgs = \{\}$

---

$\hfill$ **The Actions** $\hfill$

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

An action expression that describes the sending of message $m$.

---

**RM Actions**

$RMPrepare(rm) \triangleq$

Resource manager $rm$ prepares by sending a phase $2a$ message for ballot number 0 with value $\text{"prepared"}.$

$\quad \wedge rmState[rm] = \text{"working"}$
$\quad \wedge rmState' = [rmState \text{ EXCEPT } ![rm] = \text{"prepared"}]$
$\quad \wedge Send([type \mapsto \text{"phase2a"}, ins \mapsto rm, bal \mapsto 0, val \mapsto \text{"prepared"}])$
$\quad \wedge \text{UNCHANGED } aState$

$RMChooseToAbort(rm) \triangleq$

Resource manager $rm$ spontaneously decides to abort. It may (but need not) send a phase $2a$ message for ballot number 0 with value $\text{"aborted"}$.

$\land rmState[rm] = \text{"working"}$
$\land rmState' = [rmState \text{ EXCEPT } ![rm] = \text{"aborted"}]$
$\land Send([type \mapsto \text{"phase2a"}, ins \mapsto rm, bal \mapsto 0, val \mapsto \text{"aborted"}])$
$\land \text{UNCHANGED } aState$

$RMRcvCommitMsg(rm) \triangleq$

Resource manager $rm$ is told by the leader to commit. When this action is enabled, $rmState[rm]$ must equal either $\text{"prepared"}$ or $\text{"committed"}$. In the latter case, the action leaves the state unchanged (it is a "stuttering step").

$\land [type \mapsto \text{"Commit"}] \in msgs$
$\land rmState' = [rmState \text{ EXCEPT } ![rm] = \text{"committed"}]$
$\land \text{UNCHANGED } \langle aState, msgs \rangle$

$RMRcvAbortMsg(rm) \triangleq$

Resource manager $rm$ is told by the leader to abort. It could be in any state except $\text{"committed"}$.

$\land [type \mapsto \text{"Abort"}] \in msgs$
$\land rmState' = [rmState \text{ EXCEPT } ![rm] = \text{"aborted"}]$
$\land \text{UNCHANGED } \langle aState, msgs \rangle$

---

{\large \textbf{Leader Actions}} \vspace{.25\baselineskip}

The following actions are performed by any process that believes itself to be the current leader. Since leader selection is not assumed to be reliable, multiple processes could simultaneously consider themselves to be the leader.

$Phase1a(bal, rm) \triangleq$

If the leader times out without learning that a decision has been reached on resource manager $rm$'s prepare/abort decision, it can perform this action to initiate a new ballot $bal$. (Sending duplicate phase $1a$ messages is harmless.)

$\land Send([type \mapsto \text{"phase1a"}, ins \mapsto rm, bal \mapsto bal])$
$\land \text{UNCHANGED } \langle rmState, aState \rangle$

$Phase2a(bal, rm) \triangleq$

The action in which a leader sends a phase $2a$ message with ballot $bal > 0$ in instance $rm$, if it has received phase $1b$ messages for ballot number $bal$ from a majority of acceptors. If the leader received a phase $1b$ message from some acceptor that had sent a phase $2b$ message for this instance, then $maxbal \geq 0$ and the value $val$ the leader sends is determined by the phase $1b$ messages. (If
$val = \text{"prepared"}$, then $rm$ must have prepared.) Otherwise,
$maxbal = -1$ and the leader sends the value $\text{"aborted"}$. \vspace{.5\baselineskip}

The first conjunct asserts that the action is disabled if any commit leader has already sent a phase $2a$ message with ballot number $bal$. In practice, this is implemented by having ballot numbers partitioned among potential leaders, and having a leader record in stable storage the largest ballot number for which it sent a phase $2a$ message.

$\land \neg \exists m \in msgs : \land m.type = \text{"phase2a"}$

$$\wedge\ m.bal\ =\ bal$$
$$\wedge\ m.ins\ =\ rm$$
$$\wedge\ \exists\, MS \in Majority:$$
$$\quad \text{LET}\ mset\ \triangleq\ \{m \in msgs:\ \wedge\ m.type\ =\ \text{``phase1b''}$$
$$\wedge\ m.ins\ \ \ =\ rm$$
$$\wedge\ m.mbal\ =\ bal$$
$$\wedge\ m.acc\ \ \ \in\ MS\}$$
$$maxbal\ \triangleq\ Maximum(\{m.bal:\ m \in mset\})$$
$$val\ \triangleq\ \text{IF}\ maxbal\ =\ -1$$
$$\qquad\quad \text{THEN}\ \text{``aborted''}$$
$$\qquad\quad \text{ELSE}\ \ (\text{CHOOSE}\ m \in mset:\ m.bal = maxbal).val$$
$$\quad \text{IN}\quad \wedge\ \forall\, ac \in MS:\ \exists\, m \in mset:\ m.acc = ac$$
$$\qquad\quad \wedge\ Send([type \mapsto \text{``phase2a''},\ ins \mapsto rm,\ bal \mapsto bal,\ val \mapsto val])$$
$$\wedge\ \text{UNCHANGED}\ \langle rmState,\ aState \rangle$$

