\ * To fix: heartbeat

─────────────────── MODULE $TunableMongoDB\_Repl$ ───────────────────

EXTENDS $Naturals,\ FiniteSets,\ Sequences,\ TLC$

constants and variables

CONSTANTS $Client,\ Server,$     the set of clients and servers
           $Key,\ Value,$      the set of keys and values
           $Nil,$         model value, place holder
           $PtStop$       max physical time

VARIABLES $Primary,$       Primary node
           $Secondary,$     secondary nodes
           $Oplog,$        $oplog[s]$: $oplog$ at $server[s]$
           $Store,$        $store[s]$: data stored at $server[s]$
           $Ct,$          $Ct[s]$: cluster time at node $s$
           $Ot,$          $Ot[s]$: the last applied operation time at server $s$
           $ServerMsg,$    $ServerMsg[s]$: the channel of heartbeat $msgs$ at server $s$
           $Pt,$          $Pt[s]$: physical time at server $s$
           $Cp,$          $Cp[s]$: majority commit point at server $s$
           $State,$        $State[s]$: the latest $Ot$ of all servers that server $s$ knows
           $CurrentTerm,$    $CurrentTerm[s]$: current election term at server $s$
                              $\rightarrow$ updated in $update\_position$, heartbeat and replicate
           $ReadyToServe,$    equal to 0 before any primary is elected
           $SyncSource$     $SyncSource[s]$: sync source of server node $s$

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

ASSUME $Cardinality(Client) \geq 1$    at least one clinet
ASSUME $Cardinality(Server) \geq 2$    at least one primary and one secondary
ASSUME $Cardinality(Key) \geq 1$    at least one object
ASSUME $Cardinality(Value) \geq 2$    at least two values to update

Helpers

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

$HLCLt(x,\ y) \triangleq$ IF $x.p < y.p$
                 THEN TRUE
                ELSE IF $x.p = y.p$
                THEN IF $x.l < y.l$
                       THEN TRUE
                     ELSE FALSE
                ELSE FALSE

$HLCMin(x,\ y) \triangleq$ IF $HLCLt(x,\ y)$ THEN $x$ ELSE $y$
$HLCMax(x,\ y) \triangleq$ IF $HLCLt(x,\ y)$ THEN $y$ ELSE $x$
$HLCType \triangleq [p : Nat,\ l : Nat]$
$Min(x,\ y) \triangleq$ IF $x < y$ THEN $x$ ELSE $y$
$Max(x,\ y) \triangleq$ IF $x > y$ THEN $x$ ELSE $y$

1

$$vars \triangleq \langle Primary,\ Secondary,\ Oplog,\ Store,\ Ct,\ Ot,\ ServerMsg,$$
$$Pt,\ Cp,\ State,\ CurrentTerm,\ ReadyToServe,\ SyncSource\rangle$$

RECURSIVE $CreateState(\_,\ \_)$ init state
$CreateState(len,\ seq) \triangleq$
    IF $len = 0$ THEN $seq$
    ELSE  $CreateState(len - 1,\ Append(seq,\ [p \mapsto 0,\ l \mapsto 0]))$

$LogTerm(i,\ index) \triangleq$ IF $index = 0$ THEN $0$ ELSE  $Oplog[i][index].term$
$LastTerm(i) \triangleq CurrentTerm[i]$

Is node $i$ ahead of node $j$
$NotBehind(i,\ j) \triangleq \quad \vee\ LastTerm(i) > LastTerm(j)$
$$\vee\ \wedge\ LastTerm(i) = LastTerm(j)$$
$$\wedge\ Len(Oplog[i]) \geq Len(Oplog[j])$$

$IsMajority(servers) \triangleq Cardinality(servers) * 2 > Cardinality(Server)$

Return the maximum value from a set, or undefined if the set is empty.
$MaxVal(s) \triangleq$ CHOOSE $x \in s : \forall\, y \in s : x \geq y$
$HLCMinSet(s) \triangleq$ CHOOSE $x \in s : \forall\, y \in s : \neg HLCLt(y,\ x)$

clock

$MaxPt \triangleq$ LET $x \triangleq$ CHOOSE $s \in Server : \forall\, s1 \in Server \setminus \{s\} :$
$$Pt[s] \geq Pt[s1]$$
        IN    $Pt[x]$

$Tick(s) \triangleq Ct' =$ IF $Ct[s].p \geq Pt[s]$
                THEN $[Ct$ EXCEPT $![s] = [p \mapsto @.p,\ l \mapsto @.l + 1]]$
                ELSE  $[Ct$ EXCEPT $![s] = [p \mapsto Pt[s],\ l \mapsto 0]]$

heartbeat
Only $Primary$ node sends heartbeat once advance pt
$BroadcastHeartbeat(s) \triangleq$
    LET $msg \triangleq [type \mapsto$ "heartbeat", $s \mapsto s,\ aot \mapsto Ot[s],$
            $ct \mapsto Ct[s],\ cp \quad \mapsto Cp[s],\ term \mapsto CurrentTerm[s]]$
    IN    $ServerMsg' = [x \in Server \mapsto$ IF $x = s$ THEN $ServerMsg[x]$
                            ELSE  $Append(ServerMsg[x],\ msg)]$

Can node $i$ sync from node $j$?
$CanSyncFrom(i,\ j) \triangleq$
    $\wedge\ Len(Oplog[i]) < Len(Oplog[j])$
    $\wedge\ LastTerm(i) = LogTerm(j,\ Len(Oplog[i]))$

$Oplog$ entries needed to replicate from $j$ to $i$
$ReplicateOplog(i,\ j) \triangleq$
    LET $len\_i \triangleq Len(Oplog[i])$

$$len\_j \triangleq Len(Oplog[j])$$
IN    IF $i \in Secondary \land len\_i < len\_j$
                THEN $SubSeq(Oplog[j], len\_i + 1, len\_j)$
                ELSE $\langle\rangle$

Can node $i$ rollback its $log$ based on $j$'s $log$
$CanRollback(i, j) \triangleq \land Len(Oplog[i]) > 0$
                $\land Len(Oplog[j]) > 0$
                $\land CurrentTerm[i] < CurrentTerm[j]$
                $\land$
                    $\lor Len(Oplog[i]) > Len(Oplog[j])$
                    $\lor \land Len(Oplog[i]) \leq Len(Oplog[j])$
                        $\land CurrentTerm[i] \neq LogTerm(j, Len(Oplog[i]))$

Returns the highest common index between two divergent logs.
If there is no common index between the logs, returns 0.
$RollbackCommonPoint(i, j) \triangleq$
    LET $commonIndices \triangleq \{k \in \text{DOMAIN } Oplog[i] :$
                        $\land k \leq Len(Oplog[j])$
                        $\land Oplog[i][k] = Oplog[j][k]\}$IN
        IF $commonIndices = \{\}$ THEN 0 ELSE $MaxVal(commonIndices)$

The set of all $quorums$. This just calculates simple majorities, but the only
important property is that every quorum overlaps with every other.
$Quorum \triangleq \{i \in \text{SUBSET } (Server) : Cardinality(i) * 2 > Cardinality(Server)\}$

$QuorumAgreeInSameTerm(states) \triangleq$
    LET $quorums \triangleq \{Q \in Quorum :$
                    Make sure all nodes in quorum have actually applied some entries.
                    $\land \lor \forall s \in Q : states[s].p > 0$
                        $\lor \land \forall s \in Q : states[s].p = 0$
                            $\land \forall s \in Q : states[s].l > 0$
                    Make sure every applied entry in quorum has the same term.
                    $\land \forall s, t \in Q :$
                        $s \neq t \Rightarrow states[s].term = states[s].term\}$
    IN    $quorums$

*Init* Part

---

$InitPrimary \triangleq Primary = \{\text{CHOOSE } s \in Server : \text{TRUE}\}$
$InitSecondary \triangleq Secondary = Server \setminus Primary$
$InitOplog \triangleq Oplog = [s \in Server \mapsto \langle\rangle]$
$InitStore \triangleq Store = [n \in Server \cup Client \mapsto [k \in Key \mapsto Nil]]$
$InitCt \triangleq Ct = [n \in Server \cup Client \mapsto [p \mapsto 0, l \mapsto 0]]$
$InitOt \triangleq Ot = [n \in Server \cup Client \mapsto [p \mapsto 0, l \mapsto 0]]$
$InitServerMsg \triangleq ServerMsg = [s \in Server \mapsto \langle\rangle]$

$InitPt \;\triangleq\; Pt = [s \in Server \mapsto 1]$
$InitCp \;\triangleq\; Cp = [n \in Server \cup Client \mapsto [p \mapsto 0,\; l \mapsto 0]]$
$InitState \;\triangleq\; State = [s \in Server \mapsto [s0 \in Server \mapsto$
$$[p \mapsto 0,\; l \mapsto 0,\; term \mapsto 0]]]$$
$InitCurrentTerm \;\triangleq\; CurrentTerm = [s \in Server \mapsto 0]$
$InitReadyToServe \;\triangleq\; ReadyToServe = 0$
$InitSyncSource \;\triangleq\; SyncSource = [s \in Server \mapsto Nil]$

$Init \;\triangleq\;$
 $\land\; InitPrimary \land InitSecondary \land InitOplog \land InitStore \land InitCt$
 $\land\; InitOt \land InitPt \land InitCp$
 $\land\; InitServerMsg$
 $\land\; InitState \land InitCurrentTerm \land InitReadyToServe$
 $\land\; InitSyncSource$

Next *State* Actions
Replication Protocol: possible actions

⊢─────────────────────────────────────────────────────────────────────────┤

$TurnOnReadyToServe \;\triangleq\;$
 $\land\; ReadyToServe = 0$
 $\land\; \exists\, s \in Primary :$
  $\land\; CurrentTerm' = [CurrentTerm \text{ EXCEPT } ![s] = CurrentTerm[s] + 1]$
  $\land\; ReadyToServe' = ReadyToServe + 1$
 $\land\; \text{UNCHANGED } \langle Primary,\;\; Secondary,\; Oplog,\; Store,\; Ct,\; Ot,$
      $ServerMsg,\; Pt,\; Cp,\; State,\; SyncSource\rangle$

$Stepdown \;\triangleq\;$
   $\land\; ReadyToServe > 0$
   $\land\; \exists\, s \in Primary :$
    $\land\; Primary' = Primary \setminus \{s\}$
    $\land\; Secondary' = Secondary \cup \{s\}$
   $\land\; \text{UNCHANGED } \langle Oplog,\; Store,\; Ct,\; Ot,\; ServerMsg,$
       $Pt,\; Cp,\; State,\; CurrentTerm,$
       $ReadyToServe,\; SyncSource\rangle$

Todo: *Stepdown* when receiving a higher term heartbeat

There are majority nodes agree to elect node $i$ to become primary
$ElectPrimary \;\triangleq\;$
 $\land\; ReadyToServe > 0$
 $\land\; \exists\, i \in Server : \exists\, majorNodes \in \text{SUBSET } (Server) :$
  $\land\; \forall\, j \in majorNodes : \;\land\; NotBehind(i, j)$
        $\land\; CurrentTerm[i] \geq CurrentTerm[j]$
  $\land\; IsMajority(majorNodes)$
  voted nodes for $i$ cannot be primary anymore
  $\land\; Primary' = \text{LET } possiblePrimary \;\triangleq\; Primary \setminus majorNodes$

4

$$\text{IN} \quad possiblePrimary \cup \{i\}$$

add voted nodes into secondaries
$$\land\ Secondary' = \text{LET}\ possibleSecondary\ \triangleq\ Secondary \cup majorNodes$$
$$\text{IN} \quad possibleSecondary \setminus \{i\}$$
$$\land\ CurrentTerm' = [index \in Server \mapsto \text{IF}\ index \in (majorNodes \cup \{i\})$$
$$\text{THEN}\ CurrentTerm[i] + 1$$
$$\text{ELSE}\ CurrentTerm[index]]$$

A primary node do not have any sync source
$$\land\ SyncSource' = [SyncSource\ \text{EXCEPT}\ ![i] = Nil]$$
$$\land\ \text{UNCHANGED}\ \langle Oplog,\ Store,\ Ct,\ Ot,\ ServerMsg,\ Pt,\ Cp,\ State,\ ReadyToServe \rangle$$

$$AdvanceCp\ \triangleq$$
$$\land\ ReadyToServe > 0$$
$$\land\ \exists\, s \in Primary:$$
$$\text{LET}\ newCp\ \triangleq$$
$$\text{LET}\ quorumAgree\ \triangleq\ QuorumAgreeInSameTerm(State[s])$$
$$\text{IN} \quad \text{IF}\ Cardinality(quorumAgree) > 0$$
$$\text{THEN LET}\ QuorumSet\ \triangleq\ \text{CHOOSE}\ i \in quorumAgree : \text{TRUE}$$
$$serverInQuorum\ \triangleq\ \text{CHOOSE}\ j \in QuorumSet : \text{TRUE}$$
$$termOfQuorum\ \triangleq\ State[s][serverInQuorum].term$$
$$StateSet\ \triangleq\ \{[p \mapsto State[s][j].p,\ l \mapsto State[s][j].l] : j \in QuorumSet\}$$
$$newCommitPoint\ \triangleq\ HLCMinSet(StateSet)$$
$$oldCommitPoint\ \triangleq\ [p \mapsto Cp[s].p,\ l \mapsto Cp[s].l]$$

$newCp$ must be greater than current $Cp$ for primary to advance it
$$\text{IN} \quad \text{IF}\ termOfQuorum = CurrentTerm[s] \land HLCLt(oldCommitPoint, newCommit$$
$$\text{THEN}\ [p \mapsto newCommitPoint.p,\ l \mapsto newCommitPoint.l,\ term \mapsto termOf$$
$$\text{ELSE}\ Cp[s]$$
$$\text{ELSE}\ Cp[s]$$
$$\text{IN} \quad Cp' = [Cp\ \text{EXCEPT}\ ![s] = newCp]$$
$$\land\ \text{UNCHANGED}\ \langle Primary,\ Secondary,\ Oplog,\ Store,\ Ct,\ Ot,$$
$$ServerMsg,\ Pt,\ State,\ CurrentTerm,\ ReadyToServe,\ SyncSource \rangle$$

heartbeat oplog $Ot$ store
$$ServerTakeHeartbeat\ \triangleq$$
$$\land\ ReadyToServe > 0$$
$$\land\ \exists\, s \in Server:$$
$$\land\ Len(ServerMsg[s]) \neq 0 \quad \text{message channel is not empty}$$
$$\land\ ServerMsg[s][1].type = \text{``heartbeat''}$$
$$\land\ CurrentTerm[s] = ServerMsg[s][1].term$$
$$\land\ Ct' = [Ct\ \text{EXCEPT}\ ![s] = HLCMax(Ct[s],\ ServerMsg[s][1].ct)]$$
$$\land\ State' =$$
$$\text{LET}\ newState\ \triangleq\ [$$
$$p \mapsto ServerMsg[s][1].aot.p,$$
$$l \mapsto ServerMsg[s][1].aot.l,$$
$$term \mapsto ServerMsg[s][1].term$$

$$
\begin{aligned}
&\qquad\qquad ] \\
&\qquad\text{IN}\quad \text{LET}\ SubHbState \triangleq State[s] \\
&\qquad\qquad\qquad hb \triangleq [SubHbState\ \text{EXCEPT}\ ![ServerMsg[s][1].s] = newState] \\
&\qquad\qquad \text{IN}\quad [State\ \text{EXCEPT}\ ![s] = hb] \\
&\quad \wedge\ Cp' = \text{LET}\ newcp \triangleq
\end{aligned}
$$

primary node: compute new mcp

$$
\begin{aligned}
&\qquad\qquad \text{IF}\ s \in Primary\ \text{THEN} \\
&\qquad\qquad\qquad \text{LET}\ quorumAgree \triangleq QuorumAgreeInSameTerm(State[s])\text{IN} \\
&\qquad\qquad\qquad\quad \text{IF}\ Cardinality(quorumAgree) > 0 \\
&\qquad\qquad\qquad\qquad \text{THEN LET}\ QuorumSet \triangleq \text{CHOOSE}\ i \in quorumAgree : \text{TRUE} \\
&\qquad\qquad\qquad\qquad\qquad\qquad serverInQuorum \triangleq \text{CHOOSE}\ j \in QuorumSet : \text{TRUE} \\
&\qquad\qquad\qquad\qquad\qquad\qquad termOfQuorum \triangleq State[s][serverInQuorum].term \\
&\qquad\qquad\qquad\qquad\qquad\qquad StateSet \triangleq \{[p \mapsto State[s][j].p,\ l \mapsto State[s][j].l] : j \in QuorumSet \\
&\qquad\qquad\qquad\qquad\qquad\qquad newCommitPoint \triangleq HLCMinSet(StateSet) \\
&\qquad\qquad\qquad\qquad\qquad \text{IN}\quad \text{IF}\ termOfQuorum = CurrentTerm[s] \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{THEN} \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad [p \mapsto newCommitPoint.p,\ l \mapsto newCommitPoint.l,\ term \mapsto \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{ELSE}\quad Cp[s] \\
&\qquad\qquad\qquad\qquad \text{ELSE}\quad Cp[s]
\end{aligned}
$$

secondary node: update mcp

$$
\begin{aligned}
&\qquad\qquad\qquad \text{ELSE}\ \text{IF LET}\ msgCP \triangleq [p \mapsto ServerMsg[s][1].cp.p,\ l \mapsto ServerMsg[s][1].cp.l]\text{IN} \\
&\qquad\qquad\qquad\qquad\qquad \wedge \neg HLCLt(msgCP,\ Cp[s]) \\
&\qquad\qquad\qquad\qquad\qquad \wedge \neg HLCLt(Ot[s],\ msgCP)
\end{aligned}
$$

The term of *cp* must equal to the *CurrentTerm* of that node to advance it

$$
\begin{aligned}
&\qquad\qquad\qquad\qquad\qquad \wedge ServerMsg[s].cp.term = CurrentTerm[s] \\
&\qquad\qquad\qquad\qquad \text{THEN}\ ServerMsg[s][1].cp \\
&\qquad\qquad\qquad\qquad \text{ELSE}\quad Cp[s] \\
&\qquad\qquad \text{IN}\quad [Cp\ \text{EXCEPT}\ ![s] = newcp] \\
&\quad \wedge ServerMsg' = [ServerMsg\ \text{EXCEPT}\ ![s] = Tail(@)] \\
&\quad \wedge CurrentTerm' = [CurrentTerm\ \text{EXCEPT}\ ![s] = Max(CurrentTerm[s],\ ServerMsg[s][1].term)] \\
&\ \wedge \text{UNCHANGED}\ \langle Primary,\ Secondary,\ Oplog,\ Store,\ Ot,\ Pt, \\
&\qquad\qquad\qquad\qquad ReadyToServe,\ SyncSource\rangle
\end{aligned}
$$

$$
\begin{aligned}
&ServerTakeUpdatePosition \triangleq \\
&\quad \wedge ReadyToServe > 0 \\
&\quad \wedge \exists s \in Server : \\
&\qquad \wedge Len(ServerMsg[s]) \neq 0 \quad \text{message channel is not empty} \\
&\qquad \wedge ServerMsg[s][1].type = \text{``update\_position''} \\
&\qquad \wedge Ct' = [Ct\ \text{EXCEPT}\ ![s] = HLCMax(Ct[s],\ ServerMsg[s][1].ct)] \quad \text{update } ct \text{ accordingly} \\
&\qquad \wedge State' = \\
&\qquad \text{LET}\ newState \triangleq [ \\
&\qquad\qquad\qquad p \mapsto ServerMsg[s][1].aot.p, \\
&\qquad\qquad\qquad l \mapsto ServerMsg[s][1].aot.l, \\
&\qquad\qquad\qquad term \mapsto ServerMsg[s][1].term \\
&\qquad\qquad\qquad ]
\end{aligned}
$$

IN  LET $SubHbState \triangleq State[s]$
        $hb \triangleq [SubHbState \text{ EXCEPT } ![ServerMsg[s][1].s] = newState]$
    IN  $[State \text{ EXCEPT } ![s] = hb]$
$\land Cp' = $ LET $newcp \triangleq$
            primary node: compute new mcp
            IF $s \in Primary$ THEN
                LET $quorumAgree \triangleq QuorumAgreeInSameTerm(State[s])$IN
                    IF $Cardinality(quorumAgree) > 0$
                        THEN LET $QuorumSet \triangleq$ CHOOSE $i \in quorumAgree :$ TRUE
                                $serverInQuorum \triangleq$ CHOOSE $j \in QuorumSet :$ TRUE
                                $termOfQuorum \triangleq State[s][serverInQuorum].term$
                                $StateSet \triangleq \{[p \mapsto State[s][j].p,\ l \mapsto State[s][j].l] : j \in QuorumSet$
                                $newCommitPoint \triangleq HLCMinSet(StateSet)$
                            IN  IF $termOfQuorum = CurrentTerm[s]$
                                    THEN
                                        $[p \mapsto newCommitPoint.p,\ l \mapsto newCommitPoint.l,\ term \mapsto$
                                    ELSE  $Cp[s]$
                    ELSE  $Cp[s]$
            secondary node: update mcp
                ELSE IF LET $msgCP \triangleq [p \mapsto ServerMsg[s][1].cp.p,\ l \mapsto ServerMsg[s][1].cp.l]$IN
                            $\land \neg HLCLt(msgCP,\ Cp[s])$
                            $\land \neg HLCLt(Ot[s],\ msgCP)$
                        THEN $ServerMsg[s][1].cp$
                        ELSE  $Cp[s]$
    IN  $[Cp \text{ EXCEPT } ![s] = newcp]$
$\land CurrentTerm' = [CurrentTerm \text{ EXCEPT } ![s] = Max(CurrentTerm[s],\ ServerMsg[s][1].term)]$
$\land ServerMsg' = $ LET $newServerMsg \triangleq [ServerMsg \text{ EXCEPT } ![s] = Tail(@)]$
            IN  (LET $appendMsg \triangleq [type \mapsto \text{"update\_position"},\ s \mapsto ServerMsg[s][1].s,\ aot \mapsto Se$
                        $ct \mapsto ServerMsg[s][1].ct,\ cp \mapsto ServerMsg[s][1].cp,\ term \mapsto ServerM$
                IN  (LET $newMsg \triangleq$ IF $s \in Primary \lor SyncSource[s] = Nil$
                                    THEN $newServerMsg$ If $s$ is primary, accept the $msg$, else f
                                    ELSE $[newServerMsg \text{ EXCEPT } ![SyncSource[s]] = Append$
                    IN  $newMsg$))
$\land$ UNCHANGED $\langle Primary,\ Secondary,\ Oplog,\ Store,\ Ot,$
                $Pt,\ ReadyToServe,\ SyncSource \rangle$

$NTPSync \triangleq$  simplify $NTP$ protocal
    $\land ReadyToServe > 0$
    $\land Pt' = [s \in Server \mapsto MaxPt]$
    $\land$ UNCHANGED $\langle Primary,\ Secondary,\ Oplog,\ Store,\ Ct,\ Ot,$
                $ServerMsg,\ Cp,\ State,\ CurrentTerm,\ ReadyToServe,\ SyncSource \rangle$

$AdvancePt \triangleq$
    $\land ReadyToServe > 0$
    $\land \exists s \in Server :$

$\quad\quad \wedge s \in Primary$ <span style="background:#ccc">for simplicity</span>
$\quad\quad \wedge Pt[s] \leq PtStop$
$\quad\quad \wedge Pt' = [Pt \text{ EXCEPT } ![s] = @ + 1]$ <span style="background:#ccc">advance physical time</span>
$\quad\quad \wedge BroadcastHeartbeat(s)$ <span style="background:#ccc">broadcast heartbeat periodly</span>
$\quad \wedge \text{UNCHANGED } \langle Primary,\ Secondary,\ Oplog,\ Store,\ Ct,\ Ot,\ State,$
$\quad\quad\quad\quad\quad\quad\quad\quad Cp,\ CurrentTerm,\ ReadyToServe,\ SyncSource\rangle$

<span style="background:#ccc">Replicate oplog from node $j$ to node $i$, and update related structures accordingly</span>
$Replicate \triangleq$
$\quad \wedge ReadyToServe > 0$
$\quad \wedge \exists\, i,\, j \in Server :$
$\quad\quad \wedge CanSyncFrom(i,\, j)$ <span style="background:#ccc">$i$ can sync from $j$ only need not to rollback</span>
$\quad\quad \wedge i \in Secondary$
$\quad\quad \wedge ReplicateOplog(i,\, j) \neq \langle\rangle$
$\quad\quad \wedge Oplog' = [Oplog \text{ EXCEPT } ![i] = @ \circ ReplicateOplog(i,\, j)]$
$\quad\quad \wedge Store' = [Store \text{ EXCEPT } ![i] = Store[j]]$
$\quad\quad \wedge Ct' = [Ct \text{ EXCEPT } ![i] = HLCMax(Ct[i],\ Ct[j])]$ <span style="background:#ccc">update $Ct[i]$</span>
$\quad\quad \wedge Ot' = [Ot \text{ EXCEPT } ![i] = HLCMax(Ot[i],\ Ot[j])]$ <span style="background:#ccc">update $Ot[i]$</span>
$\quad\quad \wedge Cp' = [Cp \text{ EXCEPT } ![i] = HLCMax(Cp[i],\ Cp[j])]$ <span style="background:#ccc">update $Cp[i]$</span>
$\quad\quad \wedge CurrentTerm' = [CurrentTerm \text{ EXCEPT } ![i] = Max(CurrentTerm[i],\ CurrentTerm[j])]$ <span style="background:#ccc">update $Curren$</span>
$\quad\quad \wedge State' =$
$\quad\quad\quad \text{LET } newState \triangleq [$
$\quad\quad\quad\quad\quad\quad p \mapsto Ot[j].p,$
$\quad\quad\quad\quad\quad\quad l \mapsto Ot[j].l,$
$\quad\quad\quad\quad\quad\quad term \mapsto CurrentTerm[j]$
$\quad\quad\quad\quad\quad ]$
$\quad\quad\quad \text{IN }\quad \text{LET } SubHbState \triangleq State[i]$
$\quad\quad\quad\quad\quad\quad\quad hb \triangleq [SubHbState \text{ EXCEPT } ![j] = newState]$
$\quad\quad\quad\quad\quad \text{IN }\quad [State \text{ EXCEPT } ![i] = hb]$ <span style="background:#ccc">update $j$'s state $i$ knows</span>
$\quad\quad \wedge \text{LET } msg \triangleq [type \mapsto \text{"update\_position"},\ s \mapsto i,\ aot \mapsto Ot'[i],\ ct \mapsto Ct'[i],\ cp \mapsto Cp'[i],\ term \mapsto Curr$
$\quad\quad\quad \text{IN }\quad ServerMsg' = [ServerMsg \text{ EXCEPT } ![j] = Append(ServerMsg[j],\ msg)]$
$\quad\quad \wedge SyncSource' = [SyncSource \text{ EXCEPT } ![i] = j]$
$\quad\quad \wedge CalState' = [CalState \text{ EXCEPT } ![i] = CalState[j]]$
$\quad\quad\quad \wedge \text{UNCHANGED } \langle Primary,\ Secondary,\ Pt,\ ReadyToServe\rangle$

<span style="background:#ccc">Rollback $i$'s oplog and recover it to $j$'s state</span>
<span style="background:#ccc">Recover to $j$'s state immediately to prevent internal client request</span>
$RollbackAndRecover \triangleq$
$\quad \wedge ReadyToServe > 0$
$\quad \wedge \exists\, i,\, j \in Server :$
$\quad\quad \wedge i \in Secondary$
$\quad\quad \wedge CanRollback(i,\, j)$
$\quad\quad \wedge \text{LET } cmp \triangleq RollbackCommonPoint(i,\, j) \text{ IN}$
$\quad\quad\quad \text{LET } commonLog \triangleq SubSeq(Oplog[i],\ 1,\ cmp)$
$\quad\quad\quad\quad\quad appendLog \triangleq SubSeq(Oplog[j],\ cmp + 1,\ Len(Oplog[j]))$

8

$$\text{IN} \quad Oplog' = [Oplog \text{ EXCEPT } ![i] = commonLog \circ appendLog]$$
$$\land\ CurrentTerm' = [CurrentTerm \text{ EXCEPT } ![i] = Max(CurrentTerm[i],\ CurrentTerm[j])]\ \boxed{\text{update } Curren}$$
$$\land\ Store' = [Store \text{ EXCEPT } ![i] = \ Store[j]]$$
$$\land\ Ct' = [Ct \text{ EXCEPT } ![i] = HLCMax(Ct[i],\ Ct[j])]\ \boxed{\text{update } Ct[i]}$$
$$\land\ Ot' = [Ot \text{ EXCEPT } ![i] = HLCMax(Ot[i],\ Ot[j])]\ \boxed{\text{update } Ot[i]}$$
$$\land\ Cp' = [Cp \text{ EXCEPT } ![i] = HLCMax(Cp[i],\ Cp[j])]\ \boxed{\text{update } Cp[i]}$$
$$\land\ State' =$$

$$\text{LET } newStatei \ \triangleq\ [$$
$$\qquad p \mapsto Ot'[i].p,$$
$$\qquad l \mapsto Ot'[j].l,$$
$$\qquad term \mapsto CurrentTerm'[i]$$
$$\quad ]$$
$$\quad newStatej \ \triangleq\ [$$
$$\qquad p \mapsto Ot[j].p,$$
$$\qquad l \mapsto Ot[j].l,$$
$$\qquad term \mapsto CurrentTerm[j]$$
$$\quad ]$$
$$\text{IN} \quad \text{LET } SubHbState \ \triangleq\ State[i]$$
$$\qquad hb \ \triangleq\ [SubHbState \text{ EXCEPT } ![i] = newStatei]\ \boxed{\text{update } i\text{'s self state (used in mcp computation}}$$
$$\qquad hb1 \ \triangleq\ [hb \text{ EXCEPT } ![j] = newStatej]\ \boxed{\text{update } j\text{'s state } i \text{ knows}}$$
$$\quad \text{IN} \quad [State \text{ EXCEPT } ![i] = hb1]$$
$$\land\ \text{LET } msg \ \triangleq\ [type \mapsto \text{``update\_position''},\ s \mapsto i,\ aot \mapsto Ot'[i],\ ct \mapsto Ct'[i],\ cp \mapsto Cp'[i],\ term \mapsto Curr$$
$$\quad \text{IN} \quad ServerMsg' = [ServerMsg \text{ EXCEPT } ![j] = Append(ServerMsg[j],\ msg)]$$
$$\land\ SyncSource' = [SyncSource \text{ EXCEPT } ![i] = j]$$
$$\land\ \text{UNCHANGED } \langle Primary,\ Secondary,\ Pt,\ ReadyToServe \rangle$$

$$ClientRequest \ \triangleq$$
$$\quad \land\ ReadyToServe > 0$$
$$\quad \land\ \exists\, s \in Server,\ k \in Key,\ v \in Value :$$
$$\quad \land\ s \in Primary$$
$$\quad \land\ Tick(s)$$
$$\quad \land\ Ot' = [Ot \text{ EXCEPT } ![s] = Ct'[s]]$$
$$\quad \land\ Store' = [Store \text{ EXCEPT } ![s][k] = v]$$
$$\quad \land\ Oplog' = \text{LET } entry \ \triangleq\ [k \mapsto k,\ v \mapsto v,\ ot \mapsto Ot'[s],\ term \mapsto CurrentTerm[s]]$$
$$\qquad\qquad\qquad\quad newLog \ \triangleq\ Append(Oplog[s],\ entry)$$
$$\qquad\qquad\quad \text{IN} \quad [Oplog \text{ EXCEPT } ![s] = newLog]$$
$$\quad \land\ State' =$$
$$\quad\quad \text{LET } newState \ \triangleq\ [$$
$$\qquad\qquad p \mapsto Ot'[s].p,$$
$$\qquad\qquad l \mapsto Ot'[s].l,$$
$$\qquad\qquad term \mapsto CurrentTerm[s]$$
$$\qquad\quad ]$$
$$\quad\quad \text{IN} \quad \text{LET } SubHbState \ \triangleq\ State[s]$$
$$\qquad\qquad hb \ \triangleq\ [SubHbState \text{ EXCEPT } ![s] = newState]$$
$$\qquad\quad \text{IN} \quad [State \text{ EXCEPT } ![s] = hb]\ \boxed{\text{update } i\text{'s state}}$$

$$\land \text{UNCHANGED} \ \langle Primary, \quad Secondary, \ ServerMsg,$$
$$Pt, \ Cp,$$
$$CurrentTerm, \ ReadyToServe, \ SyncSource\rangle$$

Next state for all configurations
$Next \ \triangleq \ \lor \ Replicate$
$\quad\quad\quad \lor \ AdvancePt$
$\quad\quad\quad \lor \ AdvanceCp$
$\quad\quad\quad \lor \ ServerTakeHeartbeat$
$\quad\quad\quad \lor \ ServerTakeUpdatePosition$
$\quad\quad\quad \lor \ Stepdown$
$\quad\quad\quad \lor \ RollbackAndRecover$
$\quad\quad\quad \lor \ TurnOnReadyToServe$
$\quad\quad\quad \lor \ ElectPrimary$
$\quad\quad\quad \lor \ ClientRequest$
$\quad\quad\quad \lor \ NTPSync$

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

---

Properties to check?

$IsLogPrefix(i,\, j) \ \triangleq$
$\quad \land \ Len(Oplog[i]) \leq Len(Oplog[j])$
$\quad \land \ Oplog[i] = SubSeq(Oplog[j],\, 1,\, Len(Oplog[i]))$

If two logs have the same last *log* entry term, then one is a prefix of the other (from Will)
$LastTermsEquivalentImplyPrefixes \ \triangleq$
$\quad \forall\, i,\, j \in Server :$
$\quad\quad LogTerm(i,\, Len(Oplog[i])) = LogTerm(j,\, Len(Oplog[j])) \Rightarrow$
$\quad\quad IsLogPrefix(i,\, j) \lor IsLogPrefix(j,\, i)$

Check whether terms are incremented monotoniclly (from Will)
$TermsMonotonic \ \triangleq$
$\quad \Box[\forall\, s \in Server : CurrentTerm'[s] \geq CurrentTerm[s]]_{vars}$

Check the *log* in *Primary* node is append only (from Will)
$PrimaryAppendOnly \ \triangleq$
$\quad \Box[\forall\, s \in Server : s \in Primary \Rightarrow Len(Oplog'[s]) \geq Len(Oplog[s])]_{vars}$

Never rollback oplog before common point (from Will & Raft *Mongo*
$NeverRollbackCommonPoint \ \triangleq$
$\quad \exists\, i,\, j \in Server : CanRollback(i,\, j) \Rightarrow$
$\quad\quad \text{LET } commonPoint \ \triangleq \ RollbackCommonPoint(i,\, j)$
$\quad\quad\quad\quad\ \ lastOplog \ \triangleq \ Oplog[i][commonPoint]$
$\quad\quad \text{IN} \quad HLCLt(Cp[i],\, lastOplog.ot)$

Eventually *log* correctness (from Will
$EventuallyLogsConverge \quad \triangleq \ \Diamond\Box[\forall\, s,\, t \in Server : s \neq t \Rightarrow Oplog[s] = Oplog[t]]_{vars}$

$EventuallyLogsNonEmpty \triangleq \Diamond(\exists\, s \in Server : Len(Oplog[s]) > 0)$

(from *RaftMongo*

$TwoPrimariesInSameTerm \triangleq$
  $\exists\, i,\, j \in Server :$
    $\wedge\, i \neq j$
    $\wedge\, CurrentTerm[i] = CurrentTerm[j]$
    $\wedge\, i \in Primary$
    $\wedge\, j \in Primary$

$NoTwoPrimariesInSameTerm \triangleq \neg TwoPrimariesInSameTerm$

Check if there is any cycle of sync source path (from *RaftMongo Sync*

$SyncSourceCycleTwoNode \triangleq$
  $\exists\, s,\, t \in Server :$
    $\wedge\, s \neq t$
    $\wedge\, SyncSource[s] = t$
    $\wedge\, SyncSource[t] = s$

$BoundedSeq(s,\, n) \triangleq [1 \,..\, n \rightarrow s]$

$SyncSourcePaths \triangleq$
  $\{p \in BoundedSeq(Server,\, Cardinality(Server)) :$
    $\forall\, i \in 1 \,..\, (Len(p) - 1) : SyncSource[p[i]] = p[i+1]\}$

$SyncSourcePath(i,\, j) \triangleq$
  $\exists\, p \in SyncSourcePaths :$
    $\wedge\, Len(p) > 1$
    $\wedge\, p[1] = i$
    $\wedge\, p[Len(p)] = j$

$SyncSourceCycle \triangleq$
  $\exists\, s \in Server : SyncSourcePath(s,\, s)$

$NonTrivialSyncCycle \triangleq SyncSourceCycle \wedge \neg SyncSourceCycleTwoNode$
$NoNonTrivialSyncCycle \triangleq \neg NonTrivialSyncCycle$

\ * Modification *History*
\ * Last modified *Wed* May 04 15:41:10 *CST* 2022 by *dh*
\ * Created *Mon Apr* 18 11:38:53 *CST* 2022 by *dh*