
MODULE *ZabWithFLEAndSYNC*

This is the formal specification for the *Zab* consensus algorithm, which means *Zookeeper* Atomic Broadcast. The differences from *ZabWithFLE* is that we implement phase RECOVERY-SYNC.

Reference: *FLE*: *FastLeaderElection.java*, *Vote.java*, *QuorumPeer.java*, e.g. in <https://github.com/apache/zookeeper>.

ZAB: *QuorumPeer.java*, *Learner.java*, *Follower.java*, *LearnerHandler.java*, *Leader.java*, e.g. in <https://github.com/apache/zookeeper>.
<https://cwiki.apache.org/confluence/display/ZOOKEEPER/Zab1.0>.

EXTENDS *FastLeaderElection*

The set of requests that can go into history

CONSTANT *Value*

Zab states

CONSTANTS *ELECTION*, *DISCOVERY*, *SYNCHRONIZATION*, *BROADCAST*

Sync modes & message types

CONSTANTS *DIFF*, *TRUNC*

Message types

CONSTANTS *FOLLOWERINFO*, *LEADERINFO*, *ACKEPOCH*, *NEWLEADER*, *ACKLD*,
UPTODATE, *PROPOSAL*, *ACK*, *COMMIT*

NOTE: In production, there is no message type *ACKLD*. Server judges if counter of *ACK* is 0 to distinguish one *ACK* represents *ACKLD* or not. Here we divide *ACK* into *ACKLD* and *ACK*, to enhance readability of spec.

[*MaxTimeoutFailures*, *MaxTransactionNum*, *MaxEpoch*]

CONSTANT *Parameters*

TODO: Here we can add more constraints to decrease space, like restart, partition.

MAXEPOCH \triangleq 10

Variables in annotations mean variables defined in *FastLeaderElection*.

Variables that all servers use.

VARIABLES <i>zabState</i> ,	Current phase of server, in { <i>ELECTION</i> , <i>DISCOVERY</i> , <i>SYNCHRONIZATION</i> , <i>BROADCAST</i> }.
<i>acceptedEpoch</i> ,	Epoch of the last <i>LEADERINFO</i> packet accepted, namely <i>f.p</i> in paper.
<i>lastCommitted</i> ,	Maximum index and <i>zxid</i> known to be committed, namely 'lastCommitted' in Leader. Starts from 0, and increases monotonically before restarting.
<i>initialHistory</i>	history that server initially has before election.
<i>state</i> ,	\ * State of server, in { <i>LOOKING</i> , <i>FOLLOWING</i> , <i>LEADING</i> }.
<i>currentEpoch</i> ,	\ * Epoch of the last <i>NEWLEADER</i> packet accepted, namely <i>f.a</i> in paper.
<i>lastProcessed</i> ,	\ * Index and <i>zxid</i> of the last processed <i>txn</i> .
<i>history</i>	\ * History of servers: sequence of transactions, containing: <i>zxid</i> , value, <i>ackSid</i> , epoch.

leader : $[committedRequests + toBeApplied]$ $[outstandingProposals]$
 follower: $[committedRequests]$ $[pendingTxns]$

Variables only used for leader.

VARIABLES *learners*, Set of servers leader connects,
 namely 'learners' in Leader.
connecting, Set of learners leader has received
FOLLOWERINFO from, namely
 'connectingFollowers' in Leader.
electing, Set of learners leader has received
ACKEPOCH from, namely 'electingFollowers'
 in Leader. Set of record
 $[sid, peerLastZxid, inQuorum]$.
 And $peerLastZxid = \langle -1, -1 \rangle$ means has done
syncFollower with this *sid*.
inQuorum = TRUE means in code it is one
 element in 'electingFollowers'.
ackldRecv, Set of learners leader has received
ACK of *NEWLEADER* from, namely
 'newLeaderProposal' in Leader.
forwarding, Set of learners that are synced with
 leader, namely 'forwardingFollowers'
 in Leader.
tempMaxEpoch $(\{Maximum\ epoch\ in\ FOLLOWERINFO\} + 1)$ that
 leader has received from learners,
 namely 'epoch' in Leader.
leadingVoteSet \ * Set of voters that follow leader.

Variables only used for follower.

VARIABLES *leaderAddr*, If follower has connected with leader.
 If follower lost connection, then null.
packetsSync packets of *PROPOSAL* and *COMMIT* from leader,
 namely 'packetsNotCommitted' and
 'packetsCommitted' in *SyncWithLeader*
 in Learner.

Variables about network channel.

VARIABLE *msgs* Simulates network channel.
 $msgs[i][j]$ means the input buffer of server *j*
 from server *i*.
electionMsgs \ * Network channel in *FLE* module.

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.

Variables only used for looking.

```

VARIABLE currentVote, \ * Info of current vote, namely 'currentVote'
      \ * in QuorumPeer.
      logicalClock, \ * Election instance, namely 'logicalClock'
      \ * in FastLeaderElection.
      receiveVotes, \ * Votes from current FLE round, namely
      \ * 'recvset' in FastLeaderElection.
      outOfElection, \ * Votes from previous and current FLE round,
      \ * namely 'outofelection' in FastLeaderElection.
      recvQueue, \ * Queue of received notifications or timeout
      \ * signals.
      waitNotmsg \ * Whether waiting for new not. See line 1050
      \ * in FastLeaderElection for details.
VARIABLE idTable \ * For mapping Server to Integers,
      to compare ids between servers.
Update: we have transformed idTable from variable to function.

```

Variable used for recording data to constrain state space.

VARIABLE *recorder* Consists: members of *Parameters* and *pc*.

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

$electionVars \triangleq electionVarsL$

$leaderVars \triangleq \langle leadingVoteSet, learners, connecting, electing, \\ ackldRecv, forwarding, tempMaxEpoch \rangle$

$followerVars \triangleq \langle leaderAddr, packetsSync \rangle$

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

$msgVars \triangleq \langle msgs, electionMsgs \rangle$

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

$ServersIncNullPoint \triangleq Server \cup \{NullPoint\}$

$Zxid \triangleq \\ Seq(Nat \cup \{-1\})$

$HistoryItem \triangleq \\ [zxid : Zxid, \\ value : Value, \\ ackSid : SUBSET Server, \\ epoch : Nat]$

$$\begin{aligned}
\textit{Proposal} &\triangleq \\
&\quad [\textit{source} : \textit{Server}, \\
&\quad \textit{epoch} : \textit{Nat}, \\
&\quad \textit{zxid} : \textit{Zxid}, \\
&\quad \textit{data} : \textit{Value}] \\
\\
\textit{LastItem} &\triangleq \\
&\quad [\textit{index} : \textit{Nat}, \textit{zxid} : \textit{Zxid}] \\
\\
\textit{SyncPackets} &\triangleq \\
&\quad [\textit{notCommitted} : \textit{Seq}(\textit{HistoryItem}), \\
&\quad \textit{committed} : \textit{Seq}(\textit{Zxid})] \\
\\
\textit{Message} &\triangleq \\
&\quad [\textit{mtype} : \{\textit{FOLLOWERINFO}\}, \textit{mzxid} : \textit{Zxid}] \cup \\
&\quad [\textit{mtype} : \{\textit{LEADERINFO}\}, \textit{mzxid} : \textit{Zxid}] \cup \\
&\quad [\textit{mtype} : \{\textit{ACKEPOCH}\}, \textit{mzxid} : \textit{Zxid}, \textit{mepoch} : \textit{Nat} \cup \{-1\}] \cup \\
&\quad [\textit{mtype} : \{\textit{DIFF}\}, \textit{mzxid} : \textit{Zxid}] \cup \\
&\quad [\textit{mtype} : \{\textit{TRUNC}\}, \textit{mtruncZxid} : \textit{Zxid}] \cup \\
&\quad [\textit{mtype} : \{\textit{PROPOSAL}\}, \textit{mzxid} : \textit{Zxid}, \textit{mdata} : \textit{Value}] \cup \\
&\quad [\textit{mtype} : \{\textit{COMMIT}\}, \textit{mzxid} : \textit{Zxid}] \cup \\
&\quad [\textit{mtype} : \{\textit{NEWLEADER}\}, \textit{mzxid} : \textit{Zxid}] \cup \\
&\quad [\textit{mtype} : \{\textit{ACKLD}\}, \textit{mzxid} : \textit{Zxid}] \cup \\
&\quad [\textit{mtype} : \{\textit{ACK}\}, \textit{mzxid} : \textit{Zxid}] \cup \\
&\quad [\textit{mtype} : \{\textit{UPTODATE}\}, \textit{mzxid} : \textit{Zxid}] \\
\\
\textit{ElectionState} &\triangleq \{\textit{LOOKING}, \textit{FOLLOWING}, \textit{LEADING}\} \\
\\
\textit{ZabState} &\triangleq \{\textit{ELECTION}, \textit{DISCOVERY}, \textit{SYNCHRONIZATION}, \textit{BROADCAST}\} \\
\\
\textit{ViolationSet} &\triangleq \{\text{"stateInconsistent"}, \text{"proposalInconsistent"}, \\
&\quad \text{"commitInconsistent"}, \text{"ackInconsistent"}, \\
&\quad \text{"messageIllegal"}\} \\
\\
\textit{Electing} &\triangleq [\textit{sid} : \textit{Server}, \\
&\quad \textit{peerLastZxid} : \textit{Zxid}, \\
&\quad \textit{inQuorum} : \text{BOOLEAN}] \\
\\
\textit{Vote} &\triangleq \\
&\quad [\textit{proposedLeader} : \textit{ServersIncNullPoint}, \\
&\quad \textit{proposedZxid} : \textit{Zxid}, \\
&\quad \textit{proposedEpoch} : \textit{Nat}] \\
\\
\textit{ElectionVote} &\triangleq \\
&\quad [\textit{vote} : \textit{Vote}, \textit{round} : \textit{Nat}, \textit{state} : \textit{ElectionState}, \textit{version} : \textit{Nat}] \\
\\
\textit{ElectionMsg} &\triangleq \\
&\quad [\textit{mtype} : \{\textit{NOTIFICATION}\},
\end{aligned}$$

$m_{source} : Server,$
 $m_{state} : ElectionState,$
 $m_{round} : Nat,$
 $m_{vote} : Vote] \cup$
 $[m_{type} : \{NONE\}]$

$TypeOK \triangleq$

$\wedge \text{ zabState} \in [Server \rightarrow ZabState]$
 $\wedge \text{ acceptedEpoch} \in [Server \rightarrow Nat]$
 $\wedge \text{ lastCommitted} \in [Server \rightarrow LastItem]$
 $\wedge \text{ learners} \in [Server \rightarrow SUBSET \ Server]$
 $\wedge \text{ connecting} \in [Server \rightarrow SUBSET \ ServersIncNullPoint]$
 $\wedge \text{ electing} \in [Server \rightarrow SUBSET \ Electing]$
 $\wedge \text{ ackldRecv} \in [Server \rightarrow SUBSET \ ServersIncNullPoint]$
 $\wedge \text{ forwarding} \in [Server \rightarrow SUBSET \ Server]$
 $\wedge \text{ initialHistory} \in [Server \rightarrow Seq(HistoryItem)]$
 $\wedge \text{ tempMaxEpoch} \in [Server \rightarrow Nat]$
 $\wedge \text{ leaderAddr} \in [Server \rightarrow ServersIncNullPoint]$
 $\wedge \text{ packetsSync} \in [Server \rightarrow SyncPackets]$
 $\wedge \text{ proposalMsgsLog} \in SUBSET \ Proposal$
 $\wedge \text{ epochLeader} \in [1 \dots MAXEPOCH \rightarrow SUBSET \ Server]$
 $\wedge \text{ violatedInvariants} \in [ViolationSet \rightarrow BOOLEAN]$
 $\wedge \text{ msgs} \in [Server \rightarrow [Server \rightarrow Seq(Message)]]$
Fast Leader Election
 $\wedge \text{ electionMsgs} \in [Server \rightarrow [Server \rightarrow Seq(ElectionMsg)]]$
 $\wedge \text{ recvQueue} \in [Server \rightarrow Seq(ElectionMsg)]$
 $\wedge \text{ leadingVoteSet} \in [Server \rightarrow SUBSET \ Server]$
 $\wedge \text{ receiveVotes} \in [Server \rightarrow [Server \rightarrow ElectionVote]]$
 $\wedge \text{ currentVote} \in [Server \rightarrow Vote]$
 $\wedge \text{ outOfElection} \in [Server \rightarrow [Server \rightarrow ElectionVote]]$
 $\wedge \text{ lastProcessed} \in [Server \rightarrow LastItem]$
 $\wedge \text{ history} \in [Server \rightarrow Seq(HistoryItem)]$
 $\wedge \text{ state} \in [Server \rightarrow ElectionState]$
 $\wedge \text{ waitNotmsg} \in [Server \rightarrow BOOLEAN]$
 $\wedge \text{ currentEpoch} \in [Server \rightarrow Nat]$
 $\wedge \text{ logicalClock} \in [Server \rightarrow Nat]$

