\ * 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

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

RECURSIVE $CreateState(\_, \_)$ init state
$CreateState(len, seq) \triangleq$
$\quad$ IF $len = 0$ THEN $seq$
$\quad\quad$ 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 $j$
$NotBehind(i, j) \triangleq \quad \vee LastTerm(i) > LastTerm(j)$
$\qquad\qquad\qquad\qquad \vee \wedge LastTerm(i) = LastTerm(j)$
$\qquad\qquad\qquad\qquad\quad \wedge Len(Oplog[i]) \geq Len(Oplog[j])$

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

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 $x \in s : \forall y \in s : \neg HLCLt(y, x)$

commit point
RECURSIVE $AddState(\_, \_, \_)$
$AddState(new, state, index) \triangleq$
$\quad$ IF $index = 1 \wedge HLCLt(new, state[1])$
$\qquad$ THEN $\langle new \rangle \circ state \backslash *$ less than the first
$\quad\quad$ ELSE IF $index = Len(state) + 1$
$\qquad$ THEN $state \circ \langle new \rangle$
$\quad\quad$ ELSE IF $HLCLt(new, state[index])$
$\qquad$ THEN $SubSeq(state, 1, index - 1) \circ \langle new \rangle \circ SubSeq(state, index, Len(state))$
$\quad\quad$ ELSE $AddState(new, state, index + 1)$

RECURSIVE $RemoveState(\_, \_, \_)$
$RemoveState(old, state, index) \triangleq$
$\quad$ IF $state[index] = old$
$\qquad$ THEN $SubSeq(state, 1, index - 1) \circ SubSeq(state, index + 1, Len(state))$
$\quad\quad$ ELSE $RemoveState(old, state, index + 1)$

$AdvanceState(new, old, state) \triangleq AddState(new, RemoveState(old, state, 1), 1)$

clock

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

$Tick(s) \triangleq Ct' =$ IF $Ct[s].p \geq Pt[s]$

$$\text{THEN } [Ct \text{ EXCEPT } ![s] = [p \mapsto @.p, \, l \mapsto @.l + 1]]$$
$$\text{ELSE } [Ct \text{ EXCEPT } ![s] = [p \mapsto Pt[s], \, l \mapsto 0]]$$

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

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 $i \in Secondary \wedge len\_i < len\_j$
                      THEN $SubSeq(Oplog[j], \, len\_i + 1, \, len\_j)$
                      ELSE $\langle\rangle$

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

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 $commonIndices \triangleq \{k \in \text{DOMAIN } Oplog[i] :$
                               $\wedge k \leq Len(Oplog[j])$
                               $\wedge Oplog[i][k] = Oplog[j][k]\}$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   IF  $quorums = \{\}$ THEN $Nil$

    ELSE  CHOOSE $x \in quorums$ : TRUE

    ELSE  $quorums$

   IN    $quorums$

*Init* Part

---

$InitPrimary \;\triangleq\; Primary = \{\text{CHOOSE } s \in Server : \text{TRUE}\}$
$InitSecondary \;\triangleq\; Secondary = Server \setminus Primary$
$InitOplog \;\triangleq\; Oplog = [s \in Server \mapsto \langle\rangle]$
$InitStore \;\triangleq\; Store = [n \in Server \cup Client \mapsto [k \in Key \mapsto Nil]]$
$InitCt \;\triangleq\; Ct = [n \in Server \cup Client \mapsto [p \mapsto 0,\, l \mapsto 0]]$
$InitOt \;\triangleq\; Ot = [n \in Server \cup Client \mapsto [p \mapsto 0,\, l \mapsto 0]]$
$InitServerMsg \;\triangleq\; ServerMsg = [s \in Server \mapsto \langle\rangle]$
$InitPt \;\triangleq\; Pt = [s \in Server \mapsto 1]$
$InitCp \;\triangleq\; Cp = [n \in Server \cup Client \mapsto [p \mapsto 0,\, l \mapsto 0]]$
$InitCalState \;\triangleq\; CalState = [s \in Server \mapsto CreateState(Cardinality(Server),\, \langle\rangle)]$

create initial *state*(*for calculate*)

$InitState \;\triangleq\; State = [s \in Server \mapsto [s0 \in Server \mapsto$
$$[p \mapsto 0,\, l \mapsto 0,\, term \mapsto 0]]]$$
$InitCurrentTerm \;\triangleq\; CurrentTerm = [s \in Server \mapsto 0]$
$InitReadyToServe \;\triangleq\; ReadyToServe = 0$
$InitSyncSource \;\triangleq\; SyncSource = [s \in Server \mapsto Nil]$

