

\ * 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*
→ updated in *update_position*, heartbeat and replicate
SyncSource, *SyncSource[s]*: sync source of server node *s*

group related *vars* to optimize code

electionVars \triangleq $\langle \textit{Primary}, \textit{Secondary} \rangle$ *vars* that are related to election
storageVars \triangleq $\langle \textit{Oplog}, \textit{Store} \rangle$ *vars* that are related to storage
messageVar \triangleq $\langle \textit{ServerMsg} \rangle$ var that is related to message
serverVars \triangleq $\langle \textit{Ot}, \textit{SyncSource} \rangle$ *vars* that each server node holds for itself
learnableVars \triangleq $\langle \textit{Ct}, \textit{State}, \textit{Cp}, \textit{CurrentTerm} \rangle$ *vars* that must learn from msgs
timeVar \triangleq $\langle \textit{Pt} \rangle$ var that is used for timing

ASSUME *Cardinality(Client)* \geq 1 at least one client
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

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

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

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

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.

$MaxVal(s) \triangleq \text{CHOOSE } x \in s : \forall y \in s : x \geq y$
 $HLCMinSet(s) \triangleq \text{CHOOSE } x \in s : \forall y \in s : \neg HLClt(y, x)$

clock

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

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

heartbeat

Only *Primary* node sends heartbeat once advance pt

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

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

$ReplicateOplog(i, j) \triangleq$
 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 $commonIndices \triangleq \{k \in DOMAIN Oplog[i] :$
 $\wedge k \leq Len(Oplog[j])$
 $\wedge Oplog[i][k] = Oplog[j][k]\}$ 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.
 $Quorum \triangleq \{i \in SUBSET (Server) : Cardinality(i) * 2 > Cardinality(Server)\}$

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

compute a new common point according to new update position msg
 $ComputeNewCp(s) \triangleq$
 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  $[p \mapsto newCommitPoint.p, l \mapsto newCommitPoint.l, term \mapsto termOfQuorum]$ 
    ELSE Cp[s]
ELSE Cp[s]
secondary node: update mcp
ELSE IF Len(ServerMsg[s])  $\neq$  0 THEN
    LET msgCP  $\triangleq$   $[p \mapsto ServerMsg[s][1].cp.p, l \mapsto ServerMsg[s][1].cp.l]$  IN
    IF  $\wedge \neg HLCLt(msgCP, Cp[s])$ 
     $\wedge \neg HLCLt(Ot[s], msgCP)$ 
    THEN The term of cp must equal to the CurrentTerm of that node to advance it
     $\wedge ServerMsg[s][1].term = CurrentTerm[s]$ 
    THEN ServerMsg[s][1].cp
    ELSE Cp[s]
ELSE Cp[s]

GetNewState(s, d, np, nl, nterm)  $\triangleq$ 
    LET newSubState  $\triangleq$   $[p \mapsto np, l \mapsto nl, term \mapsto nterm]$ 
    sState  $\triangleq$  State[s]
    IN [sState EXCEPT ![d] = newSubState]

```

Init Part

```

InitPrimary  $\triangleq$  Primary = {CHOOSE  $s \in Server$  : TRUE}
InitSecondary  $\triangleq$  Secondary = Server \ 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 =  $[p \in Primary \mapsto 1] @@ [s \in Server \mapsto 0]$ 
InitSyncSource  $\triangleq$  SyncSource =  $[s \in Server \mapsto Nil]$ 

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

```

Next *State* Actions

Replication Protocol: possible actions

Stepdown \triangleq

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

Todo: *Stepdown* when receiving a higher term heartbeat

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

ElectPrimary \triangleq

$\wedge \exists i \in \text{Server} : \exists \text{majorNodes} \in \text{SUBSET } (\text{Server}) :$
 $\wedge \forall j \in \text{majorNodes} : \wedge \text{NotBehind}(i, j)$
 $\wedge \text{CurrentTerm}[i] \geq \text{CurrentTerm}[j]$
 $\wedge \text{IsMajority}(\text{majorNodes})$
 voted nodes for *i* cannot be primary anymore
 $\wedge \text{Primary}' = \text{LET } \text{possiblePrimary} \triangleq \text{Primary} \setminus \text{majorNodes}$
 $\text{IN } \text{possiblePrimary} \cup \{i\}$
 add voted nodes into secondaries
 $\wedge \text{Secondary}' = \text{LET } \text{possibleSecondary} \triangleq \text{Secondary} \cup \text{majorNodes}$
 $\text{IN } \text{possibleSecondary} \setminus \{i\}$
 $\wedge \text{CurrentTerm}' = [\text{index} \in \text{Server} \mapsto \text{IF } \text{index} \in (\text{majorNodes} \cup \{i\})$
 $\text{THEN } \text{CurrentTerm}[i] + 1$
 $\text{ELSE } \text{CurrentTerm}[\text{index}]]$

A primary node do not have any sync source

$\wedge \text{SyncSource}' = [\text{SyncSource} \text{ EXCEPT } ![i] = \text{Nil}]$
 $\wedge \text{UNCHANGED } \langle \text{storageVars}, \text{Ct}, \text{Ot}, \text{messageVar}, \text{timeVar}, \text{Cp}, \text{State} \rangle$

AdvanceCp \triangleq

$\wedge \exists s \in \text{Primary} :$
 $\text{LET } \text{newCp} \triangleq \text{ComputeNewCp}(s)$
 $\text{IN } \text{Cp}' = [\text{Cp} \text{ EXCEPT } ![s] = \text{newCp}]$
 $\wedge \text{UNCHANGED } \langle \text{electionVars}, \text{storageVars}, \text{serverVars}, \text{Ct}, \text{messageVar}, \text{timeVar}, \text{State}, \text{CurrentTerm} \rangle$

heartbeatoplogOstore

ServerTakeHeartbeat \triangleq

$\wedge \exists s \in \text{Server} :$
 $\wedge \text{Len}(\text{ServerMsg}[s]) \neq 0$ message channel is not empty
 $\wedge \text{ServerMsg}[s][1].\text{type} = \text{"heartbeat"}$
 $\wedge \text{CurrentTerm}[s] = \text{ServerMsg}[s][1].\text{term}$
 $\wedge \text{Ct}' = [\text{Ct} \text{ EXCEPT } ![s] = \text{HLCMax}(\text{Ct}[s], \text{ServerMsg}[s][1].\text{ct})]$
 $\wedge \text{State}' = \text{LET } \text{newState} \triangleq \text{GetNewState}(s, \text{ServerMsg}[s][1].s, \text{ServerMsg}[s][1].\text{aot.p}, \text{ServerMsg}[s][1].\text{ot})$
 $\text{IN } [\text{State} \text{ EXCEPT } ![s] = \text{newState}]$
 $\wedge \text{Cp}' = \text{LET } \text{newcp} \triangleq \text{ComputeNewCp}(s)$

$$\begin{aligned}
& \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 \rangle \\
ServerTakeUpdatePosition & \triangleq \\
& \wedge \exists s \in Server : \\
& \quad \wedge Len(ServerMsg[s]) \neq 0 \quad \text{message channel is not empty} \\
& \quad \wedge ServerMsg[s][1].type = \text{"update_position"} \\
& \quad \wedge Ct' = [Ct \text{ EXCEPT } ![s] = HLCMax(Ct[s], ServerMsg[s][1].ct)] \quad \text{update } ct \text{ accordingly} \\
& \quad \wedge State' = \text{LET } newState \triangleq GetNewState(s, ServerMsg[s][1].s, ServerMsg[s][1].aot.p, ServerMsg[s][1].term) \\
& \quad \quad \text{IN } [State \text{ EXCEPT } ![s] = newState] \\
& \quad \wedge Cp' = \text{LET } newcp \triangleq ComputeNewCp(s) \\
& \quad \quad \text{IN } [Cp \text{ EXCEPT } ![s] = newcp] \\
& \quad \wedge CurrentTerm' = [CurrentTerm \text{ EXCEPT } ![s] = Max(CurrentTerm[s], ServerMsg[s][1].term)] \\
& \quad \wedge ServerMsg' = \text{LET } newServerMsg \triangleq [ServerMsg \text{ EXCEPT } ![s] = Tail(@)] \\
& \quad \quad \text{IN } (\text{LET } appendMsg \triangleq [type \mapsto \text{"update_position"}, s \mapsto ServerMsg[s][1].s, aot \mapsto ServerMsg[s][1].aot.p, ct \mapsto ServerMsg[s][1].ct, cp \mapsto ServerMsg[s][1].cp, term \mapsto ServerMsg[s][1].term]) \\
& \quad \quad \quad \text{IN } (\text{LET } newMsg \triangleq \text{IF } s \in Primary \vee SyncSource[s] = Nil \\
& \quad \quad \quad \quad \text{THEN } newServerMsg \quad \text{If } s \text{ is primary, accept the } msg, \text{ else for } \\
& \quad \quad \quad \quad \text{ELSE } [newServerMsg \text{ EXCEPT } ![SyncSource[s]] = Append(appendMsg, newMsg)]) \\
& \quad \wedge \text{UNCHANGED } \langle electionVars, storageVars, serverVars, timeVar \rangle \\
NTPSync & \triangleq \quad \text{simplify } NTP \text{ protocol} \\
& \wedge Pt' = [s \in Server \mapsto MaxPt] \\
& \wedge \text{UNCHANGED } \langle electionVars, storageVars, serverVars, learnableVars, messageVar \rangle \\
AdvancePt & \triangleq \\
& \wedge \exists s \in Server : \\
& \quad \wedge s \in Primary \quad \text{for simplicity} \\
& \quad \wedge Pt[s] \leq PtStop \\
& \quad \wedge Pt' = [Pt \text{ EXCEPT } ![s] = @ + 1] \quad \text{advance physical time} \\
& \quad \wedge BroadcastHeartbeat(s) \quad \text{broadcast heartbeat periodically} \\
& \wedge \text{UNCHANGED } \langle electionVars, storageVars, serverVars, learnableVars \rangle \\
\text{Replicate oplog from node } j \text{ to node } i, \text{ and update related structures accordingly} \\
Replicate & \triangleq \\
& \wedge \exists i, j \in Server : \\
& \quad \wedge CanSyncFrom(i, j) \quad i \text{ can sync from } j \text{ 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])] \quad \text{update } Ct[i] \\
& \quad \wedge Ot' = [Ot \text{ EXCEPT } ![i] = HLCMax(Ot[i], Ot[j])] \quad \text{update } Ot[i]
\end{aligned}$$

$\wedge Cp' = [Cp \text{ EXCEPT } ![i] = HLCMax(Cp[i], Cp[j])] \text{ update } Cp[i]$
 $\wedge CurrentTerm' = [CurrentTerm \text{ EXCEPT } ![i] = Max(CurrentTerm[i], CurrentTerm[j])] \text{ update } CurrentTerm[i]$
 $\wedge State' = \text{LET } newState \triangleq GetNewState(i, j, Ot[j].p, Ot[j].l, CurrentTerm[j]) \text{ update } j\text{'s state } i \text{ knows}$
 $\quad \text{IN } [State \text{ EXCEPT } ![i] = newState]$
 $\wedge \text{LET } msg \triangleq [type \mapsto \text{"update_position"}, s \mapsto i, aot \mapsto Ot'[i], ct \mapsto Ct'[i], cp \mapsto Cp'[i], term \mapsto CurrentTerm[j]]$
 $\quad \text{IN } ServerMsg' = [ServerMsg \text{ EXCEPT } ![j] = Append(ServerMsg[j], msg)]$
 $\wedge SyncSource' = [SyncSource \text{ EXCEPT } ![i] = j]$
 $\wedge \text{UNCHANGED } \langle electionVars, timeVar \rangle$

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

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

$RollbackAndRecover \triangleq$

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

$ClientRequest \triangleq$

$\wedge \exists s \in Server, k \in Key, v \in Value :$
 $\quad \wedge s \in Primary$
 $\quad \wedge Tick(s)$
 $\quad \wedge Ot' = [Ot \text{ EXCEPT } ![s] = Ct'[s]]$
 $\quad \wedge Store' = [Store \text{ EXCEPT } ![s][k] = v]$
 $\quad \wedge Oplog' = \text{LET } entry \triangleq [k \mapsto k, v \mapsto v, ot \mapsto Ot'[s], term \mapsto CurrentTerm[s]]$
 $\quad \quad \quad newLog \triangleq Append(Oplog[s], entry)$
 $\quad \quad \text{IN } [Oplog \text{ EXCEPT } ![s] = newLog]$

$$\begin{aligned} \wedge State' = \text{LET } newState \triangleq & \text{GetNewState}(s, s, Ot'[s].p, Ot'[s].l, CurrentTerm[s]) \quad \text{update } i\text{'s state} \\ & \text{IN } [State \text{ EXCEPT } ![s] = newState] \\ \wedge \text{UNCHANGED } & \langle electionVars, messageVar, timeVar, Cp, CurrentTerm, SyncSource \rangle \end{aligned}$$

Next state for all configurations

$$\begin{aligned} Next \triangleq & \vee Replicate \\ & \vee AdvancePt \\ & \vee AdvanceCp \\ & \vee ServerTakeHeartbeat \\ & \vee ServerTakeUpdatePosition \\ & \vee Stepdown \\ & \vee RollbackAndRecover \\ & \vee ElectPrimary \\ & \vee ClientRequest \\ & \vee NTPSync \end{aligned}$$

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

Properties to check?

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

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

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

Check whether terms are incremented monotonically (from Will)

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

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

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

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

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

Eventually log correctness (from Will)

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

(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 *Tue May 24 16:55:51 CST 2022* by *dh*
\ * Created *Mon Apr 18 11:38:53 CST 2022* by *dh*