Return the maximum value from the set S

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

Return the minimum value from the set S

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

Check server state

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

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

Check if s is a quorum
 $IsQuorum(s) \triangleq s \in Quorums$

Check $zxid$ state
 $ToZxid(z) \triangleq [epoch \mapsto z[1], counter \mapsto z[2]]$
 $TxnZxidEqual(txn, z) \triangleq txn.zxid[1] = z[1] \wedge txn.zxid[2] = z[2]$
 $TxnEqual(txn1, txn2) \triangleq ZxidEqual(txn1.zxid, txn2.zxid)$
 $EpochPrecedeInTxn(txn1, txn2) \triangleq txn1.zxid[1] < txn2.zxid[1]$

Actions about recorder
 $GetParameter(p) \triangleq \text{IF } p \in \text{DOMAIN } Parameters \text{ THEN } Parameters[p] \text{ ELSE } 0$
 $RecorderGetHelper(m) \triangleq (m := recorder[m])$
 $RecorderIncHelper(m) \triangleq (m := recorder[m] + 1)$
 $RecorderIncTimeout \triangleq RecorderIncHelper("nTimeout")$
 $RecorderGetTimeout \triangleq RecorderGetHelper("nTimeout")$
 $RecorderSetTransactionNum(pc) \triangleq ("nTransaction" :=$
 IF $pc[1] = "LeaderProcessRequest"$ 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 $recorder["nTransaction"]$)
 $RecorderSetMaxEpoch(pc) \triangleq ("maxEpoch" :=$
 IF $pc[1] = "LeaderProcessFOLLOWERINFO"$ THEN
 LET $s \triangleq \text{CHOOSE } i \in Server :$
 $\forall j \in Server : acceptedEpoch'[i] \geq acceptedEpoch'[j]$
 IN $acceptedEpoch'[s]$
 ELSE $recorder["maxEpoch"]$)
 $RecorderSetPc(pc) \triangleq ("pc" := pc)$
 $RecorderSetFailure(pc) \triangleq \text{CASE } pc[1] = "Timeout" \rightarrow RecorderIncTimeout$
 $\square \quad pc[1] = "LeaderTimeout" \rightarrow RecorderIncTimeout$
 $\square \quad pc[1] = "FollowerTimeout" \rightarrow RecorderIncTimeout$
 $\square \quad \text{OTHER} \rightarrow RecorderGetTimeout$

$UpdateRecorder(pc) \triangleq recorder' = RecorderSetFailure(pc) \quad @@ RecorderSetTransactionNum(pc)$
 $\quad @@ RecorderSetMaxEpoch(pc) \quad @@ RecorderSetPc(pc) \quad @@ recorder$
 $UnchangeRecorder \triangleq UNCHANGED recorder$
 $CheckParameterHelper(n, p, Comp(-, -)) \triangleq \text{IF } p \in \text{DOMAIN } Parameters$
 $\quad \text{THEN } Comp(n, Parameters[p])$
 $\quad \text{ELSE TRUE}$
 $CheckParameterLimit(n, p) \triangleq CheckParameterHelper(n, p, \text{LAMBDA } i, j : i < j)$
 $CheckTimeout \triangleq CheckParameterLimit(recorder.nTimeout, \text{"MaxTimeoutFailures"})$
 $CheckTransactionNum \triangleq CheckParameterLimit(recorder.nTransaction, \text{"MaxTransactionNum"})$
 $CheckEpoch \triangleq CheckParameterLimit(recorder.maxEpoch, \text{"MaxEpoch"})$
 $CheckStateConstraints \triangleq CheckTimeout \wedge CheckTransactionNum \wedge CheckEpoch$

Actions about network

$PendingFOLLOWERINFO(i, j) \triangleq \wedge msgs[j][i] \neq \langle \rangle$
 $\quad \wedge msgs[j][i][1].mtype = FOLLOWERINFO$
 $PendingLEADERINFO(i, j) \triangleq \wedge msgs[j][i] \neq \langle \rangle$
 $\quad \wedge msgs[j][i][1].mtype = LEADERINFO$
 $PendingACKEPOCH(i, j) \triangleq \wedge msgs[j][i] \neq \langle \rangle$
 $\quad \wedge msgs[j][i][1].mtype = ACKEPOCH$
 $PendingNEWLEADER(i, j) \triangleq \wedge msgs[j][i] \neq \langle \rangle$
 $\quad \wedge msgs[j][i][1].mtype = NEWLEADER$
 $PendingACKLD(i, j) \triangleq \wedge msgs[j][i] \neq \langle \rangle$
 $\quad \wedge msgs[j][i][1].mtype = ACKLD$
 $PendingUPTODATE(i, j) \triangleq \wedge msgs[j][i] \neq \langle \rangle$
 $\quad \wedge msgs[j][i][1].mtype = UPTODATE$
 $PendingPROPOSAL(i, j) \triangleq \wedge msgs[j][i] \neq \langle \rangle$
 $\quad \wedge msgs[j][i][1].mtype = PROPOSAL$
 $PendingACK(i, j) \triangleq \wedge msgs[j][i] \neq \langle \rangle$
 $\quad \wedge msgs[j][i][1].mtype = ACK$
 $PendingCOMMIT(i, j) \triangleq \wedge msgs[j][i] \neq \langle \rangle$
 $\quad \wedge msgs[j][i][1].mtype = COMMIT$

Add a message to $msgs$ – add a message m to $msgs$.

$Send(i, j, m) \triangleq msgs' = [msgs \text{ EXCEPT } ![i][j] = Append(msgs[i][j], m)]$
 $SendPackets(i, j, ms) \triangleq msgs' = [msgs \text{ EXCEPT } ![i][j] = msgs[i][j] \circ ms]$
 $DiscardAndSendPackets(i, j, ms) \triangleq msgs' = [msgs \text{ EXCEPT } ![j][i] = Tail(msgs[j][i]),$
 $\quad ![i][j] = msgs[i][j] \circ ms]$

Remove a message from $msgs$ – discard head of $msgs$.

$Discard(i, j) \triangleq msgs' = \text{IF } msgs[i][j] \neq \langle \rangle \text{ THEN } [msgs \text{ EXCEPT } ![i][j] = Tail(msgs[i][j])]$
 $\quad \text{ELSE } msgs$

Leader broadcasts a message($PROPOSAL/COMMIT$) to all other servers in $forwardingFollowers$.

$Broadcast(i, m) \triangleq msgs' = [msgs \text{ EXCEPT } ![i] = [v \in Server \mapsto \text{IF } \wedge v \in forwarding[i]$
 $\quad \wedge v \neq i$
 $\quad \text{THEN } Append(msgs[i][v], m)]$

ELSE $msgs[i][v]$]

$DiscardAndBroadcast(i, j, m) \triangleq$
 $msgs' = [msgs \text{ EXCEPT } ![j][i] = Tail(msgs[j][i]),$
 $![i] = [v \in Server \mapsto \text{IF } \wedge v \in forwarding[i]$
 $\wedge v \neq i$
THEN $Append(msgs[i][v], m)$
ELSE $msgs[i][v]]$

Leader broadcasts *LEADERINFO* to all other servers in *connectingFollowers*.

$DiscardAndBroadcastLEADERINFO(i, j, m) \triangleq$
 $msgs' = [msgs \text{ EXCEPT } ![j][i] = Tail(msgs[j][i]),$
 $![i] = [v \in Server \mapsto \text{IF } \wedge v \in connecting'[i]$
 $\wedge v \in learners[i]$
 $\wedge v \neq i$
THEN $Append(msgs[i][v], m)$
ELSE $msgs[i][v]]$

Leader broadcasts *UPTODATE* to all other servers in *newLeaderProposal*.

$DiscardAndBroadcastUPTODATE(i, j, m) \triangleq$
 $msgs' = [msgs \text{ EXCEPT } ![j][i] = Tail(msgs[j][i]),$
 $![i] = [v \in Server \mapsto \text{IF } \wedge v \in ackldRecv'[i]$
 $\wedge v \in learners[i]$
 $\wedge v \neq i$
THEN $Append(msgs[i][v], m)$
ELSE $msgs[i][v]]$

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.

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

$CleanInputBufferInCluster(S) \triangleq msgs' = [s \in Server \mapsto$
 $[v \in Server \mapsto \text{IF } v \in S \text{ THEN } \langle \rangle$
ELSE $msgs[s][v]]]$

Define initial values for all variables

$InitServerVars \triangleq \wedge InitServerVarsL$
 $\wedge zabState = [s \in Server \mapsto ELECTION]$
 $\wedge acceptedEpoch = [s \in Server \mapsto 0]$
 $\wedge lastCommitted = [s \in Server \mapsto [index \mapsto 0,$
 $zxid \mapsto \langle 0, 0 \rangle]]$
 $\wedge initialHistory = [s \in Server \mapsto \langle \rangle]$

$InitLeaderVars \triangleq \wedge InitLeaderVarsL$
 $\wedge learners = [s \in Server \mapsto \{\}]$

$$\begin{aligned}
\wedge \text{connecting} &= [s \in \text{Server} \mapsto \{\}] \\
\wedge \text{electing} &= [s \in \text{Server} \mapsto \{\}] \\
\wedge \text{ackldRecv} &= [s \in \text{Server} \mapsto \{\}] \\
\wedge \text{forwarding} &= [s \in \text{Server} \mapsto \{\}] \\
\wedge \text{tempMaxEpoch} &= [s \in \text{Server} \mapsto 0]
\end{aligned}$$

$$\text{InitElectionVars} \triangleq \text{InitElectionVarsL}$$

$$\begin{aligned}
\text{InitFollowerVars} \triangleq & \wedge \text{leaderAddr} = [s \in \text{Server} \mapsto \text{NullPoint}] \\
& \wedge \text{packetsSync} = [s \in \text{Server} \mapsto \\
& \quad [\text{notCommitted} \mapsto \langle \rangle, \\
& \quad \text{committed} \mapsto \langle \rangle]]
\end{aligned}$$

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

$$\begin{aligned}
\text{InitMsgVars} \triangleq & \wedge \text{msgs} = [s \in \text{Server} \mapsto [v \in \text{Server} \mapsto \langle \rangle]] \\
& \wedge \text{electionMsgs} = [s \in \text{Server} \mapsto [v \in \text{Server} \mapsto \langle \rangle]]
\end{aligned}$$

$$\begin{aligned}
\text{InitRecorder} \triangleq & \text{recorder} = [n\text{Timeout} \mapsto 0, \\
& \quad n\text{Transaction} \mapsto 0, \\
& \quad \text{maxEpoch} \mapsto 0, \\
& \quad \text{pc} \mapsto \langle \text{"Init"} \rangle]
\end{aligned}$$

