

<hr/> <div> <div>MODULE <i>Zab</i></div> <div> This is the formal specification for the <i>Zab</i> consensus algorithm, in <i>DSN'2011</i>, which means <i>Zab pre</i> – 1.0. </div> </div> <hr/>	
EXTENDS <i>Integers, FiniteSets, Sequences, Naturals, TLC</i>	
<hr/> <div> <div>The set of servers</div> <div>CONSTANT <i>Server</i></div> <div>States of server</div> <div>CONSTANTS <i>LOOKING, FOLLOWING, LEADING</i></div> <div><i>Zab</i> states of server</div> <div>CONSTANTS <i>ELECTION, DISCOVERY, SYNCHRONIZATION, BROADCAST</i></div> <div>Message types</div> <div>CONSTANTS <i>CEPOCH, NEWEPOCH, ACKEPOCH, NEWLEADER, ACKLD, COMMITLD, PROPOSE, A</i></div> <div><i>[MaxTimeoutFailures, MaxTransactionNum, MaxEpoch, MaxRestarts]</i></div> <div>CONSTANT <i>Parameters</i></div> <div> $MAXEPOCH \triangleq 10$ $NullPoint \triangleq \text{CHOOSE } p : p \notin Server$ $Quorums \triangleq \{Q \in \text{SUBSET } Server : Cardinality(Q) * 2 > Cardinality(Server)\}$ </div> </div> <hr/>	
<div>Variables that all servers use.</div> <div>VARIABLES</div> <div> <i>state,</i> <i>zabState,</i> <i>acceptedEpoch,</i> <i>currentEpoch,</i> <i>history,</i> <i>lastCommitted</i> </div>	<div>State of server, in <i>{LOOKING, FOLLOWING, LEADING}</i>.</div> <div>Current phase of server, in <i>{ELECTION, DISCOVERY, SYNCHRONIZATION, BROADCAST}</i>.</div> <div>Epoch of the last <i>LEADERINFO</i> packet accepted, namely <i>f.p</i> in paper.</div> <div>Epoch of the last <i>NEWLEADER</i> packet accepted, namely <i>f.a</i> in paper.</div> <div>History of servers: sequence of transactions, containing: [<i>zxid</i>, value, <i>ackSid</i>, epoch].</div> <div>Maximum index and <i>zxid</i> known to be committed, namely 'lastCommitted' in Leader. Starts from 0, and increases monotonically before restarting.</div>
<div>Variables only used for leader.</div> <div>VARIABLES</div> <div> <i>learners,</i> <i>ceepochRecv,</i> <i>ackeRecv,</i> <i>ackldRecv,</i> <i>sendCounter</i> </div>	<div>Set of servers leader connects.</div> <div>Set of learners leader has received <i>CEPOCH</i> from.</div> <div>Set of record [<i>sid</i>, connected, epoch], where epoch means <i>f.p</i> from followers.</div> <div>Set of learners leader has received <i>ACEPOCH</i> from.</div> <div>Set of record [<i>sid</i>, connected, <i>peerLastEpoch</i>, <i>peerHistory</i>], to record <i>f.a</i> and <i>h(f)</i> from followers.</div> <div>Set of learners leader has received <i>ACKLD</i> from.</div> <div>Set of record [<i>sid</i>, connected].</div> <div>Count of <i>txns</i> leader has broadcast.</div>

Variables only used for follower.

VARIABLES *leaderAddr* If follower has connected with leader.
If follower lost connection, then null.

Variable representing oracle of leader.

VARIABLE *leaderOracle* Current oracle.

Variables about network channel.

VARIABLE *msgs* Simulates network channel.
msgs[i][j] means the input buffer of server *j*
from server *i*.

Variables only used in verifying properties.

VARIABLES *epochLeader*, Set of leaders in every epoch.
proposalMsgsLog, Set of all broadcast messages.
violatedInvariants Check whether there are conditions
contrary to the facts.

Variable used for recording critical data,
to constrain state space or update values.

VARIABLE *recorder* Consists: members of *Parameters* and *pc*, values.
Form is record:
[*pc*, *nTransaction*, *maxEpoch*, *nTimeout*, *nRestart*, *nClientRequest*]

$serverVars \triangleq \langle state, zabState, acceptedEpoch, currentEpoch, history, lastCommitted \rangle$

$leaderVars \triangleq \langle learners, cepochRecv, ackeRecv, ackldRecv, sendCounter \rangle$

$followerVars \triangleq leaderAddr$

$electionVars \triangleq leaderOracle$

$msgVars \triangleq msgs$

$verifyVars \triangleq \langle proposalMsgsLog, epochLeader, violatedInvariants \rangle$

$vars \triangleq \langle serverVars, leaderVars, followerVars, electionVars, msgVars, verifyVars, recorder \rangle$

Return the maximum value from the set *S*

$Maximum(S) \triangleq \text{IF } S = \{\} \text{ THEN } -1$
ELSE CHOOSE $n \in S : \forall m \in S : n \geq m$

Return the minimum value from the set *S*

$Minimum(S) \triangleq \text{IF } S = \{\} \text{ THEN } -1$
ELSE CHOOSE $n \in S : \forall m \in S : n \leq m$

Check server state

$IsLeader(s) \triangleq state[s] = LEADING$
 $IsFollower(s) \triangleq state[s] = FOLLOWING$
 $IsLooking(s) \triangleq state[s] = LOOKING$

Check if s is a quorum

$IsQuorum(s) \triangleq s \in Quorums$

$IsMyLearner(i, j) \triangleq j \in learners[i]$
 $IsMyLeader(i, j) \triangleq leaderAddr[i] = j$
 $HasNoLeader(i) \triangleq leaderAddr[i] = NullPoint$
 $HasLeader(i) \triangleq leaderAddr[i] \neq NullPoint$

FALSE: $zxid1 \leq zxid2$; TRUE: $zxid1 > zxid2$

$ZxidCompare(zxid1, zxid2) \triangleq \vee zxid1[1] > zxid2[1]$
 $\vee \wedge zxid1[1] = zxid2[1]$
 $\wedge zxid1[2] > zxid2[2]$

$ZxidEqual(zxid1, zxid2) \triangleq zxid1[1] = zxid2[1] \wedge zxid1[2] = zxid2[2]$

$TxnZxidEqual(txn, z) \triangleq txn.zxid[1] = z[1] \wedge txn.zxid[2] = z[2]$

$TxnEqual(txn1, txn2) \triangleq \wedge ZxidEqual(txn1.zxid, txn2.zxid)$
 $\wedge txn1.value = txn2.value$

$EpochPrecedeInTxn(txn1, txn2) \triangleq txn1.zxid[1] < txn2.zxid[1]$

Actions about recorder

$GetParameter(p) \triangleq \text{IF } p \in \text{DOMAIN } Parameters \text{ THEN } Parameters[p] \text{ ELSE } 0$
 $GetRecorder(p) \triangleq \text{IF } p \in \text{DOMAIN } recorder \text{ THEN } recorder[p] \text{ ELSE } 0$

$RecorderGetHelper(m) \triangleq (m := recorder[m])$

$RecorderIncHelper(m) \triangleq (m := recorder[m] + 1)$

$RecorderIncTimeout \triangleq RecorderIncHelper("nTimeout")$

$RecorderGetTimeout \triangleq RecorderGetHelper("nTimeout")$

$RecorderIncRestart \triangleq RecorderIncHelper("nRestart")$

$RecorderGetRestart \triangleq RecorderGetHelper("nRestart")$

$RecorderSetTransactionNum(pc) \triangleq ("nTransaction" :=$

IF $pc[1] = \text{"LeaderProcessRequest"}$ THEN

LET $s \triangleq \text{CHOOSE } i \in Server :$

$\forall j \in Server : Len(history'[i]) \geq Len(history'[j])$

IN $Len(history'[s])$

ELSE $recorder["nTransaction"]$)

$\triangleq ("maxEpoch" :=$

IF $pc[1] = \text{"LeaderProcessCEPOCH"}$ THEN

LET $s \triangleq \text{CHOOSE } i \in Server :$

$\forall j \in Server : acceptedEpoch'[i] \geq acceptedEpoch'[j]$

$$\begin{aligned}
& \text{IN } \text{acceptedEpoch}'[s] \\
& \text{ELSE } \text{recorder}["\text{maxEpoch}"]) \\
\text{RecorderSetRequests}(pc) & \triangleq (\text{"nClientRequest"} :> \\
& \text{IF } pc[1] = \text{"LeaderProcessRequest"} \text{ THEN} \\
& \quad \text{recorder}["\text{nClientRequest"}] + 1 \\
& \text{ELSE } \text{recorder}["\text{nClientRequest"}]) \\
\text{RecorderSetPc}(pc) & \triangleq (\text{"pc"} :> pc) \\
\text{RecorderSetFailure}(pc) & \triangleq \text{CASE } pc[1] = \text{"Timeout"} \rightarrow \text{RecorderIncTimeout} @@ \text{RecorderGetRestart} \\
& \quad \square \quad pc[1] = \text{"LeaderTimeout"} \rightarrow \text{RecorderIncTimeout} @@ \text{RecorderGetRestart} \\
& \quad \square \quad pc[1] = \text{"FollowerTimeout"} \rightarrow \text{RecorderIncTimeout} @@ \text{RecorderGetRestart} \\
& \quad \square \quad pc[1] = \text{"Restart"} \rightarrow \text{RecorderIncTimeout} @@ \text{RecorderIncRestart} \\
& \quad \square \quad \text{OTHER} \rightarrow \text{RecorderGetTimeout} @@ \text{RecorderGetRestart} \\
\text{UpdateRecorder}(pc) & \triangleq \text{recorder}' = \text{RecorderSetFailure}(pc) \quad @@ \text{RecorderSetTransactionNum}(pc) \\
& \quad @@ \text{RecorderSetMaxEpoch}(pc) \quad @@ \text{RecorderSetPc}(pc) \\
& \quad @@ \text{RecorderSetRequests}(pc) \quad @@ \text{recorder} \\
\text{UnchangeRecorder} & \triangleq \text{UNCHANGED recorder} \\
\text{CheckParameterHelper}(n, p, \text{Comp}(-, -)) & \triangleq \text{IF } p \in \text{DOMAIN Parameters} \\
& \quad \text{THEN } \text{Comp}(n, \text{Parameters}[p]) \\
& \quad \text{ELSE TRUE} \\
\text{CheckParameterLimit}(n, p) & \triangleq \text{CheckParameterHelper}(n, p, \text{LAMBDA } i, j : i < j) \\
\text{CheckTimeout} & \triangleq \text{CheckParameterLimit}(\text{recorder.nTimeout}, \quad \text{"MaxTimeoutFailures"}) \\
\text{CheckTransactionNum} & \triangleq \text{CheckParameterLimit}(\text{recorder.nTransaction}, \text{"MaxTransactionNum"}) \\
\text{CheckEpoch} & \triangleq \text{CheckParameterLimit}(\text{recorder.maxEpoch}, \quad \text{"MaxEpoch"}) \\
\text{CheckRestart} & \triangleq \wedge \text{CheckTimeout} \\
& \quad \wedge \text{CheckParameterLimit}(\text{recorder.nRestart}, \text{"MaxRestarts"}) \\
\text{CheckStateConstraints} & \triangleq \text{CheckTimeout} \wedge \text{CheckTransactionNum} \wedge \text{CheckEpoch} \wedge \text{CheckRestart}
\end{aligned}$$