$Decide\ \triangleq$

A leader can decide that *Paxos Commit* has reached a result and send a message announcing the result if it has received the necessary phase $2b$ messages.

$$\wedge\ \text{LET}\ Decided(rm,\ v)\ \triangleq$$

True iff instance \$rm\$ of the *Paxos* consensus algorithm has chosen the value \$v\$.

$$\exists\, b \in Ballot,\ MS \in Majority:$$
$$\forall\, ac \in MS:\ [type \mapsto \text{``phase2b''},\ ins \mapsto rm,$$
$$bal \mapsto b,\ val \mapsto v,\ acc \mapsto ac] \in msgs$$
$$\text{IN}\quad \vee\ \wedge\ \forall\, rm \in RM:\ Decided(rm,\ \text{``prepared''})$$
$$\qquad\quad \wedge\ Send([type \mapsto \text{``Commit''}])$$
$$\vee\ \wedge\ \exists\, rm \in RM:\ Decided(rm,\ \text{``aborted''})$$
$$\qquad\quad \wedge\ Send([type \mapsto \text{``Abort''}])$$
$$\wedge\ \text{UNCHANGED}\ \langle rmState,\ aState \rangle$$

---

{\large \textbf{Acceptor Actions}}

$Phase1b(acc)\ \triangleq$
$$\exists\, m \in msgs:$$
$$\wedge\ m.type = \text{``phase1a''}$$
$$\wedge\ aState[m.ins][acc].mbal < m.bal$$
$$\wedge\ aState' = [aState\ \text{EXCEPT}\ ![m.ins][acc].mbal = m.bal]$$
$$\wedge\ Send([type\ \mapsto\ \text{``phase1b''},$$
$$ins\ \ \mapsto\ m.ins,$$
$$mbal \mapsto\ m.bal,$$
$$bal\ \ \mapsto\ aState[m.ins][acc].bal,$$
$$val\ \ \mapsto\ aState[m.ins][acc].val,$$
$$acc\ \ \mapsto\ acc])$$
$$\wedge\ \text{UNCHANGED}\ rmState$$

$Phase2b(acc)\ \triangleq$
$$\wedge\ \exists\, m \in msgs:$$

4

$$\land m.type = \text{``phase2a''}$$
$$\land aState[m.ins][acc].mbal \leq m.bal$$
$$\land aState' = [aState \text{ EXCEPT } ![m.ins][acc].mbal = m.bal,$$
$$![m.ins][acc].bal \quad = m.bal,$$
$$![m.ins][acc].val \quad = m.val]$$
$$\land Send([type \mapsto \text{``phase2b''}, \ ins \mapsto m.ins, \ bal \mapsto m.bal,$$
$$val \mapsto m.val, \ acc \mapsto acc])$$
$$\land \text{UNCHANGED } rmState$$

---

$PCNext \ \triangleq$      The next-state action
$$\lor \exists \, rm \in RM : \lor RMPrepare(rm)$$
$$\lor RMChooseToAbort(rm)$$
$$\lor RMRcvCommitMsg(rm)$$
$$\lor RMRcvAbortMsg(rm)$$
$$\lor \exists \, bal \in Ballot \setminus \{0\}, \ rm \in RM : Phase1a(bal, rm) \lor Phase2a(bal, rm)$$
$$\lor Decide$$
$$\lor \exists \, acc \in Acceptor : Phase1b(acc) \lor Phase2b(acc)$$

---

$PCSpec \ \triangleq \ PCInit \land \Box[PCNext]_{\langle rmState, \, aState, \, msgs \rangle}$

     The complete spec of the *Paxos Commit* protocol.

THEOREM $PCSpec \Rightarrow PCTypeOK$

---

We now assert that the two-phase commit protocol implements the transaction commit protocol of module *TCommit*. The following statement defines $TC ! TCSpec$ to be the formula $TCSpec$ of module $TCommit$. (The TLA$^+$ \textsc{instance} statement must is used to rename the operators defined in module $TCommit$ to avoid possible name conflicts with operators in the current module having the same name.)

$TC \ \triangleq \ \text{INSTANCE } TCommit$

THEOREM $PCSpec \Rightarrow TC ! TCSpec$