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

$$\begin{aligned}
\text{ZabTurnToLeading}(i) \triangleq & \\
& \wedge \text{zabState}' = [\text{zabState} \text{ EXCEPT } ![i] = \text{DISCOVERY}] \\
& \wedge \text{learners}' = [\text{learners} \text{ EXCEPT } ![i] = \{i\}] \\
& \wedge \text{connecting}' = [\text{connecting} \text{ EXCEPT } ![i] = \{i\}] \\
& \wedge \text{electing}' = [\text{electing} \text{ EXCEPT } ![i] = \{[sid \mapsto i, \\
& \quad \text{peerLastZxid} \mapsto \langle -1, -1 \rangle, \\
& \quad \text{inQuorum} \mapsto \text{TRUE}] \}] \\
& \wedge \text{ackldRecv}' = [\text{ackldRecv} \text{ EXCEPT } ![i] = \{i\}] \\
& \wedge \text{forwarding}' = [\text{forwarding} \text{ EXCEPT } ![i] = \{\}] \\
& \wedge \text{initialHistory}' = [\text{initialHistory} \text{ EXCEPT } ![i] = \text{history}'[i]]
\end{aligned}$$

$$\begin{aligned}
& \wedge tempMaxEpoch' = [tempMaxEpoch \quad \text{EXCEPT } ![i] = acceptedEpoch[i] + 1] \\
ZabTurnToFollowing(i) & \triangleq \\
& \wedge zabState' = [zabState \quad \text{EXCEPT } ![i] = DISCOVERY] \\
& \wedge initialHistory' = [initialHistory \quad \text{EXCEPT } ![i] = history'[i]] \\
& \wedge packetsSync' = [packetsSync \quad \text{EXCEPT } ![i].notCommitted = \langle \rangle, \\
& \hspace{15em} ![i].committed = \langle \rangle] \\
\text{Fast Leader Election} \\
FLEReceiveNotmsg(i, j) & \triangleq \\
& \wedge ReceiveNotmsg(i, j) \\
& \wedge \text{UNCHANGED } \langle zabState, acceptedEpoch, lastCommitted, learners, connecting, \\
& \hspace{10em} initialHistory, electing, ackldRecv, forwarding, tempMaxEpoch, \\
& \hspace{10em} followerVars, verifyVars, msgs \rangle \\
& \wedge UpdateRecorder(\langle \text{"FLEReceiveNotmsg"}, i, j \rangle) \\
FLENotmsgTimeout(i) & \triangleq \\
& \wedge NotmsgTimeout(i) \\
& \wedge \text{UNCHANGED } \langle zabState, acceptedEpoch, lastCommitted, learners, connecting, \\
& \hspace{10em} initialHistory, electing, ackldRecv, forwarding, tempMaxEpoch, \\
& \hspace{10em} followerVars, verifyVars, msgs \rangle \\
& \wedge UpdateRecorder(\langle \text{"FLENotmsgTimeout"}, i \rangle) \\
FLEHandleNotmsg(i) & \triangleq \\
& \wedge HandleNotmsg(i) \\
& \wedge \text{LET } newState \triangleq state'[i] \\
& \text{IN} \\
& \vee \wedge newState = LEADING \\
& \hspace{2em} \wedge ZabTurnToLeading(i) \\
& \hspace{2em} \wedge \text{UNCHANGED } packetsSync \\
& \vee \wedge newState = FOLLOWING \\
& \hspace{2em} \wedge ZabTurnToFollowing(i) \\
& \hspace{2em} \wedge \text{UNCHANGED } \langle learners, connecting, electing, ackldRecv, \\
& \hspace{10em} forwarding, tempMaxEpoch \rangle \\
& \vee \wedge newState = LOOKING \\
& \hspace{2em} \wedge \text{UNCHANGED } \langle zabState, learners, connecting, electing, ackldRecv, \\
& \hspace{10em} forwarding, tempMaxEpoch, packetsSync, initialHistory \rangle \\
& \wedge \text{UNCHANGED } \langle lastCommitted, acceptedEpoch, leaderAddr, verifyVars, msgs \rangle \\
& \wedge UpdateRecorder(\langle \text{"FLEHandleNotmsg"}, i \rangle)
\end{aligned}$$

On the premise that $ReceiveVotes.HasQuorums = \text{TRUE}$,
corresponding to logic in line 1050 – 1055 in *FastLeaderElection*.

$$\begin{aligned}
FLEWaitNewNotmsg(i) & \triangleq \\
& \wedge WaitNewNotmsg(i) \\
& \wedge \text{UNCHANGED } \langle zabState, acceptedEpoch, lastCommitted, learners, connecting, \\
& \hspace{10em} electing, ackldRecv, forwarding, tempMaxEpoch, initialHistory,
\end{aligned}$$

$followerVars, verifyVars, msgs\rangle$
 $\wedge UpdateRecorder(\langle \text{"FLEWaitNewNotmsgEnd"}, i \rangle)$

On the premise that $ReceiveVotes.HasQuorums = \text{TRUE}$,
corresponding to logic in line 1061 – 1066 in *FastLeaderElection*.
 $FLEWaitNewNotmsgEnd(i) \triangleq$
 $\wedge WaitNewNotmsgEnd(i)$
 $\wedge \text{LET } newState \triangleq state'[i]$
IN
 $\vee \wedge newState = LEADING$
 $\wedge ZabTurnToLeading(i)$
 $\wedge \text{UNCHANGED } packetsSync$
 $\vee \wedge newState = FOLLOWING$
 $\wedge ZabTurnToFollowing(i)$
 $\wedge \text{UNCHANGED } \langle learners, connecting, electing, ackldRecv, forwarding,$
 $tempMaxEpoch \rangle$
 $\vee \wedge newState = LOOKING$
 $\wedge PrintT(\text{"Note: New state is LOOKING in FLEWaitNewNotmsgEnd,"} \circ$
 $\text{" which should not happen."})$
 $\wedge \text{UNCHANGED } \langle zabState, learners, connecting, electing, ackldRecv,$
 $forwarding, tempMaxEpoch, initialHistory, packetsSync \rangle$
 $\wedge \text{UNCHANGED } \langle lastCommitted, acceptedEpoch, leaderAddr, verifyVars, msgs \rangle$
 $\wedge UpdateRecorder(\langle \text{"FLEWaitNewNotmsgEnd"}, i \rangle)$

$InitialVotes \triangleq [vote \mapsto InitialVote,$
 $round \mapsto 0,$
 $state \mapsto LOOKING,$
 $version \mapsto 0]$

Equals to for every server in S , performing action $ZabTimeout$.
 $ZabTimeoutInCluster(S) \triangleq$
 $\wedge state' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } LOOKING \text{ ELSE } state[s]]$
 $\wedge lastProcessed' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } InitLastProcessed(s)$
 $\text{ELSE } lastProcessed[s]]$
 $\wedge logicalClock' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } logicalClock[s] + 1$
 $\text{ELSE } logicalClock[s]]$
 $\wedge currentVote' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN}$
 $[proposedLeader \mapsto s,$
 $proposedZxid \mapsto lastProcessed'[s].zxid,$
 $proposedEpoch \mapsto currentEpoch[s]]$
 $\text{ELSE } currentVote[s]]$
 $\wedge receiveVotes' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } [v \in Server \mapsto InitialVotes]$
 $\text{ELSE } receiveVotes[s]]$
 $\wedge outOfElection' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } [v \in Server \mapsto InitialVotes]$
 $\text{ELSE } outOfElection[s]]$

$$\begin{aligned}
& \wedge \text{recvQueue}' = [s \in \text{Server} \mapsto \text{IF } s \in S \text{ THEN } \langle [mtype \mapsto \text{NONE}] \rangle \\
& \quad \quad \quad \text{ELSE } \text{recvQueue}[s]] \\
& \wedge \text{waitNotmsg}' = [s \in \text{Server} \mapsto \text{IF } s \in S \text{ THEN } \text{FALSE} \text{ ELSE } \text{waitNotmsg}[s]] \\
& \wedge \text{leadingVoteSet}' = [s \in \text{Server} \mapsto \text{IF } s \in S \text{ THEN } \{\} \text{ ELSE } \text{leadingVoteSet}[s]] \\
& \wedge \text{UNCHANGED } \langle \text{electionMsgs}, \text{currentEpoch}, \text{history} \rangle \\
& \wedge \text{zabState}' = [s \in \text{Server} \mapsto \text{IF } s \in S \text{ THEN } \text{ELECTION} \text{ ELSE } \text{zabState}[s]] \\
& \wedge \text{leaderAddr}' = [s \in \text{Server} \mapsto \text{IF } s \in S \text{ THEN } \text{NullPoint} \text{ ELSE } \text{leaderAddr}[s]] \\
& \wedge \text{CleanInputBufferInCluster}(S)
\end{aligned}$$

Describe how a server transitions from *LEADING/FOLLOWING* to *LOOKING*.

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

$$\begin{aligned}
\text{LeaderShutdown}(i) & \triangleq \\
& \wedge \text{LET } \text{cluster} \triangleq \{i\} \cup \text{learners}[i] \\
& \quad \text{IN } \text{ZabTimeoutInCluster}(\text{cluster}) \\
& \wedge \text{learners}' = [\text{learners} \text{ EXCEPT } ![i] = \{\}] \\
& \wedge \text{forwarding}' = [\text{forwarding} \text{ EXCEPT } ![i] = \{\}]
\end{aligned}$$

$$\begin{aligned}
\text{RemoveElecting}(\text{set}, \text{sid}) & \triangleq \\
& \text{LET } \text{sid_electing} \triangleq \{s.\text{sid} : s \in \text{set}\} \\
& \text{IN } \text{IF } \text{sid} \notin \text{sid_electing} \text{ THEN } \text{set} \\
& \quad \text{ELSE } \text{LET } \text{info} \triangleq \text{CHOOSE } s \in \text{set} : s.\text{sid} = \text{sid} \\
& \quad \quad \text{new_info} \triangleq [\text{sid} \mapsto \text{sid}, \\
& \quad \quad \quad \text{peerLastZxid} \mapsto \langle -1, -1 \rangle, \\
& \quad \quad \quad \text{inQuorum} \mapsto \text{info.inQuorum}] \\
& \quad \text{IN } (\text{set} \setminus \{\text{info}\}) \cup \{\text{new_info}\}
\end{aligned}$$

See *removeLearnerHandler* for details.

$$\begin{aligned}
\text{RemoveLearner}(i, j) & \triangleq \\
& \wedge \text{learners}' = [\text{learners} \text{ EXCEPT } ![i] = @ \setminus \{j\}] \\
& \wedge \text{forwarding}' = [\text{forwarding} \text{ EXCEPT } ![i] = \text{IF } j \in \text{forwarding}[i] \\
& \quad \quad \quad \text{THEN } @ \setminus \{j\} \text{ ELSE } @] \\
& \wedge \text{electing}' = [\text{electing} \text{ EXCEPT } ![i] = \text{RemoveElecting}(@, j)]
\end{aligned}$$

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

$$\begin{aligned}
\text{FollowerTimeout}(i) & \triangleq \\
& \wedge \text{CheckTimeout} \quad \text{test restrictions of } \text{timeout_1} \\
& \wedge \text{IsFollower}(i) \\
& \wedge \text{HasNoLeader}(i) \\
& \wedge \text{FollowerShutdown}(i) \\
& \wedge \text{CleanInputBuffer}(i) \\
& \wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{lastCommitted}, \text{learners}, \text{connecting}, \text{electing},
\end{aligned}$$

ackldRecv, forwarding, tempMaxEpoch, initialHistory,
verifyVars, packetsSync
 $\wedge \text{UpdateRecorder}(\langle \text{"FollowerTimeout"}, i \rangle)$