Actions about network

$$\begin{aligned}
\text{PendingCEPOCH}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{CEPOCH} \\
\text{PendingNEWPOCH}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{NEWPOCH} \\
\text{PendingACKEPOCH}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{ACKEPOCH} \\
\text{PendingNEWLEADER}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{NEWLEADER} \\
\text{PendingACKLD}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{ACKLD} \\
\text{PendingCOMMITLD}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{COMMITLD} \\
\text{PendingPROPOSE}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{PROPOSE}
\end{aligned}$$

$PendingACK(i, j) \triangleq \wedge msgs[j][i] \neq \langle \rangle$
 $\wedge msgs[j][i][1].mtype = ACK$
 $PendingCOMMIT(i, j) \triangleq \wedge msgs[j][i] \neq \langle \rangle$
 $\wedge msgs[j][i][1].mtype = COMMIT$

Add a message to $msgs$ – add a message m to $msgs$.
 $Send(i, j, m) \triangleq msgs' = [msgs \text{ EXCEPT } ![i][j] = Append(msgs[i][j], m)]$

Remove a message from $msgs$ – discard head of $msgs$.
 $Discard(i, j) \triangleq msgs' = \text{IF } msgs[i][j] \neq \langle \rangle \text{ THEN } [msgs \text{ EXCEPT } ![i][j] = Tail(msgs[i][j])]$
 $\text{ELSE } msgs$

Combination of $Send$ and $Discard$ – discard head of $msgs[j][i]$ and add m into $msgs$.
 $Reply(i, j, m) \triangleq msgs' = [msgs \text{ EXCEPT } ![j][i] = Tail(msgs[j][i]),$
 $![i][j] = Append(msgs[i][j], m)]$

Shuffle input buffer.
 $Clean(i, j) \triangleq msgs' = [msgs \text{ EXCEPT } ![j][i] = \langle \rangle, ![i][j] = \langle \rangle]$
 $CleanInputBuffer(S) \triangleq msgs' = [s \in Server \mapsto$
 $[v \in Server \mapsto \text{IF } v \in S \text{ THEN } \langle \rangle$
 $\text{ELSE } msgs[s][v]]]$

Leader broadcasts a message *PROPOSE* to all other servers in Q .
 Note: In paper, Q is fuzzy. We think servers who leader broadcasts *NEWLEADER* to
 should receive every *PROPOSE*. So we consider *ackRecv* as Q .
 Since we let *ackRecv* = Q , there may exist some follower receiving *COMMIT* before
COMMITLD, and *zxid* in *COMMIT* later than *zxid* in *COMMITLD*. To avoid this situation,
 if $f \in ackRecv$ but $\notin ackldRecv$, f should not receive *COMMIT* until
 $f \in ackldRecv$ and receives *COMMITLD*.
 $Broadcast(i, m) \triangleq$
 $\text{LET } ackRecv_quorum \triangleq \{a \in ackRecv[i] : a.connected = \text{TRUE}\}$
 $sid_ackRecv \triangleq \{a.sid : a \in ackRecv_quorum\}$
 $\text{IN } msgs' = [msgs \text{ EXCEPT } ![i] = [v \in Server \mapsto \text{IF } \wedge v \in sid_ackRecv$
 $\wedge v \in learners[i]$
 $\wedge v \neq i$
 $\text{THEN } Append(msgs[i][v], m)$
 $\text{ELSE } msgs[i][v]]]$

Since leader decides to broadcasts message *COMMIT* when processing *ACK*, so
 we need to discard *ACK* and broadcast *COMMIT*.
 Here Q is *ackldRecv*, because we assume that f should not receive *COMMIT* until
 f receives *COMMITLD*.
 $DiscardAndBroadcast(i, j, m) \triangleq$
 $\text{LET } ackldRecv_quorum \triangleq \{a \in ackldRecv[i] : a.connected = \text{TRUE}\}$
 $sid_ackldRecv \triangleq \{a.sid : a \in ackldRecv_quorum\}$
 $\text{IN } msgs' = [msgs \text{ EXCEPT } ![j][i] = Tail(msgs[j][i]),$
 $![i] = [v \in Server \mapsto \text{IF } \wedge v \in sid_ackldRecv$
 $\wedge v \in learners[i]$
 $\wedge v \neq i$
 $\text{THEN } Append(msgs[i][v], m)$
 $\text{ELSE } msgs[i][v]]]$

Leader broadcasts *LEADERINFO* to all other servers in *cepcvRecv*.
 $DiscardAndBroadcastNEWPOCH(i, j, m) \triangleq$
 LET $new_cepcvRecv_quorum \triangleq \{c \in cpvRecv'[i] : c.connected = \text{TRUE}\}$
 $new_sid_cepcvRecv \triangleq \{c.sid : c \in new_cepcvRecv_quorum\}$
 IN $msgs' = [msgs \text{ EXCEPT } ![j][i] = Tail(msgs[j][i]),$
 $![i] = [v \in Server \mapsto \text{IF } \wedge v \in new_sid_cepcvRecv$
 $\wedge v \in learners[i]$
 $\wedge v \neq i$
 THEN $Append(msgs[i][v], m)$
 ELSE $msgs[i][v]]]$

Leader broadcasts *NEWLEADER* to all other servers in *ackRecv*.
 $DiscardAndBroadcastNEWLEADER(i, j, m) \triangleq$
 LET $new_ackRecv_quorum \triangleq \{a \in ackRecv'[i] : a.connected = \text{TRUE}\}$
 $new_sid_ackRecv \triangleq \{a.sid : a \in new_ackRecv_quorum\}$
 IN $msgs' = [msgs \text{ EXCEPT } ![j][i] = Tail(msgs[j][i]),$
 $![i] = [v \in Server \mapsto \text{IF } \wedge v \in new_sid_ackRecv$
 $\wedge v \in learners[i]$
 $\wedge v \neq i$
 THEN $Append(msgs[i][v], m)$
 ELSE $msgs[i][v]]]$

Leader broadcasts *COMMITLD* to all other servers in *ackldRecv*.
 $DiscardAndBroadcastCOMMITLD(i, j, m) \triangleq$
 LET $new_ackldRecv_quorum \triangleq \{a \in ackldRecv'[i] : a.connected = \text{TRUE}\}$
 $new_sid_ackldRecv \triangleq \{a.sid : a \in new_ackldRecv_quorum\}$
 IN $msgs' = [msgs \text{ EXCEPT } ![j][i] = Tail(msgs[j][i]),$
 $![i] = [v \in Server \mapsto \text{IF } \wedge v \in new_sid_ackldRecv$
 $\wedge v \in learners[i]$
 $\wedge v \neq i$
 THEN $Append(msgs[i][v], m)$
 ELSE $msgs[i][v]]]$

Define initial values for all variables

$InitServerVars \triangleq$ $\wedge state = [s \in Server \mapsto LOOKING]$
 $\wedge zabState = [s \in Server \mapsto ELECTION]$
 $\wedge acceptedEpoch = [s \in Server \mapsto 0]$
 $\wedge currentEpoch = [s \in Server \mapsto 0]$
 $\wedge history = [s \in Server \mapsto \langle \rangle]$
 $\wedge lastCommitted = [s \in Server \mapsto [index \mapsto 0,$
 $zxid \mapsto \langle 0, 0 \rangle]]$

$InitLeaderVars \triangleq$ $\wedge learners = [s \in Server \mapsto \{\}]$
 $\wedge cpvRecv = [s \in Server \mapsto \{\}]$
 $\wedge ackRecv = [s \in Server \mapsto \{\}]$
 $\wedge ackldRecv = [s \in Server \mapsto \{\}]$
 $\wedge sendCounter = [s \in Server \mapsto 0]$

$InitFollowerVars \triangleq leaderAddr = [s \in Server \mapsto NullPoint]$

$InitElectionVars \triangleq leaderOracle = NullPoint$

$InitMsgVars \triangleq msgs = [s \in Server \mapsto [v \in Server \mapsto \langle \rangle]]$

$InitVerifyVars \triangleq \begin{aligned} \wedge proposalMsgsLog &= \{\} \\ \wedge epochLeader &= [i \in 1 \dots MAXEPOCH \mapsto \{\}] \\ \wedge violatedInvariants &= [stateInconsistent \mapsto FALSE, \\ &\quad proposalInconsistent \mapsto FALSE, \\ &\quad commitInconsistent \mapsto FALSE, \\ &\quad ackInconsistent \mapsto FALSE, \\ &\quad messageIllegal \mapsto FALSE] \end{aligned}$

$InitRecorder \triangleq recorder = [nTimeout \mapsto 0, \\ nTransaction \mapsto 0, \\ maxEpoch \mapsto 0, \\ nRestart \mapsto 0, \\ pc \mapsto \langle "Init" \rangle, \\ nClientRequest \mapsto 0]$

