

---

MODULE *TunableMongoDB\_Repl*

---

EXTENDS *Naturals, FiniteSets, Sequences, TLC*

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 <i>msgs</i> 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>CalState,</i>	<i>CalState</i> : sorted <i>State[Primary]</i>
	<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>

---

ASSUME *Cardinality(Client)* ≥ 1    at least one client

ASSUME *Cardinality(Server)* ≥ 2    at least one primary and one secondary

ASSUME *Cardinality(Key)* ≥ 1    at least one object

ASSUME *Cardinality(Value)* ≥ 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  $y$

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

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.**

$MaxVal(s) \triangleq CHOOSE x \in s : \forall y \in s : x \geq y$

**commit point**

RECURSIVE  $AddState(-, -, -)$

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

RECURSIVE  $RemoveState(-, -, -)$

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

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

**clock**

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

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

ELSE  $[Ct \text{ EXCEPT } ![s] = [p \mapsto Pt[s], l \mapsto 0]]$

heartbeat

Only *Primary* node sends heartbeat once advance pt

$BroadcastHeartbeat(s) \triangleq$   
 LET  $msg \triangleq [type \mapsto \text{"heartbeat"}, s \mapsto s, aot \mapsto Ot[s],$   
 $ct \mapsto Ct[s], cp \mapsto Cp[s], term \mapsto CurrentTerm[s]]$   
 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  $MaxVal(commonIndices)$

Init Part

---

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

$$\begin{aligned}
InitOt &\triangleq Ot = [n \in Server \cup Client \mapsto [p \mapsto 0, l \mapsto 0]] \\
InitServerMsg &\triangleq ServerMsg = [s \in Server \mapsto \langle \rangle] \\
InitPt &\triangleq Pt = [s \in Server \mapsto 1] \\
InitCp &\triangleq Cp = [n \in Server \cup Client \mapsto [p \mapsto 0, l \mapsto 0]] \\
InitCalState &\triangleq CalState = [s \in Server \mapsto CreateState(Cardinality(Server), \langle \rangle)] \\
&\quad \text{create initial state(for calculate)} \\
InitState &\triangleq State = [s \in Server \mapsto [s0 \in Server \mapsto \\
&\quad [p \mapsto 0, l \mapsto 0]]] \\
InitCurrentTerm &\triangleq CurrentTerm = [s \in Server \mapsto 0] \\
InitReadyToServe &\triangleq ReadyToServe = 0 \\
InitSyncSource &\triangleq SyncSource = [s \in Server \mapsto Nil] \\
Init &\triangleq \\
&\quad \wedge InitPrimary \wedge InitSecondary \wedge InitOplog \wedge InitStore \wedge InitCt \\
&\quad \wedge InitOt \wedge InitPt \wedge InitCp \wedge InitCalState \\
&\quad \wedge InitServerMsg \\
&\quad \wedge InitState \wedge InitCurrentTerm \wedge InitReadyToServe \\
&\quad \wedge InitSyncSource
\end{aligned}$$

