

\ * To fix: heartbeat

MODULE <i>TunableMongoDB_Repl</i>	
EXTENDS <i>Naturals, FiniteSets, Sequences, TLC</i>	
constants and variables	
CONSTANTS <i>Client, Server,</i>	the set of clients and servers
<i>Key, Value,</i>	the set of keys and values
<i>Nil,</i>	model value, place holder
<i>PtStop</i>	max physical time
VARIABLES	
<i>Primary,</i>	Primary node
<i>Secondary,</i>	secondary nodes
<i>Oplog,</i>	<i>oplog[s]</i> : <i>oplog</i> at <i>server[s]</i>
<i>Store,</i>	<i>store[s]</i> : data stored at <i>server[s]</i>
<i>Ct,</i>	<i>Ct[s]</i> : cluster time at node <i>s</i>
<i>Ot,</i>	<i>Ot[s]</i> : the last applied operation time at server <i>s</i>
<i>ServerMsg,</i>	<i>ServerMsg[s]</i> : the channel of heartbeat msgs at server <i>s</i>
<i>Pt,</i>	<i>Pt[s]</i> : physical time at server <i>s</i>
<i>Cp,</i>	<i>Cp[s]</i> : majority commit point at server <i>s</i>
<i>State,</i>	<i>State[s]</i> : the latest <i>Ot</i> of all servers that server <i>s</i> knows
<i>CurrentTerm,</i>	<i>CurrentTerm[s]</i> : current election term at server <i>s</i> → updated in <i>update_position</i> , heartbeat and replicate
<i>ReadyToServe,</i>	equal to 0 before any primary is elected
<i>SyncSource</i>	<i>SyncSource[s]</i> : sync source of server node <i>s</i>
group related <i>vars</i> to optimize code	
<i>electionVars</i> $\triangleq \langle \textit{Primary}, \textit{Secondary} \rangle$	<i>vars</i> that are related to election
<i>storageVars</i> $\triangleq \langle \textit{Oplog}, \textit{Store} \rangle$	<i>vars</i> that are related to storage
<i>messageVar</i> $\triangleq \langle \textit{ServerMsg} \rangle$	var that is related to message
<i>serverVars</i> $\triangleq \langle \textit{Ot}, \textit{SyncSource} \rangle$	<i>vars</i> that each server node holds for itself
<i>learnableVars</i> $\triangleq \langle \textit{Ct}, \textit{State}, \textit{Cp}, \textit{CurrentTerm} \rangle$	<i>vars</i> that must learn from msgs
<i>timeVar</i> $\triangleq \langle \textit{Pt} \rangle$	var that is used for timing
<i>functionalVar</i> $\triangleq \langle \textit{ReadyToServe} \rangle$	var that is used for some extra function
ASSUME <i>Cardinality(Client)</i> ≥ 1	
ASSUME <i>Cardinality(Server)</i> ≥ 2	at least one primary and one secondary
ASSUME <i>Cardinality(Key)</i> ≥ 1	at least one object
ASSUME <i>Cardinality(Value)</i> ≥ 2	at least two values to update
Helpers	
<i>HLCLt</i> (<i>x, y</i>) \triangleq IF <i>x.p</i> < <i>y.p</i>	
THEN TRUE	
ELSE IF <i>x.p</i> = <i>y.p</i>	

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

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

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
 $NotBehind(i, j) \triangleq$ $\vee LastTerm(i) > LastTerm(j)$
 $\vee \wedge LastTerm(i) = LastTerm(j)$
 $\wedge Len(Oplog[i]) \geq Len(Oplog[j])$

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

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 $[Ct \text{ EXCEPT } ![s] = [p \mapsto @.p, l \mapsto @.l + 1]]$
 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 $msg \triangleq [type \mapsto \text{"heartbeat"}, s \mapsto s, aot \mapsto Ot[s],$
 $ct \mapsto Ct[s], cp \mapsto Cp[s], term \mapsto CurrentTerm[s]]$

$$\text{IN } ServerMsg' = [x \in Server \mapsto \text{IF } x = s \text{ THEN } ServerMsg[x] \\ \text{ELSE } Append(ServerMsg[x], msg)]$$

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

$$ReplicateOplog(i, j) \triangleq \\ \text{LET } len_i \triangleq Len(Oplog[i]) \\ len_j \triangleq Len(Oplog[j]) \\ \text{IN } \text{IF } i \in Secondary \wedge len_i < len_j \\ \text{THEN } SubSeq(Oplog[j], len_i + 1, len_j) \\ \text{ELSE } \langle \rangle$$

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.

$$RollbackCommonPoint(i, j) \triangleq \\ \text{LET } commonIndices \triangleq \{k \in \text{DOMAIN } Oplog[i] : \\ \wedge k \leq Len(Oplog[j]) \\ \wedge Oplog[i][k] = Oplog[j][k]\} \text{IN} \\ \text{IF } commonIndices = \{\} \text{ THEN } 0 \text{ ELSE } MaxVal(commonIndices)$$

The set of all *quorums*. This just calculates simple majorities, but the only important property is that every quorum overlaps with every other.

$$Quorum \triangleq \{i \in \text{SUBSET } (Server) : Cardinality(i) * 2 > Cardinality(Server)\}$$

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

$$\text{IN } quorums$$

Init Part

$InitPrimary \triangleq Primary = \{\text{CHOOSE } s \in Server : \text{TRUE}\}$
 $InitSecondary \triangleq Secondary = Server \setminus Primary$
 $InitOplog \triangleq Oplog = [s \in Server \mapsto \langle \rangle]$
 $InitStore \triangleq Store = [n \in Server \cup Client \mapsto [k \in Key \mapsto Nil]]$
 $InitCt \triangleq Ct = [n \in Server \cup Client \mapsto [p \mapsto 0, l \mapsto 0]]$
 $InitOt \triangleq Ot = [n \in Server \cup Client \mapsto [p \mapsto 0, l \mapsto 0]]$
 $InitServerMsg \triangleq ServerMsg = [s \in Server \mapsto \langle \rangle]$
 $InitPt \triangleq Pt = [s \in Server \mapsto 1]$
 $InitCp \triangleq Cp = [n \in Server \cup Client \mapsto [p \mapsto 0, l \mapsto 0]]$
 $InitState \triangleq State = [s \in Server \mapsto [s0 \in Server \mapsto$
 $\quad [p \mapsto 0, l \mapsto 0, term \mapsto 0]]]$
 $InitCurrentTerm \triangleq CurrentTerm = [s \in Server \mapsto 0]$
 $InitReadyToServe \triangleq ReadyToServe = 0$
 $InitSyncSource \triangleq SyncSource = [s \in Server \mapsto Nil]$

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

Next State Actions

Replication Protocol: possible actions

$TurnOnReadyToServe \triangleq$
 $\quad \wedge ReadyToServe = 0$
 $\quad \wedge \exists s \in Primary :$
 $\quad \quad \wedge CurrentTerm' = [CurrentTerm \text{ EXCEPT } ![s] = CurrentTerm[s] + 1]$
 $\quad \quad \wedge CurrentTerm' = [s \in Server \mapsto 1]$
 $\quad \quad \wedge ReadyToServe' = ReadyToServe + 1$
 $\quad \wedge \text{UNCHANGED } \langle electionVars, storageVars, serverVars, Ct, messageVar, timeVar, Cp, State \rangle$

$Stepdown \triangleq$
 $\quad \wedge ReadyToServe > 0$
 $\quad \wedge \exists s \in Primary :$
 $\quad \quad \wedge Primary' = Primary \setminus \{s\}$
 $\quad \quad \wedge Secondary' = Secondary \cup \{s\}$
 $\quad \wedge \text{UNCHANGED } \langle storageVars, serverVars, Ct, messageVar, timeVar, Cp, State, CurrentTerm, functionalV \rangle$

Todo: *Stepdown* when receiving a higher term heartbeat

There are majority nodes agree to elect node i to become primary

$ElectPrimary \triangleq$

$$\begin{aligned}
& \wedge \text{ReadyToServe} > 0 \\
& \wedge \exists i \in \text{Server} : \exists \text{majorNodes} \in \text{SUBSET}(\text{Server}) : \\
& \quad \wedge \forall j \in \text{majorNodes} : \wedge \text{NotBehind}(i, j) \\
& \quad \quad \wedge \text{CurrentTerm}[i] \geq \text{CurrentTerm}[j] \\
& \quad \wedge \text{IsMajority}(\text{majorNodes}) \\
& \quad \text{voted nodes for } i \text{ cannot be primary anymore} \\
& \quad \wedge \text{Primary}' = \text{LET } \text{possiblePrimary} \triangleq \text{Primary} \setminus \text{majorNodes} \\
& \quad \quad \text{IN } \text{possiblePrimary} \cup \{i\} \\
& \quad \text{add voted nodes into secondaries} \\
& \quad \wedge \text{Secondary}' = \text{LET } \text{possibleSecondary} \triangleq \text{Secondary} \cup \text{majorNodes} \\
& \quad \quad \text{IN } \text{possibleSecondary} \setminus \{i\} \\
& \quad \wedge \text{CurrentTerm}' = [index \in \text{Server} \mapsto \text{IF } index \in (\text{majorNodes} \cup \{i\}) \\
& \quad \quad \quad \text{THEN } \text{CurrentTerm}[i] + 1 \\
& \quad \quad \quad \text{ELSE } \text{CurrentTerm}[index]] \\
& \quad \text{A primary node do not have any sync source} \\
& \quad \wedge \text{SyncSource}' = [\text{SyncSource} \text{ EXCEPT } ![i] = \text{Nil}] \\
& \wedge \text{UNCHANGED } \langle \text{storageVars}, \text{Ct}, \text{Ot}, \text{messageVar}, \text{timeVar}, \text{Cp}, \text{State}, \text{functionalVar} \rangle \\
\\
\text{AdvanceCp} & \triangleq \\
& \wedge \text{ReadyToServe} > 0 \\
& \wedge \exists s \in \text{Primary} : \\
& \quad \text{LET } \text{newCp} \triangleq \\
& \quad \quad \text{LET } \text{quorumAgree} \triangleq \text{QuorumAgreeInSameTerm}(\text{State}[s]) \\
& \quad \quad \text{IN } \text{IF } \text{Cardinality}(\text{quorumAgree}) > 0 \\
& \quad \quad \quad \text{THEN LET } \text{QuorumSet} \triangleq \text{CHOOSE } i \in \text{quorumAgree} : \text{TRUE} \\
& \quad \quad \quad \quad \text{serverInQuorum} \triangleq \text{CHOOSE } j \in \text{QuorumSet} : \text{TRUE} \\
& \quad \quad \quad \quad \text{termOfQuorum} \triangleq \text{State}[s][\text{serverInQuorum}].\text{term} \\
& \quad \quad \quad \quad \text{StateSet} \triangleq \{[p \mapsto \text{State}[s][j].p, l \mapsto \text{State}[s][j].l] : j \in \text{QuorumSet}\} \\
& \quad \quad \quad \quad \text{newCommitPoint} \triangleq \text{HLCMinSet}(\text{StateSet}) \\
& \quad \quad \quad \quad \text{oldCommitPoint} \triangleq [p \mapsto \text{Cp}[s].p, l \mapsto \text{Cp}[s].l] \\
& \quad \quad \quad \quad \text{newCp must be greater than current Cp for primary to advance it} \\
& \quad \quad \quad \text{IN } \text{IF } \text{termOfQuorum} = \text{CurrentTerm}[s] \wedge \text{HLCLt}(\text{oldCommitPoint}, \text{newCommitPoint}) \\
& \quad \quad \quad \quad \text{THEN } [p \mapsto \text{newCommitPoint}.p, l \mapsto \text{newCommitPoint}.l, \text{term} \mapsto \text{termOfQuorum}] \\
& \quad \quad \quad \quad \text{ELSE } \text{Cp}[s] \\
& \quad \quad \quad \text{ELSE } \text{Cp}[s] \\
& \quad \quad \text{IN } \text{Cp}' = [\text{Cp} \text{ EXCEPT } ![s] = \text{newCp}] \\
& \wedge \text{UNCHANGED } \langle \text{electionVars}, \text{storageVars}, \text{serverVars}, \text{Ct}, \text{messageVar}, \text{timeVar}, \text{State}, \text{CurrentTerm}, \text{functionalVar} \rangle \\
\\
\text{ServerTakeHeartbeat} & \triangleq \\
& \wedge \text{ReadyToServe} > 0 \\
& \wedge \exists s \in \text{Server} : \\
& \quad \wedge \text{Len}(\text{ServerMsg}[s]) \neq 0 \quad \text{message channel is not empty} \\
& \quad \wedge \text{ServerMsg}[s][1].\text{type} = \text{"heartbeat"} \\
& \quad \wedge \text{CurrentTerm}[s] = \text{ServerMsg}[s][1].\text{term}
\end{aligned}$$

$$\begin{aligned}
& \wedge Ct' = [Ct \text{ EXCEPT } ![s] = HLCMax(Ct[s], ServerMsg[s][1].ct)] \\
& \wedge State' = \\
& \quad \text{LET } newState \triangleq [\\
& \quad \quad p \mapsto ServerMsg[s][1].aot.p, \\
& \quad \quad l \mapsto ServerMsg[s][1].aot.l, \\
& \quad \quad term \mapsto ServerMsg[s][1].term \\
& \quad] \\
& \text{IN } \text{LET } SubHbState \triangleq State[s] \\
& \quad hb \triangleq [SubHbState \text{ EXCEPT } ![ServerMsg[s][1].s] = newState] \\
& \quad \text{IN } [State \text{ EXCEPT } ![s] = hb] \\
& \wedge Cp' = \text{LET } newcp \triangleq \\
& \quad \text{primary node: compute new mcp} \\
& \quad \text{IF } s \in Primary \text{ THEN} \\
& \quad \quad \text{LET } quorumAgree \triangleq QuorumAgreeInSameTerm(State[s]) \text{ IN} \\
& \quad \quad \text{IF } Cardinality(quorumAgree) > 0 \\
& \quad \quad \quad \text{THEN LET } QuorumSet \triangleq \text{CHOOSE } i \in quorumAgree : \text{TRUE} \\
& \quad \quad \quad \quad serverInQuorum \triangleq \text{CHOOSE } j \in QuorumSet : \text{TRUE} \\
& \quad \quad \quad \quad termOfQuorum \triangleq State[s][serverInQuorum].term \\
& \quad \quad \quad \quad StateSet \triangleq \{[p \mapsto State[s][j].p, l \mapsto State[s][j].l] : j \in QuorumSet\} \\
& \quad \quad \quad \quad newCommitPoint \triangleq HLCMinSet(StateSet) \\
& \quad \quad \quad \text{IN } \text{IF } termOfQuorum = CurrentTerm[s] \\
& \quad \quad \quad \quad \text{THEN} \\
& \quad \quad \quad \quad \quad [p \mapsto newCommitPoint.p, l \mapsto newCommitPoint.l, term \mapsto \dots] \\
& \quad \quad \quad \quad \text{ELSE } Cp[s] \\
& \quad \quad \text{ELSE } Cp[s] \\
& \quad \quad \text{secondary node: update mcp} \\
& \quad \quad \text{ELSE IF LET } msgCP \triangleq [p \mapsto ServerMsg[s][1].cp.p, l \mapsto ServerMsg[s][1].cp.l] \text{ IN} \\
& \quad \quad \quad \wedge \neg HLCLt(msgCP, Cp[s]) \\
& \quad \quad \quad \wedge \neg HLCLt(Ot[s], msgCP) \\
& \quad \quad \quad \text{The term of } cp \text{ must equal to the } CurrentTerm \text{ of that node to advance it} \\
& \quad \quad \quad \wedge ServerMsg[s][1].term = CurrentTerm[s] \\
& \quad \quad \quad \text{THEN } ServerMsg[s][1].cp \\
& \quad \quad \quad \text{ELSE } Cp[s] \\
& \quad \quad \text{IN } [Cp \text{ EXCEPT } ![s] = newcp] \\
& \wedge ServerMsg' = [ServerMsg \text{ EXCEPT } ![s] = Tail(@)] \\
& \wedge CurrentTerm' = [CurrentTerm \text{ EXCEPT } ![s] = Max(CurrentTerm[s], ServerMsg[s][1].term)] \\
& \wedge \text{UNCHANGED } \langle electionVars, storageVars, serverVars, timeVar, functionalVar \rangle \\
\\
& ServerTakeUpdatePosition \triangleq \\
& \quad \wedge ReadyToServe > 0 \\
& \quad \wedge \exists s \in Server : \\
& \quad \quad \wedge Len(ServerMsg[s]) \neq 0 \quad \text{message channel is not empty} \\
& \quad \quad \wedge ServerMsg[s][1].type = \text{"update_position"} \\
& \quad \quad \wedge Ct' = [Ct \text{ EXCEPT } ![s] = HLCMax(Ct[s], ServerMsg[s][1].ct)] \quad \text{update } ct \text{ accordingly} \\
& \quad \quad \wedge State' =
\end{aligned}$$

```

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

 $NTPSync \triangleq$  simplify NTP protocol
 $\wedge ReadyToServe > 0$ 
 $\wedge Pt' = [s \in Server \mapsto MaxPt]$ 
 $\wedge \text{UNCHANGED } \langle electionVars, storageVars, serverVars, learnableVars, messageVar, functionalVar \rangle$ 

```

$AdvancePt \triangleq$
 $\wedge ReadyToServe > 0$
 $\wedge \exists s \in Server :$
 $\quad \wedge s \in Primary$ for simplicity
 $\quad \wedge Pt[s] \leq PtStop$
 $\quad \wedge Pt' = [Pt \text{ EXCEPT } ![s] = @ + 1]$ advance physical time
 $\quad \wedge BroadcastHeartbeat(s)$ broadcast heartbeat periodically
 $\wedge UNCHANGED \langle electionVars, storageVars, serverVars, learnableVars, functionalVar \rangle$

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

$Replicate \triangleq$
 $\wedge ReadyToServe > 0$
 $\wedge \exists i, j \in Server :$
 $\quad \wedge CanSyncFrom(i, j)$ i can sync from j only need not to rollback
 $\quad \wedge i \in Secondary$
 $\quad \wedge ReplicateOplog(i, j) \neq \langle \rangle$
 $\quad \wedge Oplog' = [Oplog \text{ EXCEPT } ![i] = @ \circ ReplicateOplog(i, j)]$
 $\quad \wedge Store' = [Store \text{ EXCEPT } ![i] = Store[j]]$
 $\quad \wedge Ct' = [Ct \text{ EXCEPT } ![i] = HLCMax(Ct[i], Ct[j])]$ update $Ct[i]$
 $\quad \wedge Ot' = [Ot \text{ EXCEPT } ![i] = HLCMax(Ot[i], Ot[j])]$ update $Ot[i]$
 $\quad \wedge Cp' = [Cp \text{ EXCEPT } ![i] = HLCMax(Cp[i], Cp[j])]$ update $Cp[i]$
 $\quad \wedge CurrentTerm' = [CurrentTerm \text{ EXCEPT } ![i] = Max(CurrentTerm[i], CurrentTerm[j])]$ update $CurrentTerm$
 $\quad \wedge State' =$
 $\quad \quad LET newState \triangleq [$
 $\quad \quad \quad p \mapsto Ot[j].p,$
 $\quad \quad \quad l \mapsto Ot[j].l,$
 $\quad \quad \quad term \mapsto CurrentTerm[j]$
 $\quad \quad]$
 $\quad IN LET SubHbState \triangleq State[i]$
 $\quad \quad hb \triangleq [SubHbState \text{ EXCEPT } ![j] = newState]$
 $\quad \quad IN [State \text{ EXCEPT } ![i] = hb]$ update j 's state i knows
 $\wedge LET msg \triangleq [type \mapsto \text{"update_position"}, s \mapsto i, aot \mapsto Ot'[i], ct \mapsto Ct'[i], cp \mapsto Cp'[i], term \mapsto CurrentTerm[j]]$
 $\quad IN ServerMsg' = [ServerMsg \text{ EXCEPT } ![j] = Append(ServerMsg[j], msg)]$
 $\quad \wedge SyncSource' = [SyncSource \text{ EXCEPT } ![i] = j]$
 $\wedge UNCHANGED \langle electionVars, timeVar, functionalVar \rangle$

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

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

$RollbackAndRecover \triangleq$
 $\wedge ReadyToServe > 0$
 $\wedge \exists i, j \in Server :$
 $\quad \wedge i \in Secondary$
 $\quad \wedge CanRollback(i, j)$
 $\quad \wedge LET cmp \triangleq RollbackCommonPoint(i, j) IN$
 $\quad \quad LET commonLog \triangleq SubSeq(Oplog[i], 1, cmp)$

$$\begin{aligned}
& \text{appendLog} \triangleq \text{SubSeq}(\text{Oplog}[j], \text{cmp} + 1, \text{Len}(\text{Oplog}[j])) \\
& \text{IN } \text{Oplog}' = [\text{Oplog} \text{ EXCEPT } ![i] = \text{commonLog} \circ \text{appendLog}] \\
& \wedge \text{CurrentTerm}' = [\text{CurrentTerm} \text{ EXCEPT } ![i] = \text{Max}(\text{CurrentTerm}[i], \text{CurrentTerm}[j])] \quad \text{update CurrentTerm}[j] \\
& \wedge \text{Store}' = [\text{Store} \text{ EXCEPT } ![i] = \text{Store}[j]] \\
& \wedge \text{Ct}' = [\text{Ct} \text{ EXCEPT } ![i] = \text{HLCMax}(\text{Ct}[i], \text{Ct}[j])] \quad \text{update Ct}[i] \\
& \wedge \text{Ot}' = [\text{Ot} \text{ EXCEPT } ![i] = \text{HLCMax}(\text{Ot}[i], \text{Ot}[j])] \quad \text{update Ot}[i] \\
& \wedge \text{Cp}' = [\text{Cp} \text{ EXCEPT } ![i] = \text{HLCMax}(\text{Cp}[i], \text{Cp}[j])] \quad \text{update Cp}[i] \\
& \wedge \text{State}' = \\
& \quad \text{LET } \text{newStatei} \triangleq [\\
& \quad \quad p \mapsto \text{Ot}'[i].p, \\
& \quad \quad l \mapsto \text{Ot}'[j].l, \\
& \quad \quad \text{term} \mapsto \text{CurrentTerm}'[i] \\
& \quad] \\
& \quad \text{newStatej} \triangleq [\\
& \quad \quad p \mapsto \text{Ot}[j].p, \\
& \quad \quad l \mapsto \text{Ot}[j].l, \\
& \quad \quad \text{term} \mapsto \text{CurrentTerm}[j] \\
& \quad] \\
& \text{IN } \text{LET } \text{SubHbState} \triangleq \text{State}[i] \\
& \quad \text{hb} \triangleq [\text{SubHbState} \text{ EXCEPT } ![i] = \text{newStatei}] \quad \text{update } i\text{'s self state (used in mcp computation)} \\
& \quad \text{hb1} \triangleq [\text{hb} \text{ EXCEPT } ![j] = \text{newStatej}] \quad \text{update } j\text{'s state } i \text{ knows} \\
& \quad \text{IN } [\text{State} \text{ EXCEPT } ![i] = \text{hb1}] \\
& \wedge \text{LET } \text{msg} \triangleq [\text{type} \mapsto \text{"update_position"}, s \mapsto i, \text{aot} \mapsto \text{Ot}'[i], \text{ct} \mapsto \text{Ct}'[i], \text{cp} \mapsto \text{Cp}'[i], \text{term} \mapsto \text{CurrentTerm}'[i]] \\
& \quad \text{IN } \text{ServerMsg}' = [\text{ServerMsg} \text{ EXCEPT } ![j] = \text{Append}(\text{ServerMsg}[j], \text{msg})] \\
& \wedge \text{SyncSource}' = [\text{SyncSource} \text{ EXCEPT } ![i] = j] \\
& \wedge \text{UNCHANGED } \langle \text{electionVars}, \text{timeVar}, \text{functionalVar} \rangle
\end{aligned}$$

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

IN $[State \text{ EXCEPT } ![s] = hb] \text{ update } i\text{'s state}$
 $\wedge \text{UNCHANGED } \langle electionVars, messageVar, timeVar, Cp, CurrentTerm, functionalVar, SyncSource \rangle$

Next state for all configurations

$Next \triangleq \vee Replicate$
 $\vee AdvancePt$
 $\vee AdvanceCp$
 $\vee ServerTakeHeartbeat$
 $\vee ServerTakeUpdatePosition$
 $\vee Stepdown$
 $\vee RollbackAndRecover$
 $\vee TurnOnReadyToServe$
 $\vee ElectPrimary$
 $\vee ClientRequest$
 $\vee NTPSync$

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

Properties to check?

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

If two logs have the same last log entry term, then one is a prefix of the other (from Will)
 $LastTermsEquivalentImpliesPrefixes \triangleq$

$\forall i, j \in Server :$
 $LogTerm(i, Len(Oplog[i])) = LogTerm(j, Len(Oplog[j])) \Rightarrow$
 $IsLogPrefix(i, j) \vee IsLogPrefix(j, i)$

Check whether terms are incremented monotonically (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*)
 $SyncSourceCycleTwoNode \triangleq$
 $\exists s, t \in Server :$
 $\wedge s \neq t$
 $\wedge SyncSource[s] = t$
 $\wedge SyncSource[t] = s$

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

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

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

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

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

\ * Modification History
 \ * Last modified *Wed May 11 22:55:31 CST 2022* by *dh*
 \ * Created *Mon Apr 18 11:38:53 CST 2022* by *dh*