$Init \triangleq \begin{aligned} &\wedge InitServerVars \\ &\wedge InitLeaderVars \\ &\wedge InitFollowerVars \\ &\wedge InitElectionVars \\ &\wedge InitVerifyVars \\ &\wedge InitMsgVars \\ &\wedge InitRecorder \end{aligned}$

Utils in state switching

$FollowerShutdown(i) \triangleq \begin{aligned} &\wedge state' = [state \text{ EXCEPT } ![i] = LOOKING] \\ &\wedge zabState' = [zabState \text{ EXCEPT } ![i] = ELECTION] \\ &\wedge leaderAddr' = [leaderAddr \text{ EXCEPT } ![i] = NullPoint] \end{aligned}$

$LeaderShutdown(i) \triangleq \begin{aligned} &\wedge LET S \triangleq learners[i] \\ &IN \quad \wedge state' = [s \in Server \mapsto IF s \in S THEN LOOKING ELSE state[s]] \\ &\quad \wedge zabState' = [s \in Server \mapsto IF s \in S THEN ELECTION ELSE zabState[s]] \\ &\quad \wedge leaderAddr' = [s \in Server \mapsto IF s \in S THEN NullPoint ELSE leaderAddr[s]] \\ &\quad \wedge CleanInputBuffer(S) \\ &\wedge learners' = [learners \text{ EXCEPT } ![i] = \{\}] \end{aligned}$

$SwitchToFollower(i) \triangleq \begin{aligned} &\wedge state' = [state \text{ EXCEPT } ![i] = FOLLOWING] \\ &\wedge zabState' = [zabState \text{ EXCEPT } ![i] = DISCOVERY] \end{aligned}$

$$\begin{aligned}
\text{SwitchToLeader}(i) &\triangleq \\
&\wedge \text{state}' = [\text{state} \text{ EXCEPT } ![i] = \text{LEADING}] \\
&\wedge \text{zabState}' = [\text{zabState} \text{ EXCEPT } ![i] = \text{DISCOVERY}] \\
&\wedge \text{learners}' = [\text{learners} \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge \text{cepochRecv}' = [\text{cepochRecv} \text{ EXCEPT } ![i] = \{[sid \mapsto i, \\
&\hspace{10em} \text{connected} \mapsto \text{TRUE}, \\
&\hspace{10em} \text{epoch} \mapsto \text{acceptedEpoch}[i]]\}] \\
&\wedge \text{ackeRecv}' = [\text{ackeRecv} \text{ EXCEPT } ![i] = \{[sid \mapsto i, \\
&\hspace{10em} \text{connected} \mapsto \text{TRUE}, \\
&\hspace{10em} \text{peerLastEpoch} \mapsto \text{currentEpoch}[i], \\
&\hspace{10em} \text{peerHistory} \mapsto \text{history}[i]]\}] \\
&\wedge \text{ackldRecv}' = [\text{ackldRecv} \text{ EXCEPT } ![i] = \{[sid \mapsto i, \\
&\hspace{10em} \text{connected} \mapsto \text{TRUE}]\}] \\
&\wedge \text{sendCounter}' = [\text{sendCounter} \text{ EXCEPT } ![i] = 0] \\
\\
\text{RemoveCepochRecv}(\text{set}, \text{sid}) &\triangleq \\
&\text{LET } \text{sid_cepochRecv} \triangleq \{s.\text{sid} : s \in \text{set}\} \\
&\text{IN IF } \text{sid} \notin \text{sid_cepochRecv} \text{ THEN } \text{set} \\
&\quad \text{ELSE LET } \text{info} \triangleq \text{CHOOSE } s \in \text{set} : s.\text{sid} = \text{sid} \\
&\quad \quad \text{new_info} \triangleq [sid \mapsto sid, \\
&\quad \quad \quad \text{connected} \mapsto \text{FALSE}, \\
&\quad \quad \quad \text{epoch} \mapsto \text{info.epoch}] \\
&\quad \text{IN } (\text{set} \setminus \{\text{info}\}) \cup \{\text{new_info}\} \\
\\
\text{RemoveAckeRecv}(\text{set}, \text{sid}) &\triangleq \\
&\text{LET } \text{sid_ackeRecv} \triangleq \{s.\text{sid} : s \in \text{set}\} \\
&\text{IN IF } \text{sid} \notin \text{sid_ackeRecv} \text{ THEN } \text{set} \\
&\quad \text{ELSE LET } \text{info} \triangleq \text{CHOOSE } s \in \text{set} : s.\text{sid} = \text{sid} \\
&\quad \quad \text{new_info} \triangleq [sid \mapsto sid, \\
&\quad \quad \quad \text{connected} \mapsto \text{FALSE}, \\
&\quad \quad \quad \text{peerLastEpoch} \mapsto \text{info.peerLastEpoch}, \\
&\quad \quad \quad \text{peerHistory} \mapsto \text{info.peerHistory}] \\
&\quad \text{IN } (\text{set} \setminus \{\text{info}\}) \cup \{\text{new_info}\} \\
\\
\text{RemoveAckldRecv}(\text{set}, \text{sid}) &\triangleq \\
&\text{LET } \text{sid_ackldRecv} \triangleq \{s.\text{sid} : s \in \text{set}\} \\
&\text{IN IF } \text{sid} \notin \text{sid_ackldRecv} \text{ THEN } \text{set} \\
&\quad \text{ELSE LET } \text{info} \triangleq \text{CHOOSE } s \in \text{set} : s.\text{sid} = \text{sid} \\
&\quad \quad \text{new_info} \triangleq [sid \mapsto sid, \\
&\quad \quad \quad \text{connected} \mapsto \text{FALSE}] \\
&\quad \text{IN } (\text{set} \setminus \{\text{info}\}) \cup \{\text{new_info}\} \\
\\
\text{RemoveLearner}(i, j) &\triangleq \\
&\wedge \text{learners}' = [\text{learners} \text{ EXCEPT } ![i] = @ \setminus \{j\}] \\
&\wedge \text{cepochRecv}' = [\text{cepochRecv} \text{ EXCEPT } ![i] = \text{RemoveCepochRecv}(@, j)] \\
&\wedge \text{ackeRecv}' = [\text{ackeRecv} \text{ EXCEPT } ![i] = \text{RemoveAckeRecv}(@, j)]
\end{aligned}$$

$$\wedge \text{ackldRecv}' = [\text{ackldRecv} \text{ EXCEPT } ![i] = \text{RemoveAckldRecv}(@, j)]$$

Actions of abnormal situations and election

$$\begin{aligned} \text{UpdateLeader}(i) &\triangleq \\ &\wedge \text{IsLooking}(i) \\ &\wedge \text{leaderOracle} \neq i \\ &\wedge \text{leaderOracle}' = i \\ &\wedge \text{SwitchToLeader}(i) \\ &\wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{currentEpoch}, \text{history}, \text{lastCommitted}, \\ &\quad \text{followerVars}, \text{verifyVars}, \text{msgVars} \rangle \\ &\wedge \text{UpdateRecorder}(\langle \text{"UpdateLeader"}, i \rangle) \end{aligned}$$

$$\begin{aligned} \text{FollowLeader}(i) &\triangleq \\ &\wedge \text{IsLooking}(i) \\ &\wedge \text{leaderOracle} \neq \text{NullPoint} \\ &\wedge \vee \wedge \text{leaderOracle} = i \\ &\quad \wedge \text{SwitchToLeader}(i) \\ &\quad \vee \wedge \text{leaderOracle} \neq i \\ &\quad \wedge \text{SwitchToFollower}(i) \\ &\quad \wedge \text{UNCHANGED } \text{leaderVars} \\ &\wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{currentEpoch}, \text{history}, \text{lastCommitted}, \\ &\quad \text{electionVars}, \text{followerVars}, \text{verifyVars}, \text{msgVars} \rangle \\ &\wedge \text{UpdateRecorder}(\langle \text{"FollowLeader"}, i \rangle) \end{aligned}$$

Follower connecting to leader fails and turns to *LOOKING*.

$$\begin{aligned} \text{FollowerTimeout}(i) &\triangleq \\ &\wedge \text{CheckTimeout} \quad \text{test restrictions of } \text{timeout_1} \\ &\wedge \text{IsFollower}(i) \\ &\wedge \text{HasNoLeader}(i) \\ &\wedge \text{FollowerShutdown}(i) \\ &\wedge \text{CleanInputBuffer}(\{i\}) \\ &\wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{currentEpoch}, \text{history}, \text{lastCommitted}, \\ &\quad \text{leaderVars}, \text{electionVars}, \text{verifyVars} \rangle \\ &\wedge \text{UpdateRecorder}(\langle \text{"FollowerTimeout"}, i \rangle) \end{aligned}$$

Leader loses support from a quorum and turns to *LOOKING*.

$$\begin{aligned} \text{LeaderTimeout}(i) &\triangleq \\ &\wedge \text{CheckTimeout} \quad \text{test restrictions of } \text{timeout_2} \\ &\wedge \text{IsLeader}(i) \\ &\wedge \neg \text{IsQuorum}(\text{learners}[i]) \\ &\wedge \text{LeaderShutdown}(i) \\ &\wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{currentEpoch}, \text{history}, \text{lastCommitted}, \\ &\quad \text{ceepochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{sendCounter}, \text{electionVars}, \\ &\quad \text{verifyVars} \rangle \\ &\wedge \text{UpdateRecorder}(\langle \text{"LeaderTimeout"}, i \rangle) \end{aligned}$$

$$Timeout(i, j) \triangleq$$
$$Restart(i) \triangleq$$

Establish connection between leader and follower.

$$ConnectAndFollowerSendCEPOCH(i, j) \triangleq$$

10

$mepoch \mapsto acceptedEpoch[j]]$ contains $f.p$
 \wedge UNCHANGED $\langle serverVars, electionVars, verifyVars, cepochRecv, ackRecv, ackldRecv, sendCounter \rangle$
 \wedge UpdateRecorder(\langle "ConnectAndFollowerSendCEPOCH", i, j \rangle)

$CepochRecvQuorumFormed(i) \triangleq$ LET $sid_cephochRecv \triangleq \{c.sid : c \in cepochRecv[i]\}$
 IN $IsQuorum(sid_cephochRecv)$

$CepochRecvBecomeQuorum(i) \triangleq$ LET $sid_cephochRecv \triangleq \{c.sid : c \in cepochRecv'[i]\}$
 IN $IsQuorum(sid_cephochRecv)$

$UpdateCepochRecv(oldSet, sid, peerEpoch) \triangleq$
 LET $sid_set \triangleq \{s.sid : s \in oldSet\}$
 IN IF $sid \in sid_set$
 THEN LET $old_info \triangleq$ CHOOSE $info \in oldSet : info.sid = sid$
 $new_info \triangleq [sid \mapsto sid,$
 $connected \mapsto TRUE,$
 $epoch \mapsto peerEpoch]$
 IN $(oldSet \setminus \{old_info\}) \cup \{new_info\}$
 ELSE LET $follower_info \triangleq [sid \mapsto sid,$
 $connected \mapsto TRUE,$
 $epoch \mapsto peerEpoch]$
 IN $oldSet \cup \{follower_info\}$

Determine new e' in this round from a quorum of *CEPOCH*.

$DetermineNewEpoch(i) \triangleq$
 LET $epoch_cephochRecv \triangleq \{c.epoch : c \in cepochRecv'[i]\}$
 IN $Maximum(epoch_cephochRecv) + 1$

Leader waits for receiving *FOLLOWERINFO* from a quorum including itself, and chooses a new epoch e' as its own epoch and broadcasts *NEWEPOCH*.

$LeaderProcessCEPOCH(i, j) \triangleq$
 \wedge CheckEpoch test restrictions of max epoch
 \wedge IsLeader(i)
 \wedge PendingCEPOCH(i, j)
 \wedge LET $msg \triangleq msgs[j][i][1]$
 $infoOk \triangleq IsMyLearner(i, j)$
 IN $\wedge infoOk$
 $\wedge \vee$ 1. has not broadcast *NEWEPOCH*
 $\wedge \neg CepochRecvQuorumFormed(i)$
 $\wedge \vee \wedge zabState[i] = DISCOVERY$
 \wedge UNCHANGED $violatedInvariants$
 $\vee \wedge zabState[i] \neq DISCOVERY$
 $\wedge PrintT("Exception: CepochRecvQuorumFormed false," \circ$
 $"while zabState not DISCOVERY.")$
 $\wedge violatedInvariants' = [violatedInvariants$
 EXCEPT $!.stateInconsistent = TRUE]$

$$\begin{aligned}
& \wedge \text{cephochRecv}' = [\text{cephochRecv} \text{ EXCEPT } ![i] = \text{UpdateCepochRecv}(@, j, \text{msg.mepoch})] \\
& \wedge \vee \begin{aligned} & \text{1.1. cepochRecv becomes quorum,} \\ & \text{then determine } e' \text{ and broadcasts NEWPOCH in } Q. \\ & \wedge \text{CepochRecvBecomeQuorum}(i) \\ & \wedge \text{acceptedEpoch}' = [\text{acceptedEpoch} \text{ EXCEPT } ![i] = \text{DetermineNewEpoch}(i)] \\ & \wedge \text{LET } m \triangleq [\text{mtype} \mapsto \text{NEWPOCH}, \\ & \quad \text{mepoch} \mapsto \text{acceptedEpoch}'[i]] \\ & \quad \text{IN } \text{DiscardAndBroadcastNEWPOCH}(i, j, m) \end{aligned} \\
& \vee \begin{aligned} & \text{1.2. cepochRecv still not quorum.} \\ & \wedge \neg \text{CepochRecvBecomeQuorum}(i) \\ & \wedge \text{Discard}(j, i) \\ & \wedge \text{UNCHANGED } \text{acceptedEpoch} \end{aligned} \\
& \vee \begin{aligned} & \text{2. has broadcast NEWPOCH} \\ & \wedge \text{CepochRecvQuorumFormed}(i) \\ & \wedge \text{cephochRecv}' = [\text{cephochRecv} \text{ EXCEPT } ![i] = \text{UpdateCepochRecv}(@, j, \text{msg.mepoch})] \\ & \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{NEWPOCH}, \\ & \quad \text{mepoch} \mapsto \text{acceptedEpoch}[i]]) \\ & \wedge \text{UNCHANGED } \langle \text{violatedInvariants}, \text{acceptedEpoch} \rangle \end{aligned} \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{zabState}, \text{currentEpoch}, \text{history}, \text{lastCommitted}, \text{learners}, \\
& \quad \text{ackRecv}, \text{ackldRecv}, \text{sendCounter}, \text{followerVars}, \\
& \quad \text{electionVars}, \text{proposalMsgsLog}, \text{epochLeader} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"LeaderProcessCEPOCH"}, i, j \rangle)
\end{aligned}$$

Follower receives *LEADERINFO*. If $\text{newEpoch} \geq \text{acceptedEpoch}$, then follower updates *acceptedEpoch* and sends *ACEPOCH* back, containing *currentEpoch* and history. After this, *zabState* turns to *SYNC*.

$$\begin{aligned}
& \text{FollowerProcessNEWPOCH}(i, j) \triangleq \\
& \wedge \text{IsFollower}(i) \\
& \wedge \text{PendingNEWPOCH}(i, j) \\
& \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \text{infoOk} \triangleq \text{IsMyLeader}(i, j) \\
& \quad \text{stateOk} \triangleq \text{zabState}[i] = \text{DISCOVERY} \\
& \quad \text{epochOk} \triangleq \text{msg.mepoch} \geq \text{acceptedEpoch}[i] \\
& \quad \text{IN } \wedge \text{infoOk} \\
& \quad \wedge \vee \begin{aligned} & \text{1. Normal case} \\ & \wedge \text{epochOk} \\ & \wedge \vee \wedge \text{stateOk} \\ & \quad \wedge \text{acceptedEpoch}' = [\text{acceptedEpoch} \text{ EXCEPT } ![i] = \text{msg.mepoch}] \\ & \quad \wedge \text{LET } m \triangleq [\text{mtype} \mapsto \text{ACEPOCH}, \\ & \quad \quad \text{mepoch} \mapsto \text{currentEpoch}[i], \\ & \quad \quad \text{mhistory} \mapsto \text{history}[i]] \\ & \quad \quad \text{IN } \text{Reply}(i, j, m) \\ & \quad \wedge \text{zabState}' = [\text{zabState} \text{ EXCEPT } ![i] = \text{SYNCHRONIZATION}] \\ & \quad \wedge \text{UNCHANGED } \text{violatedInvariants} \end{aligned} \\
& \quad \vee \wedge \neg \text{stateOk}
\end{aligned}$$

```

    ∧ PrintT("Exception: Follower receives NEWEPOCH," ∘
      " whileZabState not DISCOVERY.")
    ∧ violatedInvariants' = [violatedInvariants
      EXCEPT !.stateInconsistent = TRUE]
    ∧ Discard(j, i)
    ∧ UNCHANGED ⟨acceptedEpoch, zabState⟩
    ∧ UNCHANGED ⟨followerVars, learners, cepochRecv, ackeRecv,
      ackldRecv, state⟩
  ∨ 2. Abnormal case - go back to election
    ∧ ¬epochOk
    ∧ FollowerShutdown(i)
    ∧ LET leader ≜ leaderAddr[i]
      IN  ∧ Clean(i, leader)
          ∧ RemoveLearner(leader, i)
          ∧ UNCHANGED ⟨acceptedEpoch, violatedInvariants⟩
    ∧ UNCHANGED ⟨currentEpoch, history, lastCommitted, sendCounter,
      electionVars, proposalMsgsLog, epochLeader⟩
    ∧ UpdateRecorder(("FollowerProcessNEWEPOCH", i, j))

AckeRecvQuorumFormed(i) ≜ LET sid_akeRecv ≜ {a.sid : a ∈ ackeRecv[i]}
  IN  IsQuorum(sid_akeRecv)
AckeRecvBecomeQuorum(i) ≜ LET sid_akeRecv ≜ {a.sid : a ∈ ackeRecv'[i]}
  IN  IsQuorum(sid_akeRecv)

UpdateAckeRecv(oldSet, sid, peerEpoch, peerHistory) ≜
  LET sid_set ≜ {s.sid : s ∈ oldSet}
    follower_info ≜ [sid      ↦ sid,
                     connected ↦ TRUE,
                     peerLastEpoch ↦ peerEpoch,
                     peerHistory  ↦ peerHistory]
  IN  IF sid ∈ sid_set
      THEN LET old_info ≜ CHOOSE info ∈ oldSet : info.sid = sid
          IN  (oldSet \ {old_info}) ∪ {follower_info}
      ELSE oldSet ∪ {follower_info}

for checking invariants
RECURSIVE SetPacketsForChecking(−, −, −, −, −, −)
SetPacketsForChecking(set, src, ep, his, cur, end) ≜
  IF cur > end THEN set
  ELSE LET m_proposal ≜ [source ↦ src,
                        epoch  ↦ ep,
                        zxid   ↦ his[cur].zxid,
                        data    ↦ his[cur].value]
    IN  SetPacketsForChecking((set ∪ {m_proposal}), src, ep, his, cur + 1, end)

```

$LastZxidOfHistory(his) \triangleq$ IF $Len(his) = 0$ THEN $\langle 0, 0 \rangle$
ELSE $his[Len(his)].zxid$

TRUE: $f1.a > f2.a$ or $(f1.a = f2.a \text{ and } f1.zxid \geq f2.zxid)$
 $MoreResentOrEqual(ss1, ss2) \triangleq$ $\vee ss1.currentEpoch > ss2.currentEpoch$
 $\vee \wedge ss1.currentEpoch = ss2.currentEpoch$
 $\wedge \neg ZxidCompare(ss2.lastZxid, ss1.lastZxid)$

Determine initial history Ie' in this round from a quorum of *ACKEPOCH*.
 $DetermineInitialHistory(i) \triangleq$
LET $set \triangleq ackeRecv'[i]$
 $ss_set \triangleq \{ [sid \mapsto a.sid,$
 $currentEpoch \mapsto a.peerLastEpoch,$
 $lastZxid \mapsto LastZxidOfHistory(a.peerHistory)]$
 $: a \in set \}$
 $selected \triangleq$ CHOOSE $ss \in ss_set :$
 $\forall ss1 \in (ss_set \setminus \{ss\}) : MoreResentOrEqual(ss, ss1)$
 $info \triangleq$ CHOOSE $f \in set : f.sid = selected.sid$
IN $info.peerHistory$

RECURSIVE $InitAcksidHelper(-, -)$
 $InitAcksidHelper(txns, src) \triangleq$ IF $Len(txns) = 0$ THEN $\langle \rangle$
ELSE LET $oldTxn \triangleq txns[1]$
 $newTxn \triangleq [zxid \mapsto oldTxn.zxid,$
 $value \mapsto oldTxn.value,$
 $ackSid \mapsto \{src\},$
 $epoch \mapsto oldTxn.epoch]$
IN $\langle newTxn \rangle \circ InitAcksidHelper(Tail(txns), src)$

Atomically let all $txns$ in initial history contain self's acks.
 $InitAcksid(i, his) \triangleq InitAcksidHelper(his, i)$

Leader waits for receiving *ACKEPOCH* from a quorum, and determines *initialHistory* according to history of whom has most recent state summary from them. After this, leader's *zabState* turns to *SYNCHRONIZATION*.

$LeaderProcessACKEPOCH(i, j) \triangleq$
 $\wedge IsLeader(i)$
 $\wedge PendingACKEPOCH(i, j)$
 \wedge LET $msg \triangleq msgs[j][i][1]$
 $infoOk \triangleq IsMyLearner(i, j)$
IN $\wedge infoOk$
 $\wedge \vee$ 1. has broadcast *NEWLEADER*
 $\wedge AckeRecvQuorumFormed(i)$
 $\wedge ackeRecv' = [ackeRecv \text{ EXCEPT } ![i] = UpdateAckeRecv(@, j,$
 $msg.mepoch, msg.mhistory)]$
 \wedge LET $toSend \triangleq history[i]$ contains (Ie', Be')
 $m \triangleq [mtype \mapsto NEWLEADER,$

$$\begin{aligned}
& \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED} \langle \text{currentEpoch}, \text{history}, \text{zabState}, \\
& \quad \text{proposalMsgsLog}, \text{epochLeader} \rangle \\
& \wedge \text{UNCHANGED} \langle \text{state}, \text{acceptedEpoch}, \text{lastCommitted}, \text{learners}, \text{ceepochRecv}, \text{ackldRecv}, \\
& \quad \text{sendCounter}, \text{followerVars}, \text{electionVars} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"LeaderProcessACKEPOCH"}, i, j \rangle) \\
\hline
& \text{Follower receives NEWLEADER. Update } f.a \text{ and history.} \\
& \text{FollowerProcessNEWLEADER}(i, j) \triangleq \\
& \quad \wedge \text{IsFollower}(i) \\
& \quad \wedge \text{PendingNEWLEADER}(i, j) \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \quad \text{infoOk} \triangleq \text{IsMyLeader}(i, j) \\
& \quad \quad \text{epochOk} \triangleq \text{acceptedEpoch}[i] = \text{msg.mepoch} \\
& \quad \quad \text{stateOk} \triangleq \text{zabState}[i] = \text{SYNCHRONIZATION} \\
& \quad \text{IN } \wedge \text{infoOk} \\
& \quad \quad \wedge \vee \quad \text{1. } f.p \text{ not equals } e', \text{ starts a new iteration.} \\
& \quad \quad \quad \wedge \neg \text{epochOk} \\
& \quad \quad \quad \wedge \text{FollowerShutdown}(i) \\
& \quad \quad \quad \wedge \text{LET } \text{leader} \triangleq \text{leaderAddr}[i] \\
& \quad \quad \quad \quad \text{IN } \wedge \text{Clean}(i, \text{leader}) \\
& \quad \quad \quad \quad \quad \wedge \text{RemoveLearner}(\text{leader}, i) \\
& \quad \quad \quad \wedge \text{UNCHANGED} \langle \text{violatedInvariants}, \text{currentEpoch}, \text{history} \rangle \\
& \quad \quad \vee \quad \text{2. } f.p \text{ equals } e'. \\
& \quad \quad \quad \wedge \text{epochOk} \\
& \quad \quad \quad \wedge \vee \wedge \text{stateOk} \\
& \quad \quad \quad \quad \wedge \text{UNCHANGED } \text{violatedInvariants} \\
& \quad \quad \quad \vee \wedge \neg \text{stateOk} \\
& \quad \quad \quad \quad \wedge \text{PrintT}(\text{"Exception: Follower receives NEWLEADER,"} \circ \\
& \quad \quad \quad \quad \quad \text{" whileZabState not SYNCHRONIZATION."}) \\
& \quad \quad \quad \quad \wedge \text{violatedInvariants}' = [\text{violatedInvariants} \\
& \quad \quad \quad \quad \quad \text{EXCEPT !.stateInconsistent} = \text{TRUE}] \\
& \quad \quad \quad \quad \wedge \text{currentEpoch}' = [\text{currentEpoch} \text{ EXCEPT } ![i] = \text{acceptedEpoch}[i]] \\
& \quad \quad \quad \quad \wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{msg.mhistory}] \quad \text{no need to care ackSid} \\
& \quad \quad \quad \quad \wedge \text{LET } m \triangleq [\text{mtype} \mapsto \text{ACKLD}, \\
& \quad \quad \quad \quad \quad \text{mzxid} \mapsto \text{LastZxidOfHistory}(\text{history}'[i])] \\
& \quad \quad \quad \quad \text{IN } \text{Reply}(i, j, m) \\
& \quad \quad \quad \wedge \text{UNCHANGED} \langle \text{followerVars}, \text{state}, \text{zabState}, \text{learners}, \text{ceepochRecv}, \\
& \quad \quad \quad \quad \quad \text{ackeRecv}, \text{ackldRecv} \rangle \\
& \quad \quad \wedge \text{UNCHANGED} \langle \text{acceptedEpoch}, \text{lastCommitted}, \text{sendCounter}, \text{electionVars}, \\
& \quad \quad \quad \text{proposalMsgsLog}, \text{epochLeader} \rangle \\
& \quad \quad \wedge \text{UpdateRecorder}(\langle \text{"FollowerProcessNEWLEADER"}, i, j \rangle) \\
& \text{AckldRecvQuorumFormed}(i) \triangleq \text{LET } \text{sid_ackldRecv} \triangleq \{a.\text{sid} : a \in \text{ackldRecv}[i]\} \\
& \quad \text{IN } \text{IsQuorum}(\text{sid_ackldRecv})
\end{aligned}$$

$AckldRecvBecomeQuorum(i) \triangleq$ LET $sid_ackldRecv \triangleq \{a.sid : a \in ackldRecv'[i]\}$
IN $IsQuorum(sid_ackldRecv)$

$UpdateAckldRecv(oldSet, sid) \triangleq$
LET $sid_set \triangleq \{s.sid : s \in oldSet\}$
 $follower_info \triangleq [sid \mapsto sid,$
 $connected \mapsto \text{TRUE}]$
IN IF $sid \in sid_set$
THEN LET $old_info \triangleq$ CHOOSE $info \in oldSet : info.sid = sid$
IN $(oldSet \setminus \{old_info\}) \cup \{follower_info\}$
ELSE $oldSet \cup \{follower_info\}$

$LastZxid(i) \triangleq LastZxidOfHistory(history[i])$

RECURSIVE $UpdateAcksidHelper(-, -, -)$
 $UpdateAcksidHelper(txns, target, endZxid) \triangleq$
IF $Len(txns) = 0$ THEN $\langle \rangle$
ELSE LET $oldTxn \triangleq txns[1]$
IN IF $ZxidCompare(oldTxn.zxid, endZxid)$ THEN $txns$
ELSE LET $newTxn \triangleq [zxid \mapsto oldTxn.zxid,$
 $value \mapsto oldTxn.value,$
 $ackSid \mapsto \text{IF } target \in oldTxn.ackSid$
 $\text{THEN } oldTxn.ackSid$
 $\text{ELSE } oldTxn.ackSid \cup \{target\},$
 $epoch \mapsto oldTxn.epoch]$
IN $\langle newTxn \rangle \circ UpdateAcksidHelper(Tail(txns), target, endZxid)$

Atomically add $ackSid$ of one learner according to $zxid$ in $ACKLD$.

$UpdateAcksid(his, target, endZxid) \triangleq UpdateAcksidHelper(his, target, endZxid)$

Leader waits for receiving $ACKLD$ from a quorum including itself, and broadcasts $COMMITLD$ and turns to $BROADCAST$.

$LeaderProcessACKLD(i, j) \triangleq$
 $\wedge IsLeader(i)$
 $\wedge PendingACKLD(i, j)$
 $\wedge \text{LET } msg \triangleq msgs[j][i][1]$
 $infoOk \triangleq IsMyLearner(i, j)$
IN $\wedge infoOk$
 $\wedge \vee$ 1. has not broadcast $COMMITLD$
 $\wedge \neg AckldRecvQuorumFormed(i)$
 $\wedge \vee \wedge zabState[i] = SYNCHRONIZATION$
 $\wedge \text{UNCHANGED } violatedInvariants$
 $\vee \wedge zabState[i] \neq SYNCHRONIZATION$
 $\wedge PrintT(\text{"Exception: AckldRecvQuorumFormed false,"} \circ$
 $\text{" while zabState not SYNCHRONIZATION."})$
 $\wedge violatedInvariants' = [violatedInvariants$

```

EXCEPT !.stateInconsistent = TRUE]
 $\wedge$   $ackldRecv' = [ackldRecv \text{ EXCEPT } ![i] = UpdateAckldRecv(@, j)]$ 
 $\wedge$   $history' = [history \text{ EXCEPT } ![i] = UpdateAcksid(@, j, msg.mzxid)]$ 
 $\wedge \vee$  1.1.  $ackldRecv$  becomes quorum,
    then broadcasts COMMITLD and turns to BROADCAST.
 $\wedge AckldRecvBecomeQuorum(i)$ 
 $\wedge lastCommitted' = [lastCommitted \text{ EXCEPT }$ 
     $![i] = [index \mapsto Len(history[i]),$ 
     $zxid \mapsto LastZxid(i)]]$ 
 $\wedge zabState' = [zabState \text{ EXCEPT } ![i] = BROADCAST]$ 
 $\wedge \text{LET } m \triangleq [mtype \mapsto COMMITLD,$ 
     $mzxid \mapsto LastZxid(i)]$ 
    IN DiscardAndBroadcastCOMMITLD( $i, j, m$ )
 $\vee$  1.2.  $ackldRecv$  still not quorum.
 $\wedge \neg AckldRecvBecomeQuorum(i)$ 
 $\wedge Discard(j, i)$ 
 $\wedge \text{UNCHANGED } \langle zabState, lastCommitted \rangle$ 
 $\vee$  2. has broadcast COMMITLD
 $\wedge AckldRecvQuorumFormed(i)$ 
 $\wedge \vee \wedge zabState[i] = BROADCAST$ 
 $\wedge \text{UNCHANGED } violatedInvariants$ 
 $\vee \wedge zabState[i] \neq BROADCAST$ 
 $\wedge PrintT(\text{"Exception: AckldRecvQuorumFormed true,"} \circ$ 
     $\text{"while zabState not BROADCAST."})$ 
 $\wedge violatedInvariants' = [violatedInvariants$ 
     $\text{EXCEPT !.stateInconsistent = TRUE}]$ 
 $\wedge ackldRecv' = [ackldRecv \text{ EXCEPT } ![i] = UpdateAckldRecv(@, j)]$ 
 $\wedge history' = [history \text{ EXCEPT } ![i] = UpdateAcksid(@, j, msg.mzxid)]$ 
 $\wedge Reply(i, j, [mtype \mapsto COMMITLD,$ 
     $mzxid \mapsto lastCommitted[i].zxid])$ 
 $\wedge \text{UNCHANGED } \langle zabState, lastCommitted \rangle$ 
 $\wedge \text{UNCHANGED } \langle state, acceptedEpoch, currentEpoch, learners, cepochRecv, ackeRecv,$ 
     $sendCounter, followerVars, electionVars, proposalMsgsLog, epochLeader \rangle$ 
 $\wedge UpdateRecorder(\langle \text{"LeaderProcessACKLD"}, i, j \rangle)$ 

RECURSIVE ZxidToIndexHepler( $-, -, -, -$ )
ZxidToIndexHepler( $his, zxid, cur, appeared$ )  $\triangleq$ 
    IF  $cur > Len(his)$  THEN  $cur$ 
    ELSE IF TxnZxidEqual( $his[cur], zxid$ )
        THEN CASE  $appeared = \text{TRUE} \rightarrow -1$ 
             $\square$  OTHER  $\rightarrow Minimum(\{cur,$ 
                 $ZxidToIndexHepler(his, zxid, cur + 1, \text{TRUE})\})$ 
        ELSE ZxidToIndexHepler( $his, zxid, cur + 1, appeared$ )

return  $-1$ : this  $zxid$  appears at least twice.  $Len(his) + 1$ : does not exist.

```

$1 - \text{Len}(\text{his})$: exists and appears just once.
 $\text{ZxidToIndex}(\text{his}, \text{zxid}) \triangleq \text{IF } \text{ZxidEqual}(\text{zxid}, \langle 0, 0 \rangle) \text{ THEN } 0$
 $\quad \text{ELSE IF } \text{Len}(\text{his}) = 0 \text{ THEN } 1$
 $\quad \text{ELSE LET } \text{len} \triangleq \text{Len}(\text{his}) \text{ IN}$
 $\quad \quad \text{IF } \exists \text{idx} \in 1 \dots \text{len} : \text{TxnZxidEqual}(\text{his}[\text{idx}], \text{zxid})$
 $\quad \quad \text{THEN } \text{ZxidToIndexHepler}(\text{his}, \text{zxid}, 1, \text{FALSE})$
 $\quad \quad \text{ELSE } \text{len} + 1$

Follower receives COMMITLD. Commit all txns.
 $\text{FollowerProcessCOMMITLD}(i, j) \triangleq$
 $\quad \wedge \text{IsFollower}(i)$
 $\quad \wedge \text{PendingCOMMITLD}(i, j)$
 $\quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1]$
 $\quad \quad \text{infoOk} \triangleq \text{IsMyLeader}(i, j)$
 $\quad \quad \text{index} \triangleq \text{IF } \text{ZxidEqual}(\text{msg.mzxid}, \langle 0, 0 \rangle) \text{ THEN } 0$
 $\quad \quad \quad \text{ELSE } \text{ZxidToIndex}(\text{history}[i], \text{msg.mzxid})$
 $\quad \quad \text{logOk} \triangleq \text{index} \geq 0 \wedge \text{index} \leq \text{Len}(\text{history}[i])$
 $\quad \text{IN } \wedge \text{infoOk}$
 $\quad \quad \wedge \wedge \text{logOk}$
 $\quad \quad \quad \wedge \text{UNCHANGED } \text{violatedInvariants}$
 $\quad \quad \quad \vee \wedge \neg \text{logOk}$
 $\quad \quad \quad \wedge \text{PrintT}(\text{"Exception: zxid in COMMITLD not exists in history."})$
 $\quad \quad \quad \wedge \text{violatedInvariants}' = [\text{violatedInvariants}$
 $\quad \quad \quad \quad \text{EXCEPT !.proposalInconsistent} = \text{TRUE}]$
 $\quad \quad \quad \wedge \text{lastCommitted}' = [\text{lastCommitted} \text{ EXCEPT !}[i] = [\text{index} \mapsto \text{index},$
 $\quad \quad \quad \quad \quad \quad \quad \text{zxid} \mapsto \text{msg.mzxid}]]$
 $\quad \quad \quad \wedge \text{zabState}' = [\text{zabState} \text{ EXCEPT !}[i] = \text{BROADCAST}]$
 $\quad \quad \quad \wedge \text{Discard}(j, i)$
 $\quad \wedge \text{UNCHANGED } \langle \text{state}, \text{acceptedEpoch}, \text{currentEpoch}, \text{history}, \text{leaderVars},$
 $\quad \quad \text{followerVars}, \text{electionVars}, \text{proposalMsgsLog}, \text{epochLeader} \rangle$
 $\quad \wedge \text{UpdateRecorder}(\langle \text{"FollowerProcessCOMMITLD"}, i, j \rangle)$

$\text{IncZxid}(s, \text{zxid}) \triangleq \text{IF } \text{currentEpoch}[s] = \text{zxid}[1] \text{ THEN } \langle \text{zxid}[1], \text{zxid}[2] + 1 \rangle$
 $\quad \text{ELSE } \langle \text{currentEpoch}[s], 1 \rangle$

Leader receives client request. Note: In production, any server in traffic can receive requests and

forward it to leader if necessary. We choose to let leader be the sole one who can receive write requests, to simplify spec and keep correctness at the same time.

$\text{LeaderProcessRequest}(i) \triangleq$
 $\quad \wedge \text{CheckTransactionNum}$ test restrictions of transaction num
 $\quad \wedge \text{IsLeader}(i)$
 $\quad \wedge \text{zabState}[i] = \text{BROADCAST}$
 $\quad \wedge \text{LET } \text{request_value} \triangleq \text{GetRecorder}(\text{"nClientRequest"})$ unique value
 $\quad \quad \text{newTxn} \triangleq [\text{zxid} \mapsto \text{IncZxid}(i, \text{LastZxid}(i)),$
 $\quad \quad \quad \text{value} \mapsto \text{request_value},$

$$\begin{aligned}
& \text{ackSid} \mapsto \{i\}, \\
& \text{epoch} \mapsto \text{currentEpoch}[i] \\
& \text{IN } \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{Append}(@, \text{newTxn})] \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{zabState}, \text{acceptedEpoch}, \text{currentEpoch}, \text{lastCommitted}, \\
& \quad \text{leaderVars}, \text{followerVars}, \text{electionVars}, \text{msgVars}, \text{verifyVars} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"LeaderProcessRequest"}, i \rangle) \\
\\
& \text{Latest counter existing in history.} \\
& \text{CurrentCounter}(i) \triangleq \text{IF } \text{LastZxid}(i)[1] = \text{currentEpoch}[i] \text{ THEN } \text{LastZxid}(i)[2] \\
& \quad \text{ELSE } 0 \\
\\
& \text{Leader broadcasts PROPOSE when sendCounter < currentCounter.} \\
& \text{LeaderBroadcastPROPOSE}(i) \triangleq \\
& \quad \wedge \text{IsLeader}(i) \\
& \quad \wedge \text{zabState}[i] = \text{BROADCAST} \\
& \quad \wedge \text{sendCounter}[i] < \text{CurrentCounter}(i) \text{ there exists proposal to be sent} \\
& \quad \wedge \text{LET } \text{toSendCounter} \triangleq \text{sendCounter}[i] + 1 \\
& \quad \quad \text{toSendZxid} \triangleq \langle \text{currentEpoch}[i], \text{toSendCounter} \rangle \\
& \quad \quad \text{toSendIndex} \triangleq \text{ZxidToIndex}(\text{history}[i], \text{toSendZxid}) \\
& \quad \quad \text{toSendTxn} \triangleq \text{history}[i][\text{toSendIndex}] \\
& \quad \quad \text{m_proposal} \triangleq [\text{mtype} \mapsto \text{PROPOSE}, \\
& \quad \quad \quad \text{mzxid} \mapsto \text{toSendTxn.zxid}, \\
& \quad \quad \quad \text{mdata} \mapsto \text{toSendTxn.value}] \\
& \quad \quad \text{m_proposal_forChecking} \triangleq [\text{source} \mapsto i, \\
& \quad \quad \quad \text{epoch} \mapsto \text{currentEpoch}[i], \\
& \quad \quad \quad \text{zxid} \mapsto \text{toSendTxn.zxid}, \\
& \quad \quad \quad \text{data} \mapsto \text{toSendTxn.value}] \\
& \quad \text{IN } \wedge \text{sendCounter}' = [\text{sendCounter} \text{ EXCEPT } ![i] = \text{toSendCounter}] \\
& \quad \quad \wedge \text{Broadcast}(i, \text{m_proposal}) \\
& \quad \quad \wedge \text{proposalMsgsLog}' = \text{proposalMsgsLog} \cup \{\text{m_proposal_forChecking}\} \\
& \quad \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{learners}, \text{cepocheRecv}, \text{ackeRecv}, \text{ackldRecv}, \\
& \quad \quad \text{followerVars}, \text{electionVars}, \text{epochLeader}, \text{violatedInvariants} \rangle \\
& \quad \wedge \text{UpdateRecorder}(\langle \text{"LeaderBroadcastPROPOSE"}, i \rangle) \\
\\
& \text{IsNextZxid}(\text{curZxid}, \text{nextZxid}) \triangleq \\
& \quad \vee \text{first PROPOSAL in this epoch} \\
& \quad \quad \wedge \text{nextZxid}[2] = 1 \\
& \quad \quad \wedge \text{curZxid}[1] < \text{nextZxid}[1] \\
& \quad \vee \text{not first PROPOSAL in this epoch} \\
& \quad \quad \wedge \text{nextZxid}[2] > 1 \\
& \quad \quad \wedge \text{curZxid}[1] = \text{nextZxid}[1] \\
& \quad \quad \wedge \text{curZxid}[2] + 1 = \text{nextZxid}[2] \\
\\
& \text{Follower processes PROPOSE, saves it in history and replies ACK.} \\
& \text{FollowerProcessPROPOSE}(i, j) \triangleq \\
& \quad \wedge \text{IsFollower}(i)
\end{aligned}$$

$$\begin{aligned}
& \vee \text{hasCommitted} \\
& \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED } \langle \text{violatedInvariants}, \text{lastCommitted} \rangle \\
& \vee \wedge \text{outstanding} \\
& \wedge \neg \text{hasCommitted} \\
& \wedge \text{LeaderTryToCommit}(i, \text{index}, \text{msg.mzxid}, \text{txnAfterAddAck}, j) \\
\vee \wedge \vee \neg \text{exist} \\
& \vee \neg \text{monotonicallyInc} \\
& \wedge \text{PrintT}(\text{"Exception: No such zxid."} \circ \\
& \quad \text{" / ackIndex doesn't inc monotonically."}) \\
& \wedge \text{violatedInvariants}' = [\text{violatedInvariants} \\
& \quad \text{EXCEPT !.ackInconsistent} = \text{TRUE}] \\
& \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED } \langle \text{history}, \text{lastCommitted} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{zabState}, \text{acceptedEpoch}, \text{currentEpoch}, \text{leaderVars}, \\
& \quad \text{followerVars}, \text{electionVars}, \text{proposalMsgsLog}, \text{epochLeader} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"LeaderProcessACK"}, i, j \rangle)
\end{aligned}$$

Follower processes *COMMIT*.

$$\begin{aligned}
& \text{FollowerProcessCOMMIT}(i, j) \triangleq \\
& \wedge \text{IsFollower}(i) \\
& \wedge \text{PendingCOMMIT}(i, j) \\
& \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \text{infoOk} \triangleq \text{IsMyLeader}(i, j) \\
& \quad \text{pending} \triangleq \text{lastCommitted}[i].\text{index} < \text{Len}(\text{history}[i]) \\
& \text{IN} \wedge \text{infoOk} \\
& \wedge \vee \wedge \neg \text{pending} \\
& \quad \wedge \text{PrintT}(\text{"Warn: Committing zxid without seeing txn."}) \\
& \quad \wedge \text{UNCHANGED } \langle \text{lastCommitted}, \text{violatedInvariants} \rangle \\
& \vee \wedge \text{pending} \\
& \quad \wedge \text{LET } \text{firstElement} \triangleq \text{history}[i][\text{lastCommitted}[i].\text{index} + 1] \\
& \quad \quad \text{match} \triangleq \text{ZxidEqual}(\text{firstElement.zxid}, \text{msg.mzxid}) \\
& \text{IN} \\
& \quad \vee \wedge \neg \text{match} \\
& \quad \quad \wedge \text{PrintT}(\text{"Exception: Committing zxid not equals"} \circ \\
& \quad \quad \quad \text{" next pending txn zxid."}) \\
& \quad \quad \wedge \text{violatedInvariants}' = [\text{violatedInvariants} \text{ EXCEPT } \\
& \quad \quad \quad \text{!.commitInconsistent} = \text{TRUE}] \\
& \quad \quad \wedge \text{UNCHANGED } \text{lastCommitted} \\
& \vee \wedge \text{match} \\
& \quad \wedge \text{lastCommitted}' = [\text{lastCommitted} \text{ EXCEPT } ![i] = \\
& \quad \quad \quad [\text{index} \mapsto \text{lastCommitted}[i].\text{index} + 1, \\
& \quad \quad \quad \text{zxid} \mapsto \text{firstElement.zxid}]] \\
& \quad \wedge \text{UNCHANGED } \text{violatedInvariants} \\
& \wedge \text{Discard}(j, i)
\end{aligned}$$

$$\wedge \text{UNCHANGED } \langle \text{state}, \text{zabState}, \text{acceptedEpoch}, \text{currentEpoch}, \text{history}, \\ \text{leaderVars}, \text{followerVars}, \text{electionVars}, \text{proposalMsgsLog}, \text{epochLeader} \rangle \\ \wedge \text{UpdateRecorder}(\langle \text{"FollowerProcessCOMMIT"}, i, j \rangle)$$

Used to discard some messages which should not exist in network channel. This action should not be triggered.

$$\begin{aligned} \text{FilterNonexistentMessage}(i) &\triangleq \\ &\wedge \exists j \in \text{Server} \setminus \{i\} : \wedge \text{msgs}[j][i] \neq \langle \rangle \\ &\quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\ &\quad \text{IN} \\ &\quad \vee \wedge \text{IsLeader}(i) \\ &\quad \quad \wedge \text{LET } \text{infoOk} \triangleq \text{IsMyLearner}(i, j) \\ &\quad \quad \text{IN} \\ &\quad \quad \vee \text{msg.mtype} = \text{NEWPOCH} \\ &\quad \quad \vee \text{msg.mtype} = \text{NEWLEADER} \\ &\quad \quad \vee \text{msg.mtype} = \text{COMMITLD} \\ &\quad \quad \vee \text{msg.mtype} = \text{PROPOSE} \\ &\quad \quad \vee \text{msg.mtype} = \text{COMMIT} \\ &\quad \quad \vee \wedge \neg \text{infoOk} \\ &\quad \quad \quad \wedge \vee \text{msg.mtype} = \text{CEPOCH} \\ &\quad \quad \quad \vee \text{msg.mtype} = \text{ACKEPOCH} \\ &\quad \quad \quad \vee \text{msg.mtype} = \text{ACKLD} \\ &\quad \quad \quad \vee \text{msg.mtype} = \text{ACK} \\ &\quad \vee \wedge \text{IsFollower}(i) \\ &\quad \quad \wedge \text{LET } \text{infoOk} \triangleq \text{IsMyLeader}(i, j) \\ &\quad \quad \text{IN} \\ &\quad \quad \vee \text{msg.mtype} = \text{CEPOCH} \\ &\quad \quad \vee \text{msg.mtype} = \text{ACKEPOCH} \\ &\quad \quad \vee \text{msg.mtype} = \text{ACKLD} \\ &\quad \quad \vee \text{msg.mtype} = \text{ACK} \\ &\quad \quad \vee \wedge \neg \text{infoOk} \\ &\quad \quad \quad \wedge \vee \text{msg.mtype} = \text{NEWPOCH} \\ &\quad \quad \quad \vee \text{msg.mtype} = \text{NEWLEADER} \\ &\quad \quad \quad \vee \text{msg.mtype} = \text{COMMITLD} \\ &\quad \quad \quad \vee \text{msg.mtype} = \text{PROPOSE} \\ &\quad \quad \quad \vee \text{msg.mtype} = \text{COMMIT} \\ &\quad \vee \text{IsLooking}(i) \\ &\quad \wedge \text{Discard}(j, i) \\ &\wedge \text{violatedInvariants}' = [\text{violatedInvariants} \text{ EXCEPT } !.\text{messageIllegal} = \text{TRUE}] \\ &\wedge \text{UNCHANGED } \langle \text{serverVars}, \text{leaderVars}, \text{followerVars}, \text{electionVars}, \\ &\quad \text{proposalMsgsLog}, \text{epochLeader} \rangle \\ &\wedge \text{UnchangeRecorder} \end{aligned}$$

Defines how the variables may transition.

$\text{Next} \triangleq$

Election
 $\vee \exists i \in \text{Server} : \text{UpdateLeader}(i)$
 $\vee \exists i \in \text{Server} : \text{FollowLeader}(i)$
Abnormal situations like failure, network disconnection
 $\vee \exists i \in \text{Server} : \text{FollowerTimeout}(i)$
 $\vee \exists i \in \text{Server} : \text{LeaderTimeout}(i)$
 $\vee \exists i, j \in \text{Server} : \text{Timeout}(i, j)$
 $\vee \exists i \in \text{Server} : \text{Restart}(i)$
Zab module - Discovery and Synchronization part
 $\vee \exists i, j \in \text{Server} : \text{ConnectAndFollowerSendCEPOCH}(i, j)$
 $\vee \exists i, j \in \text{Server} : \text{LeaderProcessCEPOCH}(i, j)$
 $\vee \exists i, j \in \text{Server} : \text{FollowerProcessNEWCEPOCH}(i, j)$
 $\vee \exists i, j \in \text{Server} : \text{LeaderProcessACKEPOCH}(i, j)$
 $\vee \exists i, j \in \text{Server} : \text{FollowerProcessNEWLEADER}(i, j)$
 $\vee \exists i, j \in \text{Server} : \text{LeaderProcessACKLD}(i, j)$
 $\vee \exists i, j \in \text{Server} : \text{FollowerProcessCOMMITLD}(i, j)$
Zab module – Broadcast part
 $\vee \exists i \in \text{Server} : \text{LeaderProcessRequest}(i)$
 $\vee \exists i \in \text{Server} : \text{LeaderBroadcastPROPOSE}(i)$
 $\vee \exists i, j \in \text{Server} : \text{FollowerProcessPROPOSE}(i, j)$
 $\vee \exists i, j \in \text{Server} : \text{LeaderProcessACK}(i, j)$
 $\vee \exists i, j \in \text{Server} : \text{FollowerProcessCOMMIT}(i, j)$
An action used to judge whether there are redundant messages in network
 $\vee \exists i \in \text{Server} : \text{FilterNonexistentMessage}(i)$

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

Define safety properties of Zab.

$\text{ShouldNotBeTriggered} \triangleq \forall p \in \text{DOMAIN } \text{violatedInvariants} : \text{violatedInvariants}[p] = \text{FALSE}$

There is most one established leader for a certain epoch.

$\text{Leadership1} \triangleq \forall i, j \in \text{Server} :$
 $\quad \wedge \text{IsLeader}(i) \wedge \text{zabState}[i] \in \{\text{SYNCHRONIZATION}, \text{BROADCAST}\}$
 $\quad \wedge \text{IsLeader}(j) \wedge \text{zabState}[j] \in \{\text{SYNCHRONIZATION}, \text{BROADCAST}\}$
 $\quad \wedge \text{currentEpoch}[i] = \text{currentEpoch}[j]$
 $\quad \Rightarrow i = j$

$\text{Leadership2} \triangleq \forall \text{epoch} \in 1 \dots \text{MAXEPOCH} : \text{Cardinality}(\text{epochLeader}[\text{epoch}]) \leq 1$

PrefixConsistency: The prefix that have been committed in history in any process is the same.

$\text{PrefixConsistency} \triangleq \forall i, j \in \text{Server} :$
 $\quad \text{LET } \text{smaller} \triangleq \text{Minimum}(\{\text{lastCommitted}[i].\text{index}, \text{lastCommitted}[j].\text{index}\})$
 $\quad \text{IN } \quad \vee \text{smaller} = 0$
 $\quad \quad \vee \wedge \text{smaller} > 0$

$$\wedge \forall index \in 1 \dots smaller : \\ TxnEqual(history[i][index], history[j][index])$$

Integrity: If some follower delivers one transaction, then some primary has broadcast it.

$Integrity \triangleq \forall i \in Server :$

$$\begin{aligned} & \wedge IsFollower(i) \\ & \wedge lastCommitted[i].index > 0 \\ \Rightarrow & \forall idx \in 1 \dots lastCommitted[i].index : \exists proposal \in proposalMsgsLog : \\ & LET \quad txn_proposal \triangleq [zxid \mapsto proposal.zxid, \\ & \quad \quad \quad value \mapsto proposal.data] \\ & IN \quad TxnEqual(history[i][idx], txn_proposal) \end{aligned}$$

Agreement: If some follower f delivers transaction a and some follower f' delivers transaction b , then f' delivers a or f delivers b .

$Agreement \triangleq \forall i, j \in Server :$

$$\begin{aligned} & \wedge IsFollower(i) \wedge lastCommitted[i].index > 0 \\ & \wedge IsFollower(j) \wedge lastCommitted[j].index > 0 \\ \Rightarrow & \\ \forall idx1 \in 1 \dots lastCommitted[i].index, idx2 \in 1 \dots lastCommitted[j].index : & \\ \quad \vee \exists idx_j \in 1 \dots lastCommitted[j].index : & \\ \quad \quad TxnEqual(history[j][idx_j], history[i][idx1]) & \\ \quad \vee \exists idx_i \in 1 \dots lastCommitted[i].index : & \\ \quad \quad TxnEqual(history[i][idx_i], history[j][idx2]) & \end{aligned}$$

Total order: If some follower delivers a before b , then any process that delivers b must also deliver a and deliver a before b .

$TotalOrder \triangleq \forall i, j \in Server :$

$$\begin{aligned} & LET \quad committed1 \triangleq lastCommitted[i].index \\ & \quad \quad committed2 \triangleq lastCommitted[j].index \\ & IN \quad committed1 \geq 2 \wedge committed2 \geq 2 \\ & \quad \Rightarrow \forall idx_i1 \in 1 \dots (committed1 - 1) : \forall idx_i2 \in (idx_i1 + 1) \dots committed1 : \\ & \quad LET \quad logOk \triangleq \exists idx \in 1 \dots committed2 : \\ & \quad \quad TxnEqual(history[i][idx_i2], history[j][idx]) \\ & IN \quad \vee \neg logOk \\ & \quad \vee \wedge logOk \\ & \quad \quad \wedge \exists idx_j2 \in 1 \dots committed2 : \\ & \quad \quad \quad \wedge TxnEqual(history[i][idx_i2], history[j][idx_j2]) \\ & \quad \quad \quad \wedge \exists idx_j1 \in 1 \dots (idx_j2 - 1) : \\ & \quad \quad \quad TxnEqual(history[i][idx_i1], history[j][idx_j1]) \end{aligned}$$

Local primary order: If a primary broadcasts a before it broadcasts b , then a follower that delivers b must also deliver a before b .

$$\begin{aligned} LocalPrimaryOrder \triangleq LET \quad p_set(i, e) \triangleq \{p \in proposalMsgsLog : & \wedge p.source = i \\ & \wedge p.epoch = e\} \\ & txn_set(i, e) \triangleq \{[zxid \mapsto p.zxid, \\ & \quad \quad \quad value \mapsto p.data] : p \in p_set(i, e)\} \end{aligned}$$

$$\begin{aligned}
& \text{IN } \forall i \in \text{Server} : \forall e \in 1 \dots \text{currentEpoch}[i] : \\
& \quad \vee \text{Cardinality}(\text{txn_set}(i, e)) < 2 \\
& \quad \vee \wedge \text{Cardinality}(\text{txn_set}(i, e)) \geq 2 \\
& \quad \wedge \exists \text{txn1}, \text{txn2} \in \text{txn_set}(i, e) : \\
& \quad \quad \vee \text{TxnEqual}(\text{txn1}, \text{txn2}) \\
& \quad \quad \vee \wedge \neg \text{TxnEqual}(\text{txn1}, \text{txn2}) \\
& \quad \quad \wedge \text{LET } \text{TxnPre} \triangleq \text{IF } \text{ZxidCompare}(\text{txn1.zxid}, \text{txn2.zxid}) \text{ THEN } \text{txn2} \text{ ELSE} \\
& \quad \quad \quad \text{TxnNext} \triangleq \text{IF } \text{ZxidCompare}(\text{txn1.zxid}, \text{txn2.zxid}) \text{ THEN } \text{txn1} \text{ ELSE} \\
& \quad \quad \text{IN } \forall j \in \text{Server} : \wedge \text{lastCommitted}[j].\text{index} \geq 2 \\
& \quad \quad \quad \wedge \exists \text{idx} \in 1 \dots \text{lastCommitted}[j].\text{index} : \\
& \quad \quad \quad \quad \text{TxnEqual}(\text{history}[j][\text{idx}], \text{TxnNext}) \\
& \quad \quad \Rightarrow \exists \text{idx2} \in 1 \dots \text{lastCommitted}[j].\text{index} : \\
& \quad \quad \quad \wedge \text{TxnEqual}(\text{history}[j][\text{idx2}], \text{TxnNext}) \\
& \quad \quad \quad \wedge \text{idx2} > 1 \\
& \quad \quad \quad \wedge \exists \text{idx1} \in 1 \dots (\text{idx2} - 1) : \\
& \quad \quad \quad \quad \text{TxnEqual}(\text{history}[j][\text{idx1}], \text{TxnPre})
\end{aligned}$$

Global primary order: A follower f delivers both a with epoch e and b with epoch e' , and $e < e'$, then f must deliver a before b.

$$\begin{aligned}
\text{GlobalPrimaryOrder} & \triangleq \forall i \in \text{Server} : \text{lastCommitted}[i].\text{index} \geq 2 \\
& \Rightarrow \forall \text{idx1}, \text{idx2} \in 1 \dots \text{lastCommitted}[i].\text{index} : \\
& \quad \vee \neg \text{EpochPrecedeInTxn}(\text{history}[i][\text{idx1}], \text{history}[i][\text{idx2}]) \\
& \quad \vee \wedge \text{EpochPrecedeInTxn}(\text{history}[i][\text{idx1}], \text{history}[i][\text{idx2}]) \\
& \quad \wedge \text{idx1} < \text{idx2}
\end{aligned}$$

Primary integrity: If primary p broadcasts a and some follower f delivers b such that b has epoch smaller than epoch of p , then p must deliver b before it broadcasts a.

$$\begin{aligned}
\text{PrimaryIntegrity} & \triangleq \forall i, j \in \text{Server} : \wedge \text{IsLeader}(i) \wedge \text{IsMyLearner}(i, j) \\
& \quad \wedge \text{IsFollower}(j) \wedge \text{IsMyLeader}(j, i) \\
& \quad \wedge \text{zabState}[i] = \text{BROADCAST} \\
& \quad \wedge \text{zabState}[j] = \text{BROADCAST} \\
& \quad \wedge \text{lastCommitted}[j].\text{index} \geq 1 \\
& \Rightarrow \forall \text{idx_j} \in 1 \dots \text{lastCommitted}[j].\text{index} : \\
& \quad \vee \text{history}[j][\text{idx_j}].\text{zxid}[1] \geq \text{currentEpoch}[i] \\
& \quad \vee \wedge \text{history}[j][\text{idx_j}].\text{zxid}[1] < \text{currentEpoch}[i] \\
& \quad \quad \wedge \exists \text{idx_i} \in 1 \dots \text{lastCommitted}[i].\text{index} : \\
& \quad \quad \quad \text{TxnEqual}(\text{history}[i][\text{idx_i}], \text{history}[j][\text{idx_j}])
\end{aligned}$$

\ * Modification History
\ * Last modified Sat Dec 11 22:31:08 CST 2021 by Dell
\ * Created Thu Dec 02 20:49:23 CST 2021 by Dell