Leader loses support from a quorum and turns to *LOOKING*.
 $\text{LeaderTimeout}(i) \triangleq$
 $\wedge \text{CheckTimeout}$ test restrictions of *timeout_2*
 $\wedge \text{IsLeader}(i)$
 $\wedge \neg \text{IsQuorum}(\text{learners}[i])$
 $\wedge \text{LeaderShutdown}(i)$
 $\wedge \text{UNCHANGED} \langle \text{acceptedEpoch}, \text{lastCommitted}, \text{connecting}, \text{electing}, \text{ackldRecv},$
 $\text{tempMaxEpoch}, \text{initialHistory}, \text{verifyVars}, \text{packetsSync} \rangle$
 $\wedge \text{UpdateRecorder}(\langle \text{"LeaderTimeout"}, i \rangle)$

Timeout between leader and follower.
 $\text{Timeout}(i, j) \triangleq$
 $\wedge \text{CheckTimeout}$ test restrictions of *timeout_3*
 $\wedge \text{IsLeader}(i) \wedge \text{IsMyLearner}(i, j)$
 $\wedge \text{IsFollower}(j) \wedge \text{IsMyLeader}(j, i)$
The action of leader *i*.
 $\wedge \text{RemoveLearner}(i, j)$
The action of follower *j*.
 $\wedge \text{FollowerShutdown}(j)$
 $\wedge \text{Clean}(i, j)$
 $\wedge \text{UNCHANGED} \langle \text{acceptedEpoch}, \text{lastCommitted}, \text{connecting}, \text{ackldRecv},$
 $\text{tempMaxEpoch}, \text{initialHistory}, \text{verifyVars}, \text{packetsSync} \rangle$
 $\wedge \text{UpdateRecorder}(\langle \text{"Timeout"}, i, j \rangle)$

Establish connection between leader and follower, containing actions like *addLearnerHandler*,
findLeader, *connectToLeader*.

$\text{ConnectAndFollowerSendFOLLOWERINFO}(i, j) \triangleq$
 $\wedge \text{IsLeader}(i) \wedge \neg \text{IsMyLearner}(i, j)$
 $\wedge \text{IsFollower}(j) \wedge \text{HasNoLeader}(j) \wedge \text{MyVote}(j) = i$
 $\wedge \text{learners}' = [\text{learners} \text{ EXCEPT } ![i] = \text{learners}[i] \cup \{j\}]$
 $\wedge \text{leaderAddr}' = [\text{leaderAddr} \text{ EXCEPT } ![j] = i]$
 $\wedge \text{Send}(j, \text{leaderAddr}'[j], [\text{mtype} \mapsto \text{FOLLOWERINFO},$
 $\text{mzxid} \mapsto \langle \text{acceptedEpoch}[j], 0 \rangle])$
 $\wedge \text{UNCHANGED} \langle \text{serverVars}, \text{electionVars}, \text{leadingVoteSet}, \text{connecting},$
 $\text{electing}, \text{ackldRecv}, \text{forwarding}, \text{tempMaxEpoch},$
 $\text{verifyVars}, \text{electionMsgs}, \text{packetsSync} \rangle$
 $\wedge \text{UpdateRecorder}(\langle \text{"ConnectAndFollowerSendFOLLOWERINFO"}, i, j \rangle)$

waitingForNewEpoch in Leader

$\text{WaitingForNewEpoch}(i) \triangleq (i \in \text{connecting}[i] \wedge \text{IsQuorum}(\text{connecting}[i])) = \text{FALSE}$
 $\text{WaitingForNewEpochTurnToFalse}(i) \triangleq \wedge i \in \text{connecting}'[i]$

$\wedge IsQuorum(connecting'[i])$

Leader waits for receiving *FOLLOWERINFO* from a quorum including itself, and chooses a new epoch e' as its own epoch and broadcasts *LEADERINFO*. See *getEpochToPropose* in Leader for details.

$LeaderProcessFOLLOWERINFO(i, j) \triangleq$
 $\wedge CheckEpoch$ test restrictions of max epoch
 $\wedge IsLeader(i)$
 $\wedge PendingFOLLOWERINFO(i, j)$
 $\wedge LET\ msg \triangleq msgs[j][i][1]$
 $\quad infoOk \triangleq IsMyLearner(i, j)$
 $\quad lastAcceptedEpoch \triangleq msg.mzxid[1]$
 IN
 $\wedge infoOk$
 $\wedge \vee$ 1. has not broadcast *LEADERINFO*
 $\quad \wedge WaitingForNewEpoch(i)$
 $\quad \wedge \vee \wedge zabState[i] = DISCOVERY$
 $\quad \quad \wedge UNCHANGED\ violatedInvariants$
 $\quad \vee \wedge zabState[i] \neq DISCOVERY$
 $\quad \wedge PrintT("Exception: waitingForNewEpoch true," \circ$
 $\quad \quad "while zabState not DISCOVERY.")$
 $\quad \wedge violatedInvariants' = [violatedInvariants\ EXCEPT\ !.stateInconsistent = TRUE]$
 $\quad \wedge tempMaxEpoch' = [tempMaxEpoch\ EXCEPT\ ![i] = IF\ lastAcceptedEpoch \geq tempMaxEpoch[i]$
 $\quad \quad THEN\ lastAcceptedEpoch + 1$
 $\quad \quad ELSE\ @]$
 $\quad \wedge connecting' = [connecting\ EXCEPT\ ![i] = @ \cup \{j\}]$
 $\quad \wedge \vee \wedge WaitingForNewEpochTurnToFalse(i)$
 $\quad \quad \wedge acceptedEpoch' = [acceptedEpoch\ EXCEPT\ ![i] = tempMaxEpoch'[i]]$
 $\quad \quad \wedge LET\ newLeaderZxid \triangleq \langle acceptedEpoch'[i], 0 \rangle$
 $\quad \quad \quad m \triangleq [mtype \mapsto LEADERINFO,$
 $\quad \quad \quad \quad mxid \mapsto newLeaderZxid]$
 $\quad \quad \quad IN\ DiscardAndBroadcastLEADERINFO(i, j, m)$
 $\quad \vee \wedge \neg WaitingForNewEpochTurnToFalse(i)$
 $\quad \quad \wedge Discard(j, i)$
 $\quad \quad \wedge UNCHANGED\ acceptedEpoch$
 \vee 2. has broadcast *LEADERINFO*
 $\quad \wedge \neg WaitingForNewEpoch(i)$
 $\quad \wedge Reply(i, j, [mtype \mapsto LEADERINFO,$
 $\quad \quad \quad mxid \mapsto \langle acceptedEpoch[i], 0 \rangle])$
 $\quad \wedge UNCHANGED\ \langle tempMaxEpoch, connecting, acceptedEpoch, violatedInvariants \rangle$
 $\wedge UNCHANGED\ \langle state, currentEpoch, lastProcessed, zabState, history, lastCommitted,$
 $\quad followerVars, electionVars, initialHistory, leadingVoteSet, learners,$
 $\quad \quad \quad electing, ackldRecv, forwarding, proposalMsgsLog, epochLeader,$
 $\quad \quad \quad \quad electionMsgs \rangle$
 $\wedge UpdateRecorder(\langle "LeaderProcessFOLLOWERINFO", i, j \rangle)$

Follower receives *LEADERINFO*. If $newEpoch \geq acceptedEpoch$, then follower updates $acceptedEpoch$ and sends *ACKEPOCH* back, containing $currentEpoch$ and $lastProcessedZxid$. After this, $zabState$ turns to *SYNC*. See *registerWithLeader* in Learner for details.

$$\begin{aligned}
& \text{FollowerProcessLEADERINFO}(i, j) \triangleq \\
& \quad \wedge \text{IsFollower}(i) \\
& \quad \wedge \text{PendingLEADERINFO}(i, j) \\
& \quad \wedge \text{LET } msg \triangleq msgs[j][i][1] \\
& \quad \quad newEpoch \triangleq msg.mzxid[1] \\
& \quad \quad infoOk \triangleq \text{IsMyLeader}(i, j) \\
& \quad \quad epochOk \triangleq newEpoch \geq acceptedEpoch[i] \\
& \quad \quad stateOk \triangleq zabState[i] = \text{DISCOVERY} \\
& \quad \text{IN } \wedge infoOk \\
& \quad \quad \wedge \vee \text{1. Normal case} \\
& \quad \quad \quad \wedge epochOk \\
& \quad \quad \quad \wedge \vee \wedge stateOk \\
& \quad \quad \quad \quad \wedge \vee \wedge newEpoch > acceptedEpoch[i] \\
& \quad \quad \quad \quad \quad \wedge acceptedEpoch' = [acceptedEpoch \text{ EXCEPT } ![i] = newEpoch] \\
& \quad \quad \quad \quad \quad \wedge \text{LET } epochBytes \triangleq currentEpoch[i] \\
& \quad \quad \quad \quad \quad \quad m \triangleq [mtype \mapsto \text{ACKEPOCH}, \\
& \quad \quad \quad \quad \quad \quad \quad mzxid \mapsto lastProcessed[i].zxid, \\
& \quad \quad \quad \quad \quad \quad \quad mepoch \mapsto epochBytes] \\
& \quad \quad \quad \quad \quad \text{IN } Reply(i, j, m) \\
& \quad \quad \quad \quad \vee \wedge newEpoch = acceptedEpoch[i] \\
& \quad \quad \quad \quad \quad \wedge \text{LET } m \triangleq [mtype \mapsto \text{ACKEPOCH}, \\
& \quad \quad \quad \quad \quad \quad \quad mzxid \mapsto lastProcessed[i].zxid, \\
& \quad \quad \quad \quad \quad \quad \quad mepoch \mapsto -1] \\
& \quad \quad \quad \quad \quad \text{IN } Reply(i, j, m) \\
& \quad \quad \quad \quad \quad \wedge \text{UNCHANGED } acceptedEpoch \\
& \quad \quad \quad \quad \quad \wedge zabState' = [zabState \text{ EXCEPT } ![i] = \text{SYNCHRONIZATION}] \\
& \quad \quad \quad \quad \quad \wedge \text{UNCHANGED } violatedInvariants \\
& \quad \quad \quad \vee \wedge \neg stateOk \\
& \quad \quad \quad \quad \wedge \text{PrintT}(\text{"Exception: Follower receives LEADERINFO,"} \circ \\
& \quad \quad \quad \quad \quad \text{" whileZabState not DISCOVERY."}) \\
& \quad \quad \quad \quad \wedge violatedInvariants' = [violatedInvariants \text{ EXCEPT } !.stateInconsistent = \text{TRUE}] \\
& \quad \quad \quad \quad \wedge Discard(j, i) \\
& \quad \quad \quad \quad \wedge \text{UNCHANGED } \langle acceptedEpoch, zabState \rangle \\
& \quad \quad \quad \wedge \text{UNCHANGED } \langle varsL, leaderAddr, learners, forwarding, electing \rangle \\
& \quad \vee \text{2. Abnormal case - go back to election} \\
& \quad \quad \wedge \neg epochOk \\
& \quad \quad \wedge \text{FollowerShutdown}(i) \\
& \quad \quad \wedge \text{Clean}(i, leaderAddr[i]) \\
& \quad \quad \wedge \text{RemoveLearner}(leaderAddr[i], i) \\
& \quad \quad \wedge \text{UNCHANGED } \langle acceptedEpoch, violatedInvariants \rangle \\
& \quad \wedge \text{UNCHANGED } \langle history, lastCommitted, connecting, ackldRecv, tempMaxEpoch, \\
& \quad \quad \quad initialHistory, proposalMsgsLog, epochLeader, packetsSync \rangle
\end{aligned}$$