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

Next *State* Actions

Replication Protocol: possible actions

---

$TurnOnReadyToServe \;\triangleq\;$
    $\land\, ReadyToServe = 0$
    $\land\, \exists\, s \in Primary :$

4

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

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

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
$$\land\ Primary' = \text{LET } possiblePrimary\ \triangleq\ Primary \setminus majorNodes$$
$$\text{IN }\quad possiblePrimary \cup \{i\}$$
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
$$\land\ SyncSource' = [SyncSource \text{ EXCEPT } ![i] = Nil]$$
$$\land \text{ UNCHANGED } \langle Oplog,\ Store,\ Ct,\ Ot,\ ServerMsg,\ Pt,\ Cp,\ State,\ ReadyToServe\rangle$$

$AdvanceCp\ \triangleq$
$$\land\ ReadyToServe > 0$$
$$\land\ \exists\,s \in Primary:$$
$$Cp' = [Cp \text{ EXCEPT } ![s] = CalState[s][Cardinality(Server) \div 2 + 1]\ ]$$
$$\land \text{ UNCHANGED } \langle Primary,\ Secondary,\ Oplog,\ Store,\ Ct,\ Ot,$$
$$ServerMsg,\ Pt,\ State,\ CurrentTerm,\ ReadyToServe,\ SyncSource\rangle$$

$AdvanceCp\ \triangleq$
$$\land\ ReadyToServe > 0$$
$$\land\ \exists\,i \in Primary :$$

5

LET $quorumAgree \triangleq QuorumAgreeInSameTerm(State[i])$IN
$\quad \land quorumAgree \neq Nil$
$\quad \land$ LET $serverInQuorum \triangleq$ CHOOSE $s \in quorumAgree :$ TRUE
$\qquad termOfQuorum \triangleq State[i][serverInQuorum][3]$
$\qquad newCommitPoint \triangleq HLCMinSet(\{[p \mapsto State[i][s][1], l \mapsto State[i][s][2]] : s \in quorumAgre$
IN $\quad \land termOfQuorum = CurrentTerm[i]$
$\qquad \land$ LET $newCP \triangleq [p \mapsto newCommitPoint.p, l \mapsto newCommitPoint.l, term \mapsto termOfQuor$
$\qquad\quad Cp' = [Cp$ EXCEPT $![i] = newCP]$
$\land$ UNCHANGED $\langle Primary, Secondary, Oplog, Store, Ct, Ot,$
$\qquad\qquad ServerMsg, \ \ Pt, State, CurrentTerm, ReadyToServe, SyncSource\rangle$

heartbeatoplog$Ot$store
$ServerTakeHeartbeat \triangleq$
$\quad \land ReadyToServe > 0$
$\quad \land \exists s \in Server :$
$\qquad \land Len(ServerMsg[s]) \neq 0$ message channel is not empty
$\qquad \land ServerMsg[s][1].type = $ "heartbeat"
$\qquad \land Ct' = [Ct$ EXCEPT $![s] = HLCMax(Ct[s], ServerMsg[s][1].ct)]$
$\qquad \land State' =$
$\qquad\quad$ LET $newState \triangleq [$
$\qquad\qquad\quad p \mapsto ServerMsg[s][1].aot.p,$
$\qquad\qquad\quad l \mapsto ServerMsg[s][1].aot.l,$
$\qquad\qquad\quad term \mapsto ServerMsg[s][1].term$
$\qquad\qquad ]$
$\qquad\quad$ IN $\quad$ LET $SubHbState \triangleq State[s]$
$\qquad\qquad\qquad hb \triangleq [SubHbState$ EXCEPT $![ServerMsg[s][1].s] = newState]$
$\qquad\qquad\quad$ IN $\quad [State$ EXCEPT $![s] = hb]$
$\qquad \land Cp' = $ LET $newcp \triangleq$
$\qquad\qquad\qquad$ primary node: compute new mcp
$\qquad\qquad\qquad\quad$ IF $s \in Primary$ THEN
$\qquad\qquad\qquad\qquad$ LET $quorumAgree \triangleq QuorumAgreeInSameTerm(State[s])$IN
$\qquad\qquad\qquad\qquad\quad$ IF $Cardinality(quorumAgree) > 0$
$\qquad\qquad\qquad$ THEN LET $serverInQuorum \triangleq$ CHOOSE $i \in quorumAgree :$ TRUE
$\qquad\qquad\qquad\qquad termOfQuorum \triangleq State[s][serverInQuorum].term$
$\qquad\qquad\qquad\qquad newCommitPoint \triangleq HLCMinSet(\{[p \mapsto State[s][j].p, l \mapsto State[s][j].l] : j \in quorumAgree\})$
$\qquad\qquad\qquad\quad$ IN $\quad$ IF $termOfQuorum = CurrentTerm[s]$
$\qquad\qquad\qquad\qquad$ THEN $[p \mapsto newCommitPoint.p, l \mapsto newCommitPoint.l, term \mapsto termOfQuorum]$
$\qquad\qquad\qquad\qquad$ ELSE $Cp[s]$
$\qquad\qquad\qquad\qquad\quad$ THEN $[p \mapsto 2, l \mapsto 2, term \mapsto 2]$
$\qquad\qquad\qquad\qquad\quad$ ELSE $Cp[s]$
$\backslash *$ $\qquad\qquad\qquad\quad \land quorumAgree \neq Nil$
$\qquad\qquad\qquad\quad \land Cardinality(quorumAgree) > 0$
$\backslash *$ $\qquad\qquad\qquad\quad \land$ LET $serverInQuorum \triangleq$ CHOOSE $i \in quorumAgree:$ TRUE
$\backslash *$ $\qquad\qquad\qquad\qquad termOfQuorum \triangleq State[s][serverInQuorum].term$
$\backslash *$ $\qquad\qquad\qquad\qquad newCommitPoint \triangleq HLCMinSet(\{[p \mapsto State[s][i].p, l \mapsto State[s][i].l] : i \in quorumAgree\})$

```
\ *                              IN   IF  termOfQuorum = CurrentTerm[s]
               ∧ LET serverInQuorum ≜ CHOOSE i ∈ quorumAgree: TRUE
                   termOfQuorum ≜ State[s][serverInQuorum].term
               IN   IF  termOfQuorum = CurrentTerm[s]
\ *                        THEN  [ p ↦ newCommitPoint.p, l ↦ newCommitPoint.l, term ↦ termOfQuorum ]
                    THEN  [p ↦ 5, l ↦ 5, term ↦ termOfQuorum]
                    ELSE  Cp[s]
              secondary node: update mcp
                ELSE  IF LET  msgCP ≜ [p ↦ ServerMsg[s][1].cp.p, l ↦ ServerMsg[s][1].cp.l]IN
                            ∧ ¬HLCLt(msgCP, Cp[s])
                            ∧ ¬HLCLt(Ot[s], msgCP)
                      THEN  ServerMsg[s][1].cp
                      ELSE   Cp[s]
            IN    [Cp EXCEPT ![s] = newcp]
        ∧ ServerMsg′ = [ServerMsg EXCEPT ![s] = Tail(@)]
        ∧ CurrentTerm′ = [CurrentTerm EXCEPT ![s] = Max(CurrentTerm[s], ServerMsg[s][1].term)]
     ∧ UNCHANGED ⟨Primary, Secondary, Oplog, Store, Ot, Pt,
                    ReadyToServe, SyncSource⟩

ServerTakeUpdatePosition ≜
    ∧ ReadyToServe > 0
    ∧ ∃ s ∈ Server :
        ∧ Len(ServerMsg[s]) ≠ 0   message channel is not empty
        ∧ ServerMsg[s][1].type = "update_position"
        ∧ Ct′ = [Ct EXCEPT ![s] = HLCMax(Ct[s], ServerMsg[s][1].ct)]  update ct accordingly
        ∧ State′ =
          LET newState ≜ [
                    p ↦ ServerMsg[s][1].aot.p,
                    l ↦ ServerMsg[s][1].aot.l,
                    term ↦ ServerMsg[s][1].term
                ]
          IN   LET SubHbState ≜ State[s]
                   hb ≜ [SubHbState EXCEPT ![ServerMsg[s][1].s] = newState]
               IN   [State EXCEPT ![s] = hb]
        ∧ Cp′ = LET newcp ≜
                   primary node: compute new mcp
                   IF s ∈ Primary THEN
                       LET quorumAgree ≜ QuorumAgreeInSameTerm(State[s])IN
                           IF Cardinality(quorumAgree) > 0
                               THEN ∃ serverInQuorum ∈ quorumAgree :
                                   LET termOfQuorum ≜ State[s][serverInQuorum].term
                       StateSet ≜ {[p ↦ State[s][j].p, l ↦ State[s][j].l] : j ∈ quorumAgree}
                       newCommitPoint ≜ HLCMinSet(StateSet)
                                       newCommitPoint ≜ [p ↦ State[s][serverInQuorum].p, l ↦ State[s][se
                               IN   IF termOfQuorum = CurrentTerm[s]
```

7

$$\text{THEN } [p \mapsto newCommitPoint.p, \ l \mapsto newCommitPoint.l, \ term \mapsto term$$
$$\text{ELSE} \quad Cp[s]$$
$$\text{ELSE} \quad Cp[s]$$

$$\text{ELSE} \quad \text{IF LET } msgCP \triangleq [p \mapsto ServerMsg[s][1].cp.p, \ l \mapsto ServerMsg[s][1].cp.l]\text{IN}$$
$$\wedge \neg HLCLt(msgCP, \ Cp[s])$$
$$\wedge \neg HLCLt(Ot[s], \ msgCP)$$
$$\text{THEN } ServerMsg[s][1].cp$$
$$\text{ELSE} \quad Cp[s]$$
$$\text{IN} \quad [Cp \text{ EXCEPT } ![s] = newcp]$$
$$\wedge CurrentTerm' = [CurrentTerm \text{ EXCEPT } ![s] = Max(CurrentTerm[s], \ ServerMsg[s][1].term)]$$
$$\wedge ServerMsg' = \text{LET } newServerMsg \triangleq [ServerMsg \text{ EXCEPT } ![s] = Tail(@)]$$
$$\text{IN} \quad (\text{LET} \quad appendMsg \triangleq [type \mapsto \text{``update\_position''}, \ s \mapsto ServerMsg[s][1].s, \ aot \mapsto Se$$
$$ct \mapsto ServerMsg[s][1].ct, \ cp \mapsto ServerMsg[s][1].cp, \ term \mapsto ServerM$$
$$\text{IN} \quad (\text{LET } newMsg \triangleq \text{ IF } s \in Primary \vee SyncSource[s] = Nil$$
$$\text{THEN } newServerMsg \quad \text{If } s \text{ is primary, accept the } msg, \text{ else f}$$
$$\text{ELSE } [newServerMsg \text{ EXCEPT } ![SyncSource[s]] = Append$$
$$\text{IN} \quad newMsg))$$
$$\wedge \text{ UNCHANGED } \langle Primary, \ Secondary, \ Oplog, \ Store, \ Ot,$$
$$Pt, \ ReadyToServe, \ SyncSource \rangle$$

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

$$AdvancePt \triangleq$$
$$\wedge ReadyToServe > 0$$
$$\wedge \exists s \in Server :$$
$$\wedge s \in Primary \qquad \text{for simplicity}$$
$$\wedge Pt[s] \leq PtStop$$
$$\wedge Pt' = [Pt \text{ EXCEPT } ![s] = @ + 1] \quad \text{advance physical time}$$
$$\wedge BroadcastHeartbeat(s) \qquad \text{broadcast heartbeat periodly}$$
$$\wedge \text{ UNCHANGED } \langle Primary, \ Secondary, \ Oplog, \ Store, \ Ct, \ Ot, \ State,$$
$$Cp, \ CurrentTerm, \ ReadyToServe, \ SyncSource \rangle$$

Replicate oplog from node $j$ to node $i$, and update related structures accordingly

$$Replicate \triangleq$$
$$\wedge ReadyToServe > 0$$
$$\wedge \exists i, j \in Server :$$
$$\wedge CanSyncFrom(i, j) \quad i \text{ can sync from } j \text{ only need not to rollback}$$
$$\wedge i \in Secondary$$
$$\wedge ReplicateOplog(i, j) \neq \langle \rangle$$
$$\wedge Oplog' = [Oplog \text{ EXCEPT } ![i] = @ \circ ReplicateOplog(i, j)]$$
$$\wedge Store' = [Store \text{ EXCEPT } ![i] = Store[j]]$$

$$\land\ Ct' = [Ct \text{ EXCEPT } ![i] = HLCMax(Ct[i],\ Ct[j])] \quad \text{update } Ct[i]$$
$$\land\ Ot' = [Ot \text{ EXCEPT } ![i] = HLCMax(Ot[i],\ Ot[j])] \quad \text{update } Ot[i]$$
$$\land\ Cp' = [Cp \text{ EXCEPT } ![i] = HLCMax(Cp[i],\ Cp[j])] \quad \text{update } Cp[i]$$
$$\land\ CurrentTerm' = [CurrentTerm \text{ EXCEPT } ![i] = Max(CurrentTerm[i],\ CurrentTerm[j])] \quad \text{update } Curren$$
$$\land\ State' =$$

$\quad$ LET $newState\ \triangleq\ [$

$\qquad\qquad p \mapsto Ot[j].p,$

$\qquad\qquad l\ \mapsto Ot[j].l,$

$\qquad\qquad term \mapsto CurrentTerm[j]$

$\qquad\quad ]$

$\quad$ IN $\quad$ LET $SubHbState\ \triangleq\ State[i]$

$\qquad\qquad\quad hb\ \triangleq\ [SubHbState \text{ EXCEPT } ![j] = newState]$

$\qquad\qquad$ IN $\quad [State \text{ EXCEPT } ![i] = hb] \quad \text{update } j\text{'s state } i \text{ knows}$

$\land$ LET $msg\ \triangleq\ [type \mapsto \text{``update\_position''},\ s \mapsto i,\ aot \mapsto Ot'[i],\ ct \mapsto Ct'[i],\ cp \mapsto Cp'[i],\ term \mapsto Curr$

$\quad$ IN $\quad ServerMsg' = [ServerMsg \text{ EXCEPT } ![j] = Append(ServerMsg[j],\ msg)]$

$\land\ SyncSource' = [SyncSource \text{ EXCEPT } ![i] = j]$

$\land\ CalState' = [CalState \text{ EXCEPT } ![i] = CalState[j]]$

$\quad \land$ UNCHANGED $\langle Primary,\ Secondary,\ Pt,\ ReadyToServe \rangle$

Rollback $i$'s oplog and recover it to $j$'s state

Recover to $j$'s state immediately to prevent internal client request

$RollbackAndRecover\ \triangleq$

$\quad \land\ ReadyToServe > 0$

$\quad \land\ \exists\, i,\, j \in Server :$

$\qquad \land\ i \in Secondary$

$\qquad \land\ CanRollback(i,\ j)$

$\qquad \land$ LET $cmp\ \triangleq\ RollbackCommonPoint(i,\ j)$ IN

$\qquad\quad$ LET $commonLog\ \triangleq\ SubSeq(Oplog[i],\ 1,\ cmp)$

$\qquad\qquad\quad appendLog\ \ \triangleq\ SubSeq(Oplog[j],\ cmp + 1,\ Len(Oplog[j]))$

$\qquad\quad$ IN $\quad Oplog' = [Oplog \text{ EXCEPT } ![i] = commonLog \circ appendLog]$

$\qquad \land\ CurrentTerm' = [CurrentTerm \text{ EXCEPT } ![i] = Max(CurrentTerm[i],\ CurrentTerm[j])] \quad \text{update } Curren$

$\qquad \land\ Store' = [Store \text{ EXCEPT } ![i] = Store[j]]$

$\qquad \land\ Ct' = [Ct \text{ EXCEPT } ![i] = HLCMax(Ct[i],\ Ct[j])] \quad \text{update } Ct[i]$

$\qquad \land\ Ot' = [Ot \text{ EXCEPT } ![i] = HLCMax(Ot[i],\ Ot[j])] \quad \text{update } Ot[i]$

$\qquad \land\ Cp' = [Cp \text{ EXCEPT } ![i] = HLCMax(Cp[i],\ Cp[j])] \quad \text{update } Cp[i]$

$\qquad \land\ State' =$

$\qquad\quad$ LET $newStatei\ \triangleq\ [$

$\qquad\qquad\qquad p \mapsto Ot'[i].p,$

$\qquad\qquad\qquad l\ \mapsto Ot'[j].l,$

$\qquad\qquad\qquad term \mapsto CurrentTerm'[i]$

$\qquad\qquad\quad ]$

$\qquad\qquad newStatej\ \triangleq\ [$

$\qquad\qquad\qquad p \mapsto Ot[j].p,$

$\qquad\qquad\qquad l\ \mapsto Ot[j].l,$

$\qquad\qquad\qquad term \mapsto CurrentTerm[j]$

9

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

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

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

$$Spec \triangleq Init \wedge \square[Next]_{vars}$$

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

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

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

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

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

Eventually *log* correctness (from Will
$EventuallyLogsConverge\ \ \ \triangleq\ \Diamond\Box[\forall\, s, t \in Server : s \neq t \Rightarrow Oplog[s] = Oplog[t]]_{vars}$
$EventuallyLogsNonEmpty\ \triangleq\ \Diamond(\exists\, s \in Server : Len(Oplog[s]) > 0)$

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

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

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

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

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

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

---