Next State Actions

Replication Protocol: possible actions

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

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

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

$$\begin{aligned}
ElectPrimary &\triangleq \\
&\quad \wedge ReadyToServe > 0 \\
&\quad \wedge \exists i \in Server : \exists majorNodes \in \text{SUBSET}(Server) : \\
&\quad \quad \wedge \forall j \in majorNodes : \wedge NotBehind(i, j)
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{CurrentTerm}[i] \geq \text{CurrentTerm}[j] \\
& \wedge \text{IsMajority}(\text{majorNodes}) \\
& \text{voted nodes for } i \text{ cannot be primary anymore} \\
& \wedge \text{Primary}' = \text{LET } \text{possiblePrimary} \triangleq \text{Primary} \setminus \text{majorNodes} \\
& \quad \text{IN } \text{possiblePrimary} \cup \{i\} \\
& \text{add voted nodes into secondaries} \\
& \wedge \text{Secondary}' = \text{LET } \text{possibleSecondary} \triangleq \text{Secondary} \cup \text{majorNodes} \\
& \quad \text{IN } \text{possibleSecondary} \setminus \{i\} \\
& \wedge \text{CurrentTerm}' = [\text{index} \in \text{Server} \mapsto \text{IF } \text{index} \in (\text{majorNodes} \cup \{i\}) \\
& \quad \text{THEN } \text{CurrentTerm}[i] + 1 \\
& \quad \text{ELSE } \text{CurrentTerm}[\text{index}]] \\
& \text{A primary node do not have any sync source} \\
& \wedge \text{SyncSource}' = [\text{SyncSource} \text{ EXCEPT } ![i] = \text{Nil}] \\
& \wedge \text{UNCHANGED } \langle \text{Oplog}, \text{Store}, \text{Ct}, \text{Ot}, \text{ServerMsg}, \text{Pt}, \text{Cp}, \text{State}, \text{CalState}, \\
& \quad \text{ReadyToServe} \rangle \\
\\
\text{AdvanceCp} & \triangleq \\
& \wedge \text{ReadyToServe} > 0 \\
& \wedge \exists s \in \text{Primary} : \\
& \quad \text{Cp}' = [\text{Cp} \text{ EXCEPT } ![s] = \text{CalState}[s][\text{Cardinality}(\text{Server}) \div 2 + 1]] \\
& \wedge \text{UNCHANGED } \langle \text{Primary}, \text{Secondary}, \text{Oplog}, \text{Store}, \text{Ct}, \text{Ot}, \\
& \quad \text{ServerMsg}, \text{Pt}, \text{CalState}, \\
& \quad \text{State}, \text{CurrentTerm}, \text{ReadyToServe}, \text{SyncSource} \rangle \\
\\
& \text{heartbeatoplogOtstore} \\
\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{Ct}' = [\text{Ct} \text{ EXCEPT } ![s] = \text{HLCMax}(\text{Ct}[s], \text{ServerMsg}[s][1].\text{ct})] \\
& \quad \wedge \text{State}' = \\
& \quad \quad \text{LET } \text{SubHbState} \triangleq \text{State}[s] \\
& \quad \quad \text{hb} \triangleq [\text{SubHbState} \text{ EXCEPT } ![\text{ServerMsg}[s][1].s] = \\
& \quad \quad \quad \text{ServerMsg}[s][1].\text{aot}] \\
& \quad \quad \text{IN } [\text{State} \text{ EXCEPT } ![s] = \text{hb}] \\
& \quad \wedge \text{CalState}' = \text{LET } \text{newcal} \triangleq \\
& \quad \quad \text{IF } s \in \text{Primary} \quad \text{primary node: update CalState} \\
& \quad \quad \quad \text{THEN } [\text{CalState} \text{ EXCEPT } ![s] = \\
& \quad \quad \quad \quad \text{AdvanceState}(\text{ServerMsg}[s][1].\text{aot}, \text{State}[s][\text{ServerMsg}[s][1].s], \text{CalState})] \\
& \quad \quad \quad \text{ELSE } \text{CalState} \\
& \quad \quad \text{IN } \text{newcal} \\
& \quad \wedge \text{Cp}' = \text{LET } \text{newcp} \triangleq \\
& \quad \quad \text{primary node: compute new mcp} \\
& \quad \quad \text{IF } s \in \text{Primary} \text{ THEN } \text{CalState}'[s][\text{Cardinality}(\text{Server}) \div 2 + 1]
\end{aligned}$$

```

secondary node: update mcp
ELSE IF  $\neg HLCLt(ServerMsg[s][1].cp, Cp[s])$ 
       $\wedge \neg HLCLt(Ot[s], ServerMsg[s][1].cp)$ 
      THEN  $ServerMsg[s][1].cp$ 
      ELSE  $Cp[s]$ 
IN  $[Cp \text{ EXCEPT } ![s] = newcp]$ 
 $\wedge ServerMsg' = [ServerMsg \text{ EXCEPT } ![s] = Tail(@)]$ 
 $\wedge CurrentTerm' = [CurrentTerm \text{ EXCEPT } ![s] = Max(CurrentTerm[s], ServerMsg[s][1].term)]$ 
 $\wedge \text{UNCHANGED } \langle Primary, Secondary, Oplog, Store, Ot, Pt,$ 
       $ReadyToServe, SyncSource \rangle$ 

```

$ServerTakeUpdatePosition \triangleq$

$\wedge ReadyToServe > 0$

$\wedge \exists s \in Server :$

$\wedge Len(ServerMsg[s]) \neq 0$  message channel is not empty

$\wedge ServerMsg[s][1].type = \text{"update\_position"}$

$\wedge Ct' = [Ct \text{ EXCEPT } ![s] = HLCMax(Ct[s], ServerMsg[s][1].ct)]$  update  $ct$  accordingly

$\wedge State' =$

LET  $SubHbState \triangleq State[s]$

$hb \triangleq [SubHbState \text{ EXCEPT } ![ServerMsg[s][1].s] =$   
 $ServerMsg[s][1].aot]$

IN  $[State \text{ EXCEPT } ![s] = hb]$

$\wedge CalState' = \text{LET } newcal \triangleq$

IF  $s \in Primary$  primary node: update  $CalState$

THEN  $[CalState \text{ EXCEPT } ![s] =$

$AdvanceState(ServerMsg[s][1].aot, State[s][ServerMsg[s][1].s], CalState[s])$

ELSE  $CalState$

IN  $newcal$

$\wedge Cp' = \text{LET } newcp \triangleq$

primary node: compute new mcp

IF  $s \in Primary$  THEN  $CalState'[s][Cardinality(Server) \div 2 + 1]$

secondary node: update mcp

ELSE IF  $\neg HLCLt(ServerMsg[s][1].cp, Cp[s])$

$\wedge \neg HLCLt(Ot[s], ServerMsg[s][1].cp)$

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  $appendMsg$

ELSE  $[newServerMsg \text{ EXCEPT } ![SyncSource[s]] = AppendMsg[s]]$

IN  $newMsg$ )

$\wedge$  UNCHANGED  $\langle \text{Primary}, \text{Secondary}, \text{Oplog}, \text{Store}, \text{Ot}, \text{Pt}, \text{ReadyToServe}, \text{SyncSource} \rangle$

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

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

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

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

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

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

$\text{RollbackAndRecover} \triangleq$

$$\begin{aligned}
& \wedge \text{ReadyToServe} > 0 \\
& \wedge \exists i, j \in \text{Server} : \\
& \quad \wedge i \in \text{Secondary} \\
& \quad \wedge \text{CanRollback}(i, j) \\
& \quad \wedge \text{LET } \text{cmp} \triangleq \text{RollbackCommonPoint}(i, j) \text{ IN} \\
& \quad \quad \text{LET } \text{commonLog} \triangleq \text{SubSeq}(\text{Oplog}[i], 1, \text{cmp}) \\
& \quad \quad \quad \text{appendLog} \triangleq \text{SubSeq}(\text{Oplog}[j], \text{cmp} + 1, \text{Len}(\text{Oplog}[j])) \\
& \quad \quad \text{IN } \text{Oplog}' = [\text{Oplog} \text{ EXCEPT } ![i] = \text{commonLog} \circ \text{appendLog}] \\
& \quad \text{CurrentTerm}' = [\text{CurrentTerm} \text{ EXCEPT } ![i] = \text{Max}(\text{CurrentTerm}[i], \text{CurrentTerm}[j])] \quad \text{update CurrentTerm} \\
& \quad \text{Store}' = [\text{Store} \text{ EXCEPT } ![i] = \text{Store}[j]] \\
& \quad \text{Ct}' = [\text{Ct} \text{ EXCEPT } ![i] = \text{HLCMax}(\text{Ct}[i], \text{Ct}[j])] \quad \text{update Ct}[i] \\
& \quad \text{Ot}' = [\text{Ot} \text{ EXCEPT } ![i] = \text{HLCMax}(\text{Ot}[i], \text{Ot}[j])] \quad \text{update Ot}[i] \\
& \quad \text{Cp}' = [\text{Cp} \text{ EXCEPT } ![i] = \text{HLCMax}(\text{Cp}[i], \text{Cp}[j])] \quad \text{update Cp}[i] \\
& \quad \text{State}' = \text{LET } \text{SubHbState} \triangleq \text{State}[i] \\
& \quad \quad \text{hb} \triangleq [\text{SubHbState} \text{ EXCEPT } ![j] = \text{Ot}[j]] \\
& \quad \quad \text{IN } [\text{State} \text{ EXCEPT } ![i] = \text{hb}] \quad \text{update j's state i knows} \\
& \quad \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]] \\
& \quad \text{IN } \text{ServerMsg}' = [\text{ServerMsg} \text{ EXCEPT } ![j] = \text{Append}(\text{ServerMsg}[j], \text{msg})] \\
& \quad \text{SyncSource}' = [\text{SyncSource} \text{ EXCEPT } ![i] = j] \\
& \quad \wedge \text{UNCHANGED } \langle \text{Primary}, \text{Secondary}, \text{Pt}, \text{CalState}, \\
& \quad \quad \text{ReadyToServe} \rangle
\end{aligned}$$

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

$$\begin{aligned}
& \text{Next state for all configurations} \\
& \text{Next} \triangleq \vee \text{Replicate} \\
& \quad \vee \text{AdvancePt} \\
& \quad \vee \text{ServerTakeHeartbeat}
\end{aligned}$$



$\vee \text{ServerTakeUpdatePosition}$   
 $\vee \text{Stepdown}$   
 $\vee \text{RollbackAndRecover}$   
 $\vee \text{TurnOnReadyToServe}$   
 $\vee \text{ElectPrimary}$   
 $\vee \text{ClientRequest}$   
 $\vee \text{NTPSync}$

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

---

Properties to check?

$$\begin{aligned}
\text{IsLogPrefix}(i, j) &\triangleq \\
&\wedge \text{Len}(\text{Oplog}[i]) \leq \text{Len}(\text{Oplog}[j]) \\
&\wedge \text{Oplog}[i] = \text{SubSeq}(\text{Oplog}[j], 1, \text{Len}(\text{Oplog}[i]))
\end{aligned}$$

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

$$\begin{aligned}
\text{LastTermsEquivalentImplyPrefixes} &\triangleq \\
&\forall i, j \in \text{Server} : \\
&\quad \text{LogTerm}(i, \text{Len}(\text{Oplog}[i])) = \text{LogTerm}(j, \text{Len}(\text{Oplog}[j])) \Rightarrow \\
&\quad \text{IsLogPrefix}(i, j) \vee \text{IsLogPrefix}(j, i)
\end{aligned}$$

Check whether terms are incremented monotonically (from Will)

$$\begin{aligned}
\text{TermsMonotonic} &\triangleq \\
&\Box[\forall s \in \text{Server} : \text{CurrentTerm}'[s] \geq \text{CurrentTerm}[s]]_{\text{vars}}
\end{aligned}$$

Check the *log* in *Primary* node is append only (from Will)

$$\begin{aligned}
\text{PrimaryAppendOnly} &\triangleq \\
&\Box[\forall s \in \text{Server} : s \in \text{Primary} \Rightarrow \text{Len}(\text{Oplog}'[s]) \geq \text{Len}(\text{Oplog}[s])]_{\text{vars}}
\end{aligned}$$

Never rollback oplog before common point (from Will & Raft *Mongo*)

$$\begin{aligned}
\text{NeverRollbackCommonPoint} &\triangleq \\
&\exists i, j \in \text{Server} : \text{CanRollback}(i, j) \Rightarrow \\
&\quad \text{LET } \text{commonPoint} \triangleq \text{RollbackCommonPoint}(i, j) \\
&\quad \text{lastOplog} \triangleq \text{Oplog}[i][\text{commonPoint}] \\
&\quad \text{IN } \text{HLCLt}(\text{Cp}[i], \text{lastOplog.ot})
\end{aligned}$$

Eventually *log* correctness (from Will)

$$\begin{aligned}
\text{EventuallyLogsConverge} &\triangleq \Diamond\Box[\forall s, t \in \text{Server} : s \neq t \Rightarrow \text{Oplog}[s] = \text{Oplog}[t]]_{\text{vars}} \\
\text{EventuallyLogsNonEmpty} &\triangleq \Diamond(\exists s \in \text{Server} : \text{Len}(\text{Oplog}[s]) > 0)
\end{aligned}$$

(from *RaftMongo*)

$$\begin{aligned}
\text{TwoPrimariesInSameTerm} &\triangleq \\
&\exists i, j \in \text{Server} : \\
&\quad \wedge i \neq j \\
&\quad \wedge \text{CurrentTerm}[i] = \text{CurrentTerm}[j] \\
&\quad \wedge i \in \text{Primary}
\end{aligned}$$

$$\wedge j \in \text{Primary}$$

$$\text{NoTwoPrimariesInSameTerm} \triangleq \neg \text{TwoPrimariesInSameTerm}$$

$$\text{SyncSourceCycleTwoNode} \triangleq \text{Check if there is any cycle of sync source path (from RaftMongo Sync)}$$

$$\exists s, t \in \text{Server} :$$

$$\wedge s \neq t$$

$$\wedge \text{SyncSource}[s] = t$$

$$\wedge \text{SyncSource}[t] = s$$

$$\text{BoundedSeq}(s, n) \triangleq [1 \dots n \rightarrow s]$$

$$\text{SyncSourcePaths} \triangleq$$

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

$$\text{SyncSourcePath}(i, j) \triangleq$$

$$\exists p \in \text{SyncSourcePaths} :$$

$$\wedge \text{Len}(p) > 1$$

$$\wedge p[1] = i$$

$$\wedge p[\text{Len}(p)] = j$$

$$\text{SyncSourceCycle} \triangleq$$

$$\exists s \in \text{Server} : \text{SyncSourcePath}(s, s)$$

$$\text{NonTrivialSyncCycle} \triangleq \text{SyncSourceCycle} \wedge \neg \text{SyncSourceCycleTwoNode}$$

$$\text{NoNonTrivialSyncCycle} \triangleq \neg \text{NonTrivialSyncCycle}$$

---

\ \* Modification History  
 \ \* Last modified Fri Apr 22 10:10:51 CST 2022 by dh  
 \ \* Created Mon Apr 18 11:38:53 CST 2022 by dh