$\wedge \text{UpdateRecorder}(\langle \text{"FollowerProcessLEADERINFO"}, i, j \rangle)$

```

RECURSIVE UpdateAckSidHelper(-, -, -, -)
UpdateAckSidHelper(his, cur, end, target)  $\triangleq$ 
  IF cur > end THEN his
  ELSE LET curTxn  $\triangleq$  [zxid  $\mapsto$  his[1].zxid,
                       value  $\mapsto$  his[1].value,
                       ackSid  $\mapsto$  IF target  $\in$  his[1].ackSid THEN his[1].ackSid
                                   ELSE his[1].ackSid  $\cup$  {target},
                       epoch  $\mapsto$  his[1].epoch]
  IN  $\langle$ curTxn $\rangle \circ \text{UpdateAckSidHelper}(\text{Tail}(\text{his}), \text{cur} + 1, \text{end}, \text{target})$ 

```

There originally existed one bug in *LeaderProcessACK* when *monotonicallyInc* = FALSE, and it is we did not add *ackSid* of history in SYNC. So we update *ackSid* in *syncFollower*.

```

UpdateAckSid(his, lastSeenIndex, target)  $\triangleq$ 
  IF Len(his) = 0  $\vee$  lastSeenIndex = 0 THEN his
  ELSE UpdateAckSidHelper(his, 1, Minimum({Len(his), lastSeenIndex}), target)

```

return -1: this *zxid* appears at least twice; *Len*(*his*) + 1: does not exist;
 1 \neg *Len*(*his*): exists and appears just once.

```

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

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

```

Find index *idx* which meets:

history[*idx*].*zxid* \leq *zxid* < *history*[*idx* + 1].*zxid*

```

RECURSIVE IndexOfZxidHelper(-, -, -, -)
IndexOfZxidHelper(his, zxid, cur, end)  $\triangleq$ 
  IF cur > end THEN end
  ELSE IF ZxidCompare(his[cur].zxid, zxid) THEN cur - 1
  ELSE IndexOfZxidHelper(his, zxid, cur + 1, end)

```

```

IndexOfZxid(his, zxid)  $\triangleq$  IF Len(his) = 0 THEN 0

```



```

ELSE LET  $idx \triangleq ZxidToIndex(his, zxid)$ 
       $len \triangleq Len(his)$ 
IN
  IF  $idx \leq len$  THEN  $idx$ 
  ELSE  $IndexOfZxidHelper(his, zxid, 1, len)$ 

RECURSIVE  $queuePackets(-, -, -, -, -)$ 
 $queuePackets(queue, his, cur, committed, end) \triangleq$ 
  IF  $cur > end$  THEN  $queue$ 
  ELSE CASE  $cur > committed \rightarrow$ 
    LET  $m\_proposal \triangleq [mtype \mapsto PROPOSAL,$ 
       $mzxid \mapsto his[cur].zxid,$ 
       $mdata \mapsto his[cur].value]$ 
    IN  $queuePackets(Append(queue, m\_proposal), his, cur + 1, committed, end)$ 
  □  $cur \leq committed \rightarrow$ 
    LET  $m\_proposal \triangleq [mtype \mapsto PROPOSAL,$ 
       $mzxid \mapsto his[cur].zxid,$ 
       $mdata \mapsto his[cur].value]$ 
       $m\_commit \triangleq [mtype \mapsto COMMIT,$ 
       $mzxid \mapsto his[cur].zxid]$ 
       $newQueue \triangleq queue \circ \langle m\_proposal, m\_commit \rangle$ 
    IN  $queuePackets(newQueue, his, cur + 1, committed, end)$ 

RECURSIVE  $setPacketsForChecking(-, -, -, -, -, -)$ 
 $setPacketsForChecking(set, src, ep, his, cur, end) \triangleq$ 
  IF  $cur > end$  THEN  $set$ 
  ELSE LET  $m\_proposal \triangleq [source \mapsto src,$ 
     $epoch \mapsto ep,$ 
     $zxid \mapsto his[cur].zxid,$ 
     $data \mapsto his[cur].value]$ 
  IN  $setPacketsForChecking((set \cup \{m\_proposal\}), src, ep, his, cur + 1, end)$ 

See queueCommittedProposals in LearnerHandler and startForwarding in Leader for details.
For proposals in committedLog and toBeApplied, send  $\langle PROPOSAL, COMMIT \rangle$ . For
proposals in outstandingProposals, send PROPOSAL only.

 $StartForwarding(i, j, lastSeenZxid, lastSeenIndex, mode, needRemoveHead) \triangleq$ 
   $\wedge$  LET  $lastCommittedIndex \triangleq$  IF  $zabState[i] = BROADCAST$ 
    THEN  $lastCommitted[i].index$ 
    ELSE  $Len(initialHistory[i])$ 
     $lastProposedIndex \triangleq Len(history[i])$ 
     $queue\_origin \triangleq$  IF  $lastSeenIndex \geq lastProposedIndex$ 
      THEN  $\langle \rangle$ 
      ELSE  $queuePackets(\langle \rangle, history[i],$ 
         $lastSeenIndex + 1, lastCommittedIndex,$ 
         $lastProposedIndex)$ 
     $set\_forChecking \triangleq$  IF  $lastSeenIndex \geq lastProposedIndex$ 

```

```

      THEN {}
      ELSE setPacketsForChecking({}, i,
        acceptedEpoch[i], history[i],
        lastSeenIndex + 1, lastProposedIndex)
m_trunc  $\triangleq$  [mtype  $\mapsto$  TRUNC, mtruncZxid  $\mapsto$  lastSeenZxid]
m_diff  $\triangleq$  [mtype  $\mapsto$  DIFF, mzxid  $\mapsto$  lastSeenZxid]
newLeaderZxid  $\triangleq$   $\langle$ acceptedEpoch[i], 0 $\rangle$ 
m_newleader  $\triangleq$  [mtype  $\mapsto$  NEWLEADER,
  mzxid  $\mapsto$  newLeaderZxid]
queue_toSend  $\triangleq$  CASE mode = TRUNC  $\rightarrow$  ( $\langle$ m_trunc $\rangle$   $\circ$  queue_origin)  $\circ$   $\langle$ m_newleader $\rangle$ 
  □ OTHER  $\rightarrow$  ( $\langle$ m_diff $\rangle$   $\circ$  queue_origin)  $\circ$   $\langle$ m_newleader $\rangle$ 
IN  $\wedge \vee \wedge$  needRemoveHead
   $\wedge$  DiscardAndSendPackets(i, j, queue_toSend)
   $\vee \wedge \neg$  needRemoveHead
   $\wedge$  SendPackets(i, j, queue_toSend)
   $\wedge$  proposalMsgsLog' = proposalMsgsLog  $\cup$  set_forChecking
 $\wedge$  forwarding' = [forwarding EXCEPT ![i] = @  $\cup$  {j}]
 $\wedge$  history' = [history EXCEPT ![i] = UpdateAckSid(@, lastSeenIndex, j)]

```

Leader syncs with follower using *DIFF/TRUNC/PROPOSAL/COMMIT* ... See *syncFollower* in *LearnerHandler* for details.

```

SyncFollower(i, j, peerLastZxid, needRemoveHead)  $\triangleq$ 
  LET IsPeerNewEpochZxid  $\triangleq$  peerLastZxid[2] = 0
  lastProcessedZxid  $\triangleq$  lastProcessed[i].zxid
  maxCommittedLog  $\triangleq$  IF zabState[i] = BROADCAST
    THEN lastCommitted[i].zxid
    ELSE LET totalLen  $\triangleq$  Len(initialHistory[i])
      IN IF totalLen = 0 THEN  $\langle$ 0, 0 $\rangle$ 
      ELSE history[i][totalLen].zxid

```

Hypothesis: 1. *minCommittedLog* : *zxid* of head of history, so no SNAP.
 2. *maxCommittedLog* = *lastCommitted*, to compress state space.
 3. merge *queueCommittedProposals*, *startForwarding* and
 sending *NEWLEADER* into *StartForwarding*.

```

IN  $\vee$  case1. peerLastZxid = lastProcessedZxid
  DIFF + StartForwarding(lastProcessedZxid)
   $\wedge$  ZxidEqual(peerLastZxid, lastProcessedZxid)
   $\wedge$  StartForwarding(i, j, peerLastZxid, lastProcessed[i].index,
    DIFF, needRemoveHead)
 $\vee \wedge \neg$  ZxidEqual(peerLastZxid, lastProcessedZxid)
   $\wedge \vee$  case2. peerLastZxid > maxCommittedLog
    TRUNC + StartForwarding(maxCommittedLog)
     $\wedge$  ZxidCompare(peerLastZxid, maxCommittedLog)
     $\wedge$  LET maxCommittedIndex  $\triangleq$  IF zabState[i] = BROADCAST
      THEN lastCommitted[i].index

```

ELSE $Len(initialHistory[i])$
 IN $StartForwarding(i, j, maxCommittedLog, maxCommittedIndex, TRUNC, needRemoveHead)$
 \vee $case3. minCommittedLog \leq peerLastZxid \leq maxCommittedLog$
 $\wedge \neg ZxidCompare(peerLastZxid, maxCommittedLog)$
 $\wedge LET lastSeenIndex \triangleq ZxidToIndex(history[i], peerLastZxid)$
 $exist \triangleq \wedge lastSeenIndex \geq 0$
 $\wedge lastSeenIndex \leq Len(history[i])$
 $lastIndex \triangleq IF exist THEN lastSeenIndex$
 $ELSE IndexOfZxid(history[i], peerLastZxid)$
 $Maximum\ zxid\ that\ < peerLastZxid$
 $lastZxid \triangleq IF exist THEN peerLastZxid$
 $ELSE IF lastIndex = 0 THEN \langle 0, 0 \rangle$
 $ELSE history[i][lastIndex].zxid$
 IN
 \vee $case\ 3.1.\ peerLastZxid\ exists\ in\ history$
 $DIFF + StartForwarding$
 $\wedge exist$
 $\wedge StartForwarding(i, j, peerLastZxid, lastSeenIndex, DIFF, needRemoveHead)$
 \vee $case\ 3.2.\ peerLastZxid\ does\ not\ exist\ in\ history$
 $TRUNC + StartForwarding$
 $\wedge \neg exist$
 $\wedge StartForwarding(i, j, lastZxid, lastIndex, TRUNC, needRemoveHead)$
 we will not have case 4 where $peerLastZxid < minCommittedLog$, because $minCommittedLog$ default value is 1 in our spec.

compare state summary of two servers
 $IsMoreRecentThan(ss1, ss2) \triangleq \vee ss1.currentEpoch > ss2.currentEpoch$
 $\vee \wedge ss1.currentEpoch = ss2.currentEpoch$
 $\wedge ZxidCompare(ss1.lastZxid, ss2.lastZxid)$

$electionFinished$ in Leader
 $ElectionFinished(i, set) \triangleq \wedge i \in set$
 $\wedge IsQuorum(set)$

There may exist some follower shuts down and connects again, while it has sent $ACKEPOCH$ or updated $currentEpoch$ last time. This means sid of this follower has existed in $electingFollower$ but its $info$ is old. So we need to make sure each sid in $electingFollower$ is unique and $latest(newest)$.

$UpdateElecting(oldSet, sid, peerLastZxid, inQuorum) \triangleq$
 $LET sid_electing \triangleq \{s.sid : s \in oldSet\}$
 IN $IF sid \in sid_electing$
 $THEN LET old_info \triangleq CHOOSE info \in oldSet : info.sid = sid$

$$\begin{aligned}
& \text{follower_info} \triangleq \\
& \quad [sid \mapsto sid, \\
& \quad \text{peerLastZxid} \mapsto \text{peerLastZxid}, \\
& \quad \text{inQuorum} \mapsto (\text{inQuorum} \vee \text{old_info.inQuorum})] \\
& \text{IN } (\text{oldSet} \setminus \{\text{old_info}\}) \cup \{\text{follower_info}\} \\
& \text{ELSE LET } \text{follower_info} \triangleq \\
& \quad [sid \mapsto sid, \\
& \quad \text{peerLastZxid} \mapsto \text{peerLastZxid}, \\
& \quad \text{inQuorum} \mapsto \text{inQuorum}] \\
& \text{IN } \text{oldSet} \cup \{\text{follower_info}\}
\end{aligned}$$

$$\begin{aligned}
& \text{LeaderTurnToSynchronization}(i) \triangleq \\
& \quad \wedge \text{currentEpoch}' = [\text{currentEpoch} \text{ EXCEPT } ![i] = \text{acceptedEpoch}[i]] \\
& \quad \wedge \text{zabState}' = [\text{zabState} \text{ EXCEPT } ![i] = \text{SYNCHRONIZATION}]
\end{aligned}$$

Leader waits for receiving *ACKEPOCH* from a quorum, and check whether it has most recent state summary from them. After this, leader's *zabState* turns to *SYNCHRONIZATION*. See *waitForEpochAck* in Leader for details.

$$\begin{aligned}
& \text{LeaderProcessACKEPOCH}(i, j) \triangleq \\
& \quad \wedge \text{IsLeader}(i) \\
& \quad \wedge \text{PendingACKEPOCH}(i, j) \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \quad \text{infoOk} \triangleq \text{IsMyLearner}(i, j) \\
& \quad \quad \text{leaderStateSummary} \triangleq [\text{currentEpoch} \mapsto \text{currentEpoch}[i], \\
& \quad \quad \quad \text{lastZxid} \mapsto \text{lastProcessed}[i].\text{zxid}] \\
& \quad \quad \text{followerStateSummary} \triangleq [\text{currentEpoch} \mapsto \text{msg.mepoch}, \\
& \quad \quad \quad \text{lastZxid} \mapsto \text{msg.mzxid}] \\
& \quad \quad \text{logOk} \triangleq \text{whether follower is no more up-to-date than leader} \\
& \quad \quad \quad \neg \text{IsMoreRecentThan}(\text{followerStateSummary}, \text{leaderStateSummary}) \\
& \quad \quad \text{electing_quorum} \triangleq \{e \in \text{electing}[i] : e.\text{inQuorum} = \text{TRUE}\} \\
& \quad \quad \text{sid_electing} \triangleq \{s.\text{sid} : s \in \text{electing_quorum}\} \\
& \text{IN } \wedge \text{infoOk} \\
& \quad \wedge \vee \text{electionFinished} = \text{true, jump out of } \text{waitForEpochAck}. \\
& \quad \quad \text{Different from code, here we still need to record } \text{info} \\
& \quad \quad \text{into electing, to help us perform } \text{syncFollower} \text{ afterwards.} \\
& \quad \quad \text{Since electing already meets quorum, it does not break} \\
& \quad \quad \text{consistency between code and spec.} \\
& \quad \quad \wedge \text{ElectionFinished}(i, \text{sid_electing}) \\
& \quad \quad \wedge \text{electing}' = [\text{electing} \text{ EXCEPT } ![i] = \text{UpdateElecting}(@, j, \text{msg.mzxid}, \text{FALSE})] \\
& \quad \quad \wedge \text{Discard}(j, i) \\
& \quad \quad \wedge \text{UNCHANGED } \langle \text{varsL}, \text{zabState}, \text{forwarding}, \text{leaderAddr}, \\
& \quad \quad \quad \text{learners}, \text{epochLeader}, \text{violatedInvariants} \rangle \\
& \vee \wedge \neg \text{ElectionFinished}(i, \text{sid_electing}) \\
& \quad \wedge \vee \wedge \text{zabState}[i] = \text{DISCOVERY} \\
& \quad \quad \wedge \text{UNCHANGED } \text{violatedInvariants} \\
& \quad \vee \wedge \text{zabState}[i] \neq \text{DISCOVERY}
\end{aligned}$$

$\wedge \text{PrintT}(\text{"Exception: electionFinished false,"} \circ$
 $\quad \text{" while zabState not DISCOVERY."})$
 $\wedge \text{violatedInvariants}' = [\text{violatedInvariants} \text{ EXCEPT}$
 $\quad \quad \quad \text{!.stateInconsistent} = \text{TRUE}]$
 $\wedge \vee \wedge \text{followerStateSummary.currentEpoch} = -1$
 $\wedge \text{electing}' = [\text{electing} \text{ EXCEPT } ![i] \quad = \text{UpdateElecting}(@, j,$
 $\quad \quad \quad \text{msg.mzxid, FALSE})]$
 $\wedge \text{Discard}(j, i)$
 $\wedge \text{UNCHANGED } \langle \text{varsL, zabState, forwarding, leaderAddr,}$
 $\quad \quad \quad \text{learners, epochLeader} \rangle$
 $\vee \wedge \text{followerStateSummary.currentEpoch} > -1$
 $\wedge \vee \text{normal follower}$
 $\wedge \text{logOk}$
 $\wedge \text{electing}' = [\text{electing} \text{ EXCEPT } ![i] =$
 $\quad \quad \quad \text{UpdateElecting}(@, j, \text{msg.mzxid, TRUE})]$
 $\wedge \text{LET } \text{new_electing_quorum} \triangleq \{e \in \text{electing}'[i] : e.\text{inQuorum} = \text{TRUE}\}$
 $\quad \quad \quad \text{new_sid_electing} \triangleq \{s.\text{sid} : s \in \text{new_electing_quorum}\}$
 IN
 $\vee \text{electionFinished} = \text{true, jump out of waitForEpochAck,}$
 $\quad \quad \quad \text{update currentEpoch and zabState.}$
 $\wedge \text{ElectionFinished}(i, \text{new_sid_electing})$
 $\wedge \text{LeaderTurnToSynchronization}(i)$
 $\wedge \text{LET } \text{newLeaderEpoch} \triangleq \text{acceptedEpoch}[i]$
 $\quad \quad \quad \text{IN } \text{epochLeader}' = [\text{epochLeader} \text{ EXCEPT } ![newLeaderEpoch]$
 $\quad \quad \quad = @ \cup \{i\}] \text{ for checking invariants}$
 $\vee \text{there still exists electionFinished} = \text{false.}$
 $\wedge \neg \text{ElectionFinished}(i, \text{new_sid_electing})$
 $\wedge \text{UNCHANGED } \langle \text{currentEpoch, zabState, epochLeader} \rangle$
 $\wedge \text{Discard}(j, i)$
 $\wedge \text{UNCHANGED } \langle \text{state, lastProcessed, electionVars, leadingVoteSet,}$
 $\quad \quad \quad \text{electionMsgs, leaderAddr, learners, history, forwarding} \rangle$
 $\vee \text{Exists follower more recent than leader}$
 $\wedge \neg \text{logOk}$
 $\wedge \text{LeaderShutdown}(i)$
 $\wedge \text{UNCHANGED } \langle \text{electing, epochLeader} \rangle$
 $\wedge \text{UNCHANGED } \langle \text{acceptedEpoch, lastCommitted, connecting, ackldRecv,}$
 $\quad \quad \quad \text{tempMaxEpoch, initialHistory, packetsSync, proposalMsgsLog} \rangle$
 $\wedge \text{UpdateRecorder}(\langle \text{"LeaderProcessACKEPOCH", } i, j \rangle)$

Strip *syncFollower* from *LeaderProcessACKEPOCH*.
 Only when *electionFinished* = true and there exists some
learnerHandler has not perform *syncFollower*, this
 action will be called.

$\text{LeaderSyncFollower}(i) \triangleq$
 $\quad \wedge \text{IsLeader}(i)$

$\wedge \text{LET } \text{electing_quorum} \triangleq \{e \in \text{electing}[i] : e.\text{inQuorum} = \text{TRUE}\}$
 $\text{electionFinished} \triangleq \text{ElectionFinished}(i, \{s.\text{sid} : s \in \text{electing_quorum}\})$
 $\text{toSync} \triangleq \{s \in \text{electing}[i] : \wedge \neg \text{ZxidEqual}(s.\text{peerLastZxid}, \langle -1, -1 \rangle)$
 $\quad \wedge s.\text{sid} \in \text{learners}[i]\}$
 $\text{canSync} \triangleq \text{toSync} \neq \{\}$
 IN
 $\wedge \text{electionFinished}$
 $\wedge \text{canSync}$
 $\wedge \text{LET } \text{chosen} \triangleq \text{CHOOSE } s \in \text{toSync} : \text{TRUE}$
 $\quad \text{newChosen} \triangleq [\text{sid} \mapsto \text{chosen}.\text{sid},$
 $\quad \text{peerLastZxid} \mapsto \langle -1, -1 \rangle, \quad \langle -1, -1 \rangle \text{ means has handled.}$
 $\quad \text{inQuorum} \mapsto \text{chosen}.\text{inQuorum}]$
 IN $\wedge \text{SyncFollower}(i, \text{chosen}.\text{sid}, \text{chosen}.\text{peerLastZxid}, \text{FALSE})$
 $\wedge \text{electing}' = [\text{electing} \text{ EXCEPT } ![i] = (@ \setminus \{\text{chosen}\}) \cup \{\text{newChosen}\}]$
 $\wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{lastProcessed}, \text{zabState}, \text{acceptedEpoch},$
 $\quad \text{lastCommitted}, \text{initialHistory}, \text{electionVars}, \text{leadingVoteSet},$
 $\quad \text{learners}, \text{connecting}, \text{ackldRecv}, \text{tempMaxEpoch}, \text{followerVars},$
 $\quad \text{epochLeader}, \text{violatedInvariants}, \text{electionMsgs} \rangle$
 $\wedge \text{UpdateRecorder}(\langle \text{"LeaderSyncFollower"}, i \rangle)$
 $\text{TruncateLog}(\text{his}, \text{index}) \triangleq \text{IF } \text{index} \leq 0 \text{ THEN } \langle \rangle$
 $\quad \text{ELSE } \text{SubSeq}(\text{his}, 1, \text{index})$

Follower receives *DIFF/TRUNC*, and then may receives
PROPOSAL, COMMIT, NEWLEADER, and *UPTODATE*. See *syncWithLeader* in Learner
 for details.

$\text{FollowerProcessSyncMessage}(i, j) \triangleq$
 $\wedge \text{IsFollower}(i)$
 $\wedge \text{msgs}[j][i] \neq \langle \rangle$
 $\wedge \text{msgs}[j][i][1].\text{mtype} = \text{DIFF} \vee \text{msgs}[j][i][1].\text{mtype} = \text{TRUNC}$
 $\wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1]$
 $\quad \text{infoOk} \triangleq \text{IsMyLeader}(i, j)$
 $\quad \text{stateOk} \triangleq \text{zabState}[i] = \text{SYNCHRONIZATION}$
 IN $\wedge \text{infoOk}$
 $\wedge \vee$ Follower should receive packets in *SYNC*.
 $\quad \wedge \neg \text{stateOk}$
 $\quad \wedge \text{PrintT}(\text{"Exception: Follower receives DIFF/TRUNC,"} \circ$
 $\quad \quad \text{"whileZabState not SYNCHRONIZATION."})$
 $\quad \wedge \text{violatedInvariants}' = [\text{violatedInvariants} \text{ EXCEPT } !.\text{stateInconsistent} = \text{TRUE}]$
 $\quad \wedge \text{UNCHANGED } \langle \text{history}, \text{initialHistory}, \text{lastProcessed}, \text{lastCommitted} \rangle$
 $\vee \wedge \text{stateOk}$
 $\quad \wedge \vee \wedge \text{msg.mtype} = \text{DIFF}$
 $\quad \quad \wedge \text{UNCHANGED } \langle \text{history}, \text{initialHistory}, \text{lastProcessed}, \text{lastCommitted},$
 $\quad \quad \text{violatedInvariants} \rangle$
 $\quad \vee \wedge \text{msg.mtype} = \text{TRUNC}$
 $\quad \quad \wedge \text{LET } \text{truncZxid} \triangleq \text{msg}.\text{mtruncZxid}$

```

    truncIndex  $\triangleq$  ZxidToIndex(history[i], truncZxid)
  IN
     $\vee \wedge$  truncIndex > Len(history[i])
     $\wedge$  PrintT("Exception: TRUNC error.")
     $\wedge$  violatedInvariants' = [violatedInvariants EXCEPT
      !.proposalInconsistent = TRUE]
     $\wedge$  UNCHANGED  $\langle$ history, initialHistory, lastProcessed, lastCommitted $\rangle$ 
     $\vee \wedge$  truncIndex  $\leq$  Len(history[i])
     $\wedge$  history' = [history EXCEPT
      ! [i] = TruncateLog(history[i], truncIndex)]
     $\wedge$  initialHistory' = [initialHistory EXCEPT ! [i] = history'[i]]
     $\wedge$  lastProcessed' = [lastProcessed EXCEPT
      ! [i] = [index  $\mapsto$  truncIndex,
      zxid  $\mapsto$  truncZxid]]
     $\wedge$  lastCommitted' = [lastCommitted EXCEPT
      ! [i] = [index  $\mapsto$  truncIndex,
      zxid  $\mapsto$  truncZxid]]
     $\wedge$  UNCHANGED violatedInvariants
   $\wedge$  Discard(j, i)
   $\wedge$  UNCHANGED  $\langle$ state, currentEpoch, zabState, acceptedEpoch, electionVars,
    leaderVars, tempMaxEpoch, followerVars,
    proposalMsgsLog, epochLeader, electionMsgs $\rangle$ 
   $\wedge$  UpdateRecorder(("FollowerProcessSyncMessage", i, j))

  See lastProposed in Leader for details.
  LastProposed(i)  $\triangleq$  IF Len(history[i]) = 0 THEN [index  $\mapsto$  0,
    zxid  $\mapsto$   $\langle$ 0, 0 $\rangle$ ]
  ELSE
    LET lastIndex  $\triangleq$  Len(history[i])
    entry  $\triangleq$  history[i][lastIndex]
  IN [index  $\mapsto$  lastIndex,
    zxid  $\mapsto$  entry.zxid]

  See lastQueued in Learner for details.
  LastQueued(i)  $\triangleq$  IF  $\neg$ IsFollower(i)  $\vee$  zabState[i]  $\neq$  SYNCHRONIZATION
  THEN LastProposed(i)
  ELSE condition: IsFollower(i)  $\wedge$  zabState = SYNCHRONIZATION
    LET packetsInSync  $\triangleq$  packetsSync[i].notCommitted
    lenSync  $\triangleq$  Len(packetsInSync)
    totalLen  $\triangleq$  Len(history[i]) + lenSync
  IN IF lenSync = 0 THEN LastProposed(i)
    ELSE [index  $\mapsto$  totalLen,
      zxid  $\mapsto$  packetsInSync[lenSync].zxid]

  IsNextZxid(curZxid, nextZxid)  $\triangleq$ 
     $\vee$  first PROPOSAL in this epoch

```

$$\begin{array}{l} \wedge nextZxid[2] = 1 \\ \wedge curZxid[1] < nextZxid[1] \\ \vee \text{ not first } PROPOSAL \text{ in this epoch} \\ \wedge nextZxid[2] > 1 \\ \wedge curZxid[1] = nextZxid[1] \\ \wedge curZxid[2] + 1 = nextZxid[2] \end{array}$$

```

FollowerProcessPROPOSALInSync(i, j)  $\triangleq$ 
   $\wedge$  IsFollower(i)
   $\wedge$  PendingPROPOSAL(i, j)
   $\wedge$  zabState[i] = SYNCHRONIZATION
   $\wedge$  LET msg  $\triangleq$  msgs[j][i][1]
    infoOk  $\triangleq$  IsMyLeader(i, j)
    isNext  $\triangleq$  IsNextZxid(LastQueued(i).zxid, msg.mzxid)
    newTxn  $\triangleq$  [zxid  $\mapsto$  msg.mzxid,
                value  $\mapsto$  msg.mdata,
                ackSid  $\mapsto$  {}, follower do not consider ackSid
                epoch  $\mapsto$  acceptedEpoch[i]] epoch of this round

  IN  $\wedge$  infoOk
     $\wedge \vee \wedge$  isNext
       $\wedge$  packetsSync' = [packetsSync EXCEPT ![i].notCommitted
                        = Append(packetsSync[i].notCommitted, newTxn)]
     $\vee \wedge \neg$  isNext
       $\wedge$  PrintT("Warn: Follower receives PROPOSAL,"  $\circ$ 
                " while zxid != lastQueued + 1.")
       $\wedge$  UNCHANGED packetsSync
    logRequest  $\rightarrow$  SyncRequestProcessor  $\rightarrow$  SendAckRequestProcessor  $\rightarrow$  reply ack
    So here we do not need to send ack to leader.
   $\wedge$  Discard(j, i)
   $\wedge$  UNCHANGED {serverVars, electionVars, leaderVars, leaderAddr,
                verifyVars, electionMsgs}
   $\wedge$  UpdateRecorder(("FollowerProcessPROPOSALInSync", i, j))

RECURSIVE IndexOfFirstTxnWithEpoch(–, –, –, –)
IndexOfFirstTxnWithEpoch(his, epoch, cur, end)  $\triangleq$ 
  IF cur > end THEN cur
  ELSE IF his[cur].epoch = epoch THEN cur
  ELSE IndexOfFirstTxnWithEpoch(his, epoch, cur + 1, end)

LastCommitted(i)  $\triangleq$  IF zabState[i] = BROADCAST THEN lastCommitted[i]
  ELSE CASE IsLeader(i)  $\rightarrow$ 
    LET lastInitialIndex  $\triangleq$  Len(initialHistory[i])
    IN IF lastInitialIndex = 0 THEN [index  $\mapsto$  0,
                                     zxid  $\mapsto$  <0, 0>]
    ELSE [index  $\mapsto$  lastInitialIndex,
          zxid  $\mapsto$  history[i][lastInitialIndex].zxid]

```



```

□ IsFollower(i) →
  LET completeHis  $\triangleq$  history[i] ∘ packetsSync[i].notCommitted
    packetsCommitted  $\triangleq$  packetsSync[i].committed
    lenCommitted  $\triangleq$  Len(packetsCommitted)
  IN IF lenCommitted = 0 return last one in initial history
    THEN LET lastInitialIndex  $\triangleq$  Len(initialHistory[i])
      IN IF lastInitialIndex = 0
        THEN [index ↦ 0,
              zxid ↦ ⟨0, 0⟩]
        ELSE [index ↦ lastInitialIndex,
              zxid ↦ completeHis[lastInitialIndex].zxid]
      ELSE return tail of packetsCommitted
    LET committedIndex  $\triangleq$  ZxidToIndex(completeHis,
                                         packetsCommitted[lenCommitted])
    IN [index ↦ committedIndex,
        zxid ↦ packetsCommitted[lenCommitted]]
□ OTHER → lastCommitted[i]

```

```

TxnWithIndex(i, idx)  $\triangleq$  IF  $\neg$ IsFollower(i) ∨ zabState[i] ≠ SYNCHRONIZATION
  THEN history[i][idx]
  ELSE LET completeHis  $\triangleq$  history[i] ∘ packetsSync[i].notCommitted
    IN completeHis[idx]

```

To simplify specification, we assume *snapshotNeeded* = false and *writeToTxnLog* = true. So here we just call *packetsCommitted.add*.

```

FollowerProcessCOMMITInSync(i, j)  $\triangleq$ 
  ∧ IsFollower(i)
  ∧ PendingCOMMIT(i, j)
  ∧ zabState[i] = SYNCHRONIZATION
  ∧ LET msg  $\triangleq$  msgs[j][i][1]
    infoOk  $\triangleq$  IsMyLeader(i, j)
    committedIndex  $\triangleq$  LastCommitted(i).index + 1
    exist  $\triangleq$  ∧ committedIndex ≤ LastQueued(i).index
      ∧ IsNextZxid(LastCommitted(i).zxid, msg.mzxid)
    match  $\triangleq$  ZxidEqual(msg.mzxid, TxnWithIndex(i, committedIndex).zxid)
  IN ∧ infoOk
    ∧ ∨ ∧ exist
      ∧ ∨ ∧ match
        ∧ packetsSync' = [packetsSync EXCEPT !i].committed
          = Append(packetsSync[i].committed, msg.mzxid)
        ∧ UNCHANGED violatedInvariants
      ∨ ∧  $\neg$ match
        ∧ PrintT(“Warn: Follower receives COMMIT,” ∘
          “ but zxid not the next committed zxid in COMMIT.”)
        ∧ violatedInvariants' = [violatedInvariants EXCEPT

```

```

    !.commitInconsistent = TRUE]
    ∧ UNCHANGED packetsSync
  ∨ ∧ ¬exist
    ∧ PrintT("Warn: Follower receives COMMIT," ∘
      " but no packets with its zxid exists.")
    ∧ violatedInvariants' = [violatedInvariants EXCEPT
      !.commitInconsistent = TRUE]
    ∧ UNCHANGED packetsSync
  ∧ Discard(j, i)
  ∧ UNCHANGED ⟨serverVars, electionVars, leaderVars,
    leaderAddr, epochLeader, proposalMsgsLog, electionMsgs⟩
  ∧ UpdateRecorder(("FollowerProcessCOMMITInSync", i, j))

RECURSIVE ACKInBatches(-, -)
ACKInBatches(queue, packets) ≜
  IF packets = ⟨⟩ THEN queue
  ELSE LET head ≜ packets[1]
    newPackets ≜ Tail(packets)
    m_ack ≜ [mtype ↦ ACK,
      mzxid ↦ head.zxid]
  IN ACKInBatches(Append(queue, m_ack), newPackets)

Update currentEpoch, and logRequest every packets in packetsNotCommitted and clear it. As
syncProcessor will be called in logRequest, we have to reply acks here.
FollowerProcessNEWLEADER(i, j) ≜
  ∧ IsFollower(i)
  ∧ PendingNEWLEADER(i, j)
  ∧ LET msg ≜ msgs[j][i][1]
    infoOk ≜ IsMyLeader(i, j)
    packetsInSync ≜ packetsSync[i].notCommitted
    m_ackld ≜ [mtype ↦ ACKLD,
      mzxid ↦ msg.mzxid]
    ms_ack ≜ ACKInBatches(⟨⟩, packetsInSync)
    queue_toSend ≜ ⟨m_ackld⟩ ∘ ms_ack send ACK – NEWLEADER first.
  IN
    ∧ infoOk
    ∧ currentEpoch' = [currentEpoch EXCEPT ![i] = acceptedEpoch[i]]
    ∧ history' = [history EXCEPT ![i] = @ ∘ packetsInSync]
    ∧ packetsSync' = [packetsSync EXCEPT ![i].notCommitted = ⟨⟩]
    ∧ DiscardAndSendPackets(i, j, queue_toSend)
  ∧ UNCHANGED ⟨state, lastProcessed, zabState, acceptedEpoch, lastCommitted,
    electionVars, leaderVars, initialHistory, leaderAddr, verifyVars,
    electionMsgs⟩
  ∧ UpdateRecorder(("FollowerProcessNEWLEADER", i, j))

quorumFormed in Leader
QuorumFormed(i) ≜ i ∈ ackldRecv[i] ∧ IsQuorum(ackldRecv[i])

```

$QuorumFormedTurnToTrue(i) \triangleq i \in ackldRecv'[i] \wedge IsQuorum(ackldRecv'[i])$

$UpdateElectionVote(i, epoch) \triangleq UpdateProposal(i, currentVote[i].proposedLeader, currentVote[i].proposedZxid, epoch)$

See *startZkServer* in Leader for details.

$StartZkServer(i) \triangleq$
 LET $latest \triangleq LastProposed(i)$
 IN $\wedge lastCommitted' = [lastCommitted \text{ EXCEPT } ![i] = latest]$
 $\wedge lastProcessed' = [lastProcessed \text{ EXCEPT } ![i] = latest]$
 $\wedge UpdateElectionVote(i, acceptedEpoch[i])$

$LeaderTurnToBroadcast(i) \triangleq$
 $\wedge StartZkServer(i)$
 $\wedge zabState' = [zabState \text{ EXCEPT } ![i] = BROADCAST]$

Leader waits for receiving quorum of *ACK* whose lower bits of *zxid* is 0, and broadcasts *UPTODATE*. See *waitForNewLeaderAck* for details.

$LeaderProcessACKLD(i, j) \triangleq$
 $\wedge IsLeader(i)$
 $\wedge PendingACKLD(i, j)$
 \wedge LET $msg \triangleq msgs[j][i][1]$
 $infoOk \triangleq IsMyLearner(i, j)$
 $match \triangleq ZxidEqual(msg.mzxid, \langle acceptedEpoch[i], 0 \rangle)$
 $currentZxid \triangleq \langle acceptedEpoch[i], 0 \rangle$
 $m_uptodate \triangleq [mtype \mapsto UPTODATE,$
 $mzxid \mapsto currentZxid]$ not important
 IN $\wedge infoOk$
 $\wedge \vee$ just reply *UPTODATE*.
 $\wedge QuorumFormed(i)$
 $\wedge Reply(i, j, m_uptodate)$
 \wedge UNCHANGED $\langle ackldRecv, zabState, lastCommitted, lastProcessed,$
 $currentVote, violatedInvariants \rangle$
 $\vee \wedge \neg QuorumFormed(i)$
 $\wedge \vee \wedge match$
 $\wedge ackldRecv' = [ackldRecv \text{ EXCEPT } ![i] = @ \cup \{j\}]$
 $\wedge \vee$ jump out of *waitForNewLeaderAck*, and do *startZkServer*,
 $setZabState$, and reply *UPTODATE*.
 $\wedge QuorumFormedTurnToTrue(i)$
 $\wedge LeaderTurnToBroadcast(i)$
 $\wedge DiscardAndBroadcastUPTODATE(i, j, m_uptodate)$
 \vee still wait in *waitForNewLeaderAck*.
 $\wedge \neg QuorumFormedTurnToTrue(i)$
 $\wedge Discard(j, i)$
 \wedge UNCHANGED $\langle zabState, lastCommitted, lastProcessed, currentVote \rangle$
 \wedge UNCHANGED *violatedInvariants*

```

    ∨ ∧ ¬match
    ∧ PrintT("Exception: NEWLEADER ACK is from a different epoch. ")
    ∧ violatedInvariants' = [violatedInvariants EXCEPT
        !.ackInconsistent = TRUE]
    ∧ Discard(j, i)
    ∧ UNCHANGED ⟨ackldRecv, zabState, lastCommitted,
        lastProcessed, currentVote⟩
    ∧ UNCHANGED ⟨state, currentEpoch, acceptedEpoch, history, logicalClock, receiveVotes,
        outOfElection, recvQueue, waitNotmsg, leadingVoteSet, learners, connecting,
        electing, forwarding, tempMaxEpoch, initialHistory, followerVars,
        proposalMsgsLog, epochLeader, electionMsgs⟩
    ∧ UpdateRecorder(("LeaderProcessACKLD", i, j))

TxnsWithPreviousEpoch(i) ≜
    LET completeHis ≜ IF ¬IsFollower(i) ∨ zabState[i] ≠ SYNCHRONIZATION
        THEN history[i]
        ELSE history[i] ∘ packetsSync[i].notCommitted
    end ≜ Len(completeHis)
    first ≜ IndexOfFirstTxnWithEpoch(completeHis, acceptedEpoch[i], 1, end)
    IN IF first > end THEN completeHis
        ELSE SubSeq(completeHis, 1, first - 1)

TxnsRcvWithCurEpoch(i) ≜
    LET completeHis ≜ IF ¬IsFollower(i) ∨ zabState[i] ≠ SYNCHRONIZATION
        THEN history[i]
        ELSE history[i] ∘ packetsSync[i].notCommitted
    end ≜ Len(completeHis)
    first ≜ IndexOfFirstTxnWithEpoch(completeHis, acceptedEpoch[i], 1, end)
    IN IF first > end THEN ⟨⟩
        ELSE SubSeq(completeHis, first, end) completeHis[first : end]

Txns received in current epoch but not committed.
See pendingTxns in FollowerZooKeeper for details.
PendingTxns(i) ≜ IF ¬IsFollower(i) ∨ zabState[i] ≠ SYNCHRONIZATION
    THEN SubSeq(history[i], lastCommitted[i].index + 1, Len(history[i]))
    ELSE LET packetsCommitted ≜ packetsSync[i].committed
        completeHis ≜ history[i] ∘ packetsSync[i].notCommitted
    IN IF Len(packetsCommitted) = 0
        THEN SubSeq(completeHis, Len(initialHistory[i]) + 1, Len(completeHis))
        ELSE SubSeq(completeHis, LastCommitted(i).index + 1, Len(completeHis))

CommittedTxns(i) ≜ IF ¬IsFollower(i) ∨ zabState[i] ≠ SYNCHRONIZATION
    THEN SubSeq(history[i], 1, lastCommitted[i].index)
    ELSE LET packetsCommitted ≜ packetsSync[i].committed
        completeHis ≜ history[i] ∘ packetsSync[i].notCommitted
    IN IF Len(packetsCommitted) = 0 THEN initialHistory[i]

```

ELSE $SubSeq(completeHis, 1, LastCommitted(i).index)$

Each $zxid$ of $packetsCommitted$ equals to $zxid$ of
corresponding txn in $txns$.

RECURSIVE $TxnsAndCommittedMatch(-, -)$

$TxnsAndCommittedMatch(txns, packetsCommitted) \triangleq$

LET $len1 \triangleq Len(txns)$

$len2 \triangleq Len(packetsCommitted)$

IN IF $len2 = 0$ THEN TRUE

ELSE IF $len1 < len2$ THEN FALSE

ELSE $\wedge ZxidEqual(txns[len1].zxid, packetsCommitted[len2])$

$\wedge TxnsAndCommittedMatch(SubSeq(txns, 1, len1 - 1),$

$SubSeq(packetsCommitted, 1, len2 - 1))$

$FollowerLogRequestInBatches(i, leader, ms_ack, packetsNotCommitted) \triangleq$

$\wedge history' = [history \text{ EXCEPT } ![i] = @ \circ packetsNotCommitted]$

$\wedge DiscardAndSendPackets(i, leader, ms_ack)$

Since $commit$ will call $commitProcessor.commit$, which will finally
update $lastProcessed$, we update it here atomically.

$FollowerCommitInBatches(i) \triangleq$

LET $committedTxns \triangleq CommittedTxns(i)$

$packetsCommitted \triangleq packetsSync[i].committed$

$match \triangleq TxnsAndCommittedMatch(committedTxns, packetsCommitted)$

IN

$\vee \wedge match$

$\wedge lastCommitted' = [lastCommitted \text{ EXCEPT } ![i] = LastCommitted(i)]$

$\wedge lastProcessed' = [lastProcessed \text{ EXCEPT } ![i] = lastCommitted'[i]]$

$\wedge \text{UNCHANGED } violatedInvariants$

$\vee \wedge \neg match$

$\wedge PrintT(\text{"Warn: Committing zxid without see txn. /"} \circ$

$\text{"Committing zxid != pending txn zxid."})$

$\wedge violatedInvariants' = [violatedInvariants \text{ EXCEPT }$

$!.commitInconsistent = \text{TRUE}]$

$\wedge \text{UNCHANGED } \langle lastCommitted, lastProcessed \rangle$

Follower jump out of $outerLoop$ here, and log the stuff that came in between snapshot and
uptodate, which means calling $logRequest$ and $commit$ to clear $packetsNotCommitted$ and
 $packetsCommitted$.

$FollowerProcessUPTODATE(i, j) \triangleq$

$\wedge IsFollower(i)$

$\wedge PendingUPTODATE(i, j)$

$\wedge \text{LET } msg \triangleq msgs[j][i][1]$

$infoOk \triangleq IsMyLeader(i, j)$

$packetsNotCommitted \triangleq packetsSync[i].notCommitted$

$ms_ack \triangleq ACKInBatches(\langle \rangle, packetsNotCommitted)$

Try to commit one operation, called by *LeaderProcessAck*.

See *tryToCommit* in Leader for details.

commitProcessor.commit \rightarrow *processWrite* \rightarrow *toBeApplied.processRequest*
 \rightarrow *finalProcessor.processRequest*, finally *processTxn* will be implemented

and *lastProcessed* will be updated. So we update it here.

$$LeaderTryToCommit(s, index, zxid, newTxn, follower) \triangleq$$
$$\text{LET } allTrnsBeforeCommitted \triangleq lastCommitted[s].index \geq index - 1$$

Only when all proposals before *xxid* has been committed, this proposal can be permitted to be committed.

$$hasAllQuorums \triangleq IsQuorum(newTxn.ackSid)$$

In order to be committed, a proposal must be accepted by a quorum.

$$ordered \stackrel{\Delta}{=} \overline{lastCommitted[s].index} + 1 = index$$

Commit proposals in order.

IN	$\vee \wedge$	Current conditions do not satisfy committing the proposal.
----	---------------	--

$$\vee \neg allTxnsBeforeCommitted$$
$$\vee \neg hasAllQuorums$$
$$\wedge Discard(follower, s)$$
$$\wedge \text{UNCHANGED } \langle \textit{violatedInvariants}, \textit{lastCommitted}, \textit{lastProcessed} \rangle$$
$$\vee \wedge allTxnsBeforeCommitted$$
 $\wedge hasAllQuorums$ $\wedge \vee \wedge \neg$ ordered
$$\wedge PrintT(\text{"Warn: Committing xzid "} \circ ToString(xzid) \circ \text{" not first."})$$
$$\wedge \textit{violatedInvariants}' = [\textit{violatedInvariants} \text{ EXCEPT}$$

```
!.commitInconsistent = TRUE]
```

 $\vee \wedge \textit{ordered}$

\wedge UNCHANGED *violatedInvariants*

$$\wedge LeaderCommit(s, follower, index, zxid)$$
$$\wedge lastProcessed' = [lastProcessed \text{ EXCEPT } ![s] = [index \mapsto index, \\ \phantom{\wedge lastProcessed' = [lastProcessed \text{ EXCEPT } ![s] = [} xid \mapsto xid]]$$

Leader Keeps a count of acks for a particular proposal, and try to commit the proposal. See case *Leader.ACK* in *LearnerHandler*, *processRequest* in *AckRequestProcessor*, