─────────────────── MODULE *ZeusReliableCommit* ───────────────────

Specification of *Zeus*'s reliable commit protocol presented in the *Zeus* paper
that appears in *Eurosys′*21.
This module includes everything but the pipelining optimization presented in the paper.

Model check passed [@ 21*st* of *Jan* 2021] with the following parameters:
$R\_NODES = \{0, 1, 2\}$
$R\_MAX\_EPOCH = 4$
$R\_MAX\_VERSION = 4$

EXTENDS *Integers*

CONSTANTS $R\_NODES$,
$\qquad\qquad R\_MAX\_EPOCH$,
$\qquad\qquad R\_MAX\_VERSION$

VARIABLES $rMsgs$,
$\qquad\qquad rKeyState$,
$\qquad\qquad rKeySharers$,
$\qquad\qquad rKeyVersion$,
$\qquad\qquad rKeyRcvedACKs$,
$\qquad\qquad rKeyLastWriter$,
$\qquad\qquad rNodeEpochID$,
$\qquad\qquad rAliveNodes$,
$\qquad\qquad rEpochID$

$vars \triangleq \langle rMsgs, rKeyState, rKeySharers, rKeyVersion, rKeyRcvedACKs,$
$\qquad\qquad rKeyLastWriter, rNodeEpochID, rAliveNodes, rEpochID \rangle$

─────────────────────────────────────────────────────────────────

The consistent invariant: all alive nodes in valid state should have the same value / *TS*
$RConsistentInvariant \triangleq$
$\quad \forall\, k,\, s \in rAliveNodes : \quad \lor rKeyState[k] \neq$ "valid"
$\qquad\qquad\qquad\qquad\qquad\qquad \lor rKeyState[s] \neq$ "valid"
$\qquad\qquad\qquad\qquad\qquad\qquad \lor rKeyVersion[k] = rKeyVersion[s]$

$RMaxVersionDistanceInvariant \triangleq$ this does not hold w/ the pipelining optimization
$\quad \forall\, k,\, s \in rAliveNodes :$
$\qquad\qquad\qquad\qquad\qquad \lor rKeyVersion[k] \leq rKeyVersion[s] + 1$
$\qquad\qquad\qquad\qquad\qquad \lor rKeyVersion[s] \leq rKeyVersion[k] + 1$

$RSingleOnwerInvariant \triangleq$
$\quad \forall\, k,\, s \in rAliveNodes :$
$\qquad\qquad\qquad\qquad\qquad \lor rKeySharers[k] \neq$ "owner"
$\qquad\qquad\qquad\qquad\qquad \lor rKeySharers[s] \neq$ "owner"
$\qquad\qquad\qquad\qquad\qquad \lor k = s$

$ROnwerOnlyWriterInvariant \triangleq$

$\forall\, k \in rAliveNodes :$
$$\lor\ rKeyState[k] \neq \text{``write''}$$
$$\lor\ rKeySharers[k] = \text{``owner''}$$

$ROnwerHighestVersionInvariant\ \triangleq$ owner has the highest version among alive nodes
$\quad \forall\, k,\, s \in rAliveNodes :$
$$\lor\ \land\ rKeySharers[s] \neq \text{``owner''}$$
$$\qquad \land\ rKeySharers[k] \neq \text{``owner''}$$
$$\lor$$
$$\qquad \land\ rKeySharers[k] = \text{``owner''}$$
$$\qquad \land\ rKeyVersion[k] \geq rKeyVersion[s]$$
$$\lor$$
$$\qquad \land\ rKeySharers[s] = \text{``owner''}$$
$$\qquad \land\ rKeyVersion[s] \geq rKeyVersion[k]$$

---

$RMessage\ \triangleq$ Messages exchanged by the Protocol
$\quad [type : \{\,\text{``INV''},\ \text{``ACK''}\,\},\ sender \quad\ : R\_NODES,$
$\qquad\qquad\qquad\qquad\qquad epochID \quad : 0\,..\,R\_MAX\_EPOCH,$
$\qquad\qquad\qquad\qquad\qquad version \quad\ : 0\,..\,R\_MAX\_VERSION]$
$\qquad\ \cup$
$\quad [type : \{\,\text{``VAL''}\,\}, \qquad\quad epochID \quad : 0\,..\,R\_MAX\_EPOCH,$
$\qquad\qquad\qquad\qquad\qquad version \quad\ : 0\,..\,R\_MAX\_VERSION]$

$RTypeOK\ \triangleq$ The type correctness invariant
$\quad \land \quad rMsgs \qquad\qquad \subseteq RMessage$
$\quad \land \quad rAliveNodes \qquad\ \subseteq R\_NODES$
$\quad \land \quad \forall\, n \in R\_NODES : rKeyRcvedACKs[n] \subseteq (R\_NODES \setminus \{n\})$
$\quad \land \quad rNodeEpochID \quad\ \in [R\_NODES \to 0\,..\,R\_MAX\_EPOCH]$
$\quad \land \quad rKeyLastWriter \quad \in [R\_NODES \to R\_NODES]$
$\quad \land \quad rKeyVersion \qquad \in [R\_NODES \to 0\,..\,R\_MAX\_VERSION]$
$\quad \land \quad rKeySharers \qquad \in [R\_NODES \to \{\,\text{``owner''},\ \text{``reader''},\ \text{``non-sharer''}\,\}]$
$\quad \land \quad rKeyState \qquad\quad \in [R\_NODES \to \{\,\text{``valid''},\ \text{``invalid''},\ \text{``write''},\ \text{``replay''}\,\}]$

$RInit\ \triangleq$ The initial predicate
$\quad \land\ rMsgs \qquad\qquad = \{\}$
$\quad \land\ rEpochID \qquad\quad = 0$
$\quad \land\ rAliveNodes \qquad = R\_NODES$
$\quad \land\ rKeyVersion \qquad = [n \in R\_NODES \mapsto 0]$
$\quad \land\ rNodeEpochID \quad\ = [n \in R\_NODES \mapsto 0]$
$\quad \land\ rKeyRcvedACKs\ = [n \in R\_NODES \mapsto \{\}]$
$\quad \land\ rKeySharers \qquad = [n \in R\_NODES \mapsto \text{``reader''}]$
$\quad \land\ rKeyState \qquad\quad = [n \in R\_NODES \mapsto \text{``valid''}]$
$\quad \land\ rKeyLastWriter \quad = [n \in R\_NODES \mapsto \textsc{choose}\ k \in R\_NODES :$

$$\forall\, m \in R\_NODES : k \leq m]$$

---

$RNoChanges\_in\_membership \;\triangleq\; \text{UNCHANGED} \;\langle rAliveNodes,\, rEpochID \rangle$

$RNoChanges\_but\_membership \;\triangleq$
$\quad$ UNCHANGED $\langle rMsgs,\, rKeyState,\, rKeyVersion,$
$\qquad\qquad\qquad rKeyRcvedACKs,\, rKeyLastWriter,$
$\qquad\qquad\qquad rKeySharers,\, rNodeEpochID \rangle$

$RNoChanges \;\triangleq$
$\quad \wedge\; RNoChanges\_in\_membership$
$\quad \wedge\; RNoChanges\_but\_membership$

---

A buffer maintaining all network messages. Messages are only appended to
this variable (not \* removed once delivered) intentionally to check
protocol's tolerance in dublicates and reorderings
$RSend(m) \;\triangleq\; rMsgs' = rMsgs \cup \{m\}$

Check if all acknowledgments for a write have been received
$RAllACKsRcved(n) \;\triangleq\; (rAliveNodes \setminus \{n\}) \subseteq rKeyRcvedACKs[n]$

$RIsAlive(n) \;\triangleq\; n \in rAliveNodes$

$RNodeFailure(n) \;\triangleq\;$ Emulate a node failure
$\quad$ Make sure that there are atleast 3 alive nodes before killing a node
$\quad \wedge\, \exists\, k,\, m \in rAliveNodes : \;\wedge\; k \;\neq n$
$\qquad\qquad\qquad\qquad\qquad\quad\; \wedge\; m \neq n$
$\qquad\qquad\qquad\qquad\qquad\quad\; \wedge\; m \neq k$
$\quad \wedge\, rEpochID' = rEpochID + 1$
$\quad \wedge\, rAliveNodes' = rAliveNodes \setminus \{n\}$
$\quad \wedge\, RNoChanges\_but\_membership$

---

$RNewOwner(n) \;\triangleq$
$\quad \wedge\, \forall\, k \in rAliveNodes :$
$\qquad \wedge\, rKeySharers[k] \qquad\quad \neq \text{``owner''}$
$\qquad \wedge\, \vee\; \wedge\; rKeyState[k] \quad\; = \text{``valid''} \qquad$ all alive replicas are in valid state
$\qquad\qquad\quad\; \wedge\; rKeySharers[k] \;= \text{``reader''} \qquad$ and there is not alive owner
$\qquad\qquad \vee\; \wedge\; rKeySharers[k] \;= \text{``non-sharer''} \quad$ and there is not alive owner
$\quad \wedge\, rKeySharers' \qquad\qquad\; = [rKeySharers \quad\; \text{EXCEPT} \;\,![n] = \text{``owner''}]$
$\quad \wedge\, \text{UNCHANGED} \;\langle rMsgs,\, rKeyState,\, rKeyVersion,\, rKeyRcvedACKs,$
$\qquad\qquad\qquad\qquad rKeyLastWriter,\, rAliveNodes,\, rNodeEpochID,\, rEpochID \rangle$

$ROverthrowOwner(n) \;\triangleq$
$\quad \exists\, k \in rAliveNodes :$

$\quad\quad\quad \land\ rKeyState[k] \quad = \text{``valid''}$
$\quad\quad\quad \land\ rKeySharers[k] = \text{``owner''}$
$\quad\quad\quad \land\ rKeySharers' \quad = [rKeySharers \text{ EXCEPT } ![n] = \text{``owner''},$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad ![k] = \text{``reader''}]$
$\quad\quad\quad \land\ \text{UNCHANGED } \langle rMsgs,\ rKeyState,\ rKeyVersion,\ rKeyRcvedACKs,$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad rKeyLastWriter,\ rAliveNodes,\ rNodeEpochID,\ rEpochID\rangle$

$RGetOwnership(n) \triangleq$
$\quad\quad \land\ rKeySharers[n] \neq \text{``owner''}$
$\quad\quad \land\ \forall\, x \in rAliveNodes : rNodeEpochID[x] = rEpochID \quad \boxed{TODO \text{ may move this to } RNewOwner}$
$\quad\quad \land\ \lor\ ROverthrowOwner(n)$
$\quad\quad\quad\ \lor\ RNewOwner(n)$

---

$RRead(n) \triangleq \quad \boxed{\text{Execute a read}}$
$\quad\quad \land\ rNodeEpochID[n] = rEpochID$
$\quad\quad \land\ rKeyState[n] \quad\quad = \text{``valid''}$
$\quad\quad \land\ RNoChanges$

$RRcvInv(n) \triangleq \quad \boxed{\text{Process a received invalidation}}$
$\ \exists\, m \in rMsgs :$
$\quad\quad \land\ m.type \quad\quad\ = \text{``INV''}$
$\quad\quad \land\ m.epochID \ = rEpochID$
$\quad\quad \land\ m.sender \ \neq n$
$\quad\quad \land\ m.sender \ \in rAliveNodes$
$\quad\quad \boxed{\text{always acknowledge a received invalidation (irrelevant to the timestamp)}}$
$\quad\quad \land\ RSend([type \quad\quad \mapsto \text{``ACK''},$
$\quad\quad\quad\quad\quad\quad epochID \quad\ \mapsto rEpochID,$
$\quad\quad\quad\quad\quad\quad sender \quad\quad \mapsto n,$
$\quad\quad\quad\quad\quad\quad version \quad\ \mapsto m.version])$
$\quad\quad \land\ \lor\ m.version \quad\quad\quad > rKeyVersion[n]$
$\quad\quad\quad\quad \land\ rKeyState[n] \quad \in \{\text{``valid''}, \text{``invalid''}, \text{``replay''}\}$
$\quad\quad\quad\quad \land\ rKeyState' \quad\quad = [rKeyState \text{ EXCEPT } ![n] = \text{``invalid''}]$
$\quad\quad\quad\quad \land\ rKeyVersion' \quad = [rKeyVersion \text{ EXCEPT } ![n] \ = m.version]$
$\quad\quad\quad\quad \land\ rKeyLastWriter' = [rKeyLastWriter \text{ EXCEPT } ![n] = m.sender]$
$\quad\quad\quad \lor\ m.version \quad\quad\quad \leq rKeyVersion[n]$
$\quad\quad\quad\quad \land\ \text{UNCHANGED } \langle rKeyState,\ rKeyVersion,\ rKeyLastWriter\rangle$
$\quad\quad \land\ \text{UNCHANGED } \langle rAliveNodes,\ rKeySharers,\ rKeyRcvedACKs,\ rNodeEpochID,\ rEpochID\rangle$

$RRcvVal(n) \triangleq \quad \boxed{\text{Process a received validation}}$
$\quad\ \exists\, m \in rMsgs :$
$\quad\quad \land\ rKeyState[n] \neq \text{``valid''}$
$\quad\quad \land\ m.type \quad\quad\ = \text{``VAL''}$
$\quad\quad \land\ m.epochID \quad = rEpochID$
$\quad\quad \land\ m.version \quad\ = rKeyVersion[n]$
$\quad\quad \land\ rKeyState' \quad = [rKeyState \text{ EXCEPT } ![n] = \text{``valid''}]$

4

$\land$ UNCHANGED $\langle rMsgs,\ rKeyVersion,\ rKeyLastWriter,\ rKeySharers,$
$\qquad\qquad\qquad rAliveNodes,\ rKeyRcvedACKs,\ rNodeEpochID,\ rEpochID\rangle$

$RReaderActions(n) \triangleq$     Actions of a write follower
     $\lor\ RRead(n)$
     $\lor\ RRcvInv(n)$
     $\lor\ RRcvVal(n)$

---

$RWrite(n) \triangleq$
     $\land\ rNodeEpochID[n] \quad = \quad rEpochID$
     $\land\ rKeySharers[n] \quad\ \in \quad \{\text{"owner"}\}$
     $\land\ rKeyState[n] \qquad\ \in\ \{\text{"valid"}\}$ May add invalid state here as well
     $\land\ rKeyVersion[n] \quad\ < \quad R\_MAX\_VERSION$
     $\land\ rKeyLastWriter' \quad = \quad [rKeyLastWriter \quad \text{EXCEPT}\ ![n] \quad = n]$
     $\land\ rKeyRcvedACKs' \quad = \quad [rKeyRcvedACKs \quad\ \text{EXCEPT}\ ![n] = \{\}]$
     $\land\ rKeyState' \qquad\ = \quad [rKeyState \qquad\ \text{EXCEPT}\ ![n] \qquad = \text{"write"}]$
     $\land\ rKeyVersion' \qquad = \quad [rKeyVersion \qquad\ \text{EXCEPT}\ ![n] \qquad = rKeyVersion[n] + 1]$
     $\land\ RSend([type \qquad\quad \mapsto \text{"INV"},$
             $epochID \qquad\ \mapsto rEpochID,$
             $sender \qquad\quad \mapsto n,$
             $version \qquad\quad \mapsto rKeyVersion[n] + 1])$
     $\land$ UNCHANGED $\langle rAliveNodes,\ rKeySharers,\ rNodeEpochID,\ rEpochID\rangle$

$RRcvAck(n) \triangleq$     Process a received acknowledment
     $\exists\, m \in rMsgs :$
        $\land\ m.type \qquad\quad = \qquad \text{"ACK"}$
        $\land\ m.epochID \qquad = \qquad rEpochID$
        $\land\ m.sender \qquad\ \neq \qquad n$
        $\land\ m.version \qquad\ = \qquad rKeyVersion[n]$
        $\land\ m.sender \qquad\ \notin\ rKeyRcvedACKs[n]$
        $\land\ rKeyState[n]\ \in \{\text{"write"},\ \text{"replay"}\}$
        $\land\ rKeyRcvedACKs' = \quad [rKeyRcvedACKs \qquad \text{EXCEPT}\ ![n] =$
                                $rKeyRcvedACKs[n] \cup \{m.sender\}]$
        $\land$ UNCHANGED $\langle rMsgs,\ rKeyState,\ rKeyVersion,\ rKeyLastWriter,$
                  $rAliveNodes,\ rKeySharers,\ rNodeEpochID,\ rEpochID\rangle$

$RSendVals(n) \triangleq$    Send validations once received acknowledments from all alive nodes
     $\land\ rKeyState[n] \quad\ \in \quad \{\text{"write"},\ \text{"replay"}\}$
     $\land\ RAllACKsRcved(n)$
     $\land\ rKeyState' \qquad\ = \quad [rKeyState\ \text{EXCEPT}\ ![n] = \text{"valid"}]$
     $\land\ RSend([type \qquad\quad \mapsto \text{"VAL"},$
           $epochID \qquad\ \mapsto rEpochID,$
           $version \qquad\quad \mapsto rKeyVersion[n]])$
     $\land$ UNCHANGED $\langle rKeyRcvedACKs,\ rKeyVersion,\ rKeyLastWriter,$

5

$$\langle rAliveNodes,\ rKeySharers,\ rNodeEpochID,\ rEpochID \rangle$$

$ROwnerActions(n) \triangleq$ Actions of a read/write coordinator
 $\lor\ RRead(n)$
 $\lor\ RWrite(n)$
 $\lor\ RRcvAck(n)$
 $\lor\ RSendVals(n)$

---

$RWriteReplay(n) \triangleq$ Execute a write-replay
 $\land\ rKeyLastWriter' \quad = \quad [rKeyLastWriter \quad \text{EXCEPT} \ ![n] \quad = n]$
 $\land\ rKeyRcvedACKs' \quad = \quad [rKeyRcvedACKs \quad \text{EXCEPT} \ ![n] = \{\}]$
 $\land\ rKeyState' \qquad\quad = \quad [rKeyState \qquad\quad \text{EXCEPT} \ ![n] \quad = \text{"replay"}]$
 $\land\ RSend([type \qquad\quad \mapsto \text{"INV"},$
     $sender \qquad\quad \mapsto n,$
     $epochID \qquad\quad \mapsto rEpochID,$
     $version \qquad\quad \mapsto rKeyVersion[n]])$
 $\land\ \text{UNCHANGED}\ \langle rKeyVersion,\ rKeySharers,\ rAliveNodes,\ rNodeEpochID,\ rEpochID \rangle$

$RLocalWriteReplay(n) \triangleq$
 $\land\ \lor\ rKeySharers[n] = \text{"owner"}$
  $\lor\ rKeyState[n] \quad = \text{"replay"}$
 $\land\ RWriteReplay(n)$

$RFailedNodeWriteReplay(n) \triangleq$
 $\land\ \neg RIsAlive(rKeyLastWriter[n])$
 $\land\ rKeyState[n] \qquad = \text{"invalid"}$
 $\land\ RWriteReplay(n)$

$RUpdateLocalEpochID(n) \triangleq$
 $\land\ rKeyState[n] \qquad = \text{"valid"}$
 $\land\ rNodeEpochID' \quad = [rNodeEpochID \ \text{EXCEPT} \ ![n] = rEpochID]$
 $\land\ \text{UNCHANGED}\ \langle rMsgs,\ rKeyState,\ rKeyVersion,\ rKeyRcvedACKs,$
       $rKeyLastWriter,\ rKeySharers,\ rAliveNodes,\ rEpochID \rangle$

$RReplayActions(n) \triangleq$
 $\land\ rNodeEpochID[n] < rEpochID$
 $\land\ \lor\ RLocalWriteReplay(n)$
  $\lor\ RFailedNodeWriteReplay(n)$
  $\lor\ RUpdateLocalEpochID(n)$

---

$RNext \triangleq$ Modeling protocol (Owner and Reader actions while emulating failures)
 $\exists\, n \in rAliveNodes :$
   $\lor\ RReaderActions(n)$
   $\lor\ ROwnerActions(n)$

$\lor\ RReplayActions(n)$
$\lor\ RGetOwnership(n)$
$\lor\ RNodeFailure(n)$ <span style="background-color:#d3d3d3">emulate node failures</span>

<span style="background-color:#d3d3d3">The complete definition of the algorithm</span>

$Spec\ \triangleq\ RInit\ \land\ \Box[RNext]_{vars}$

$Invariants\ \triangleq\ \land\ (\Box RTypeOK)$
$\qquad\qquad\quad\ \land\ (\Box RConsistentInvariant)$
$\qquad\qquad\quad\ \land\ (\Box RSingleOnwerInvariant)$
$\qquad\qquad\quad\ \land\ (\Box ROnwerOnlyWriterInvariant)$
$\qquad\qquad\quad\ \land\ (\Box RMaxVersionDistanceInvariant)$
$\qquad\qquad\quad\ \land\ (\Box ROnwerHighestVersionInvariant)$

THEOREM $Spec\ \Rightarrow\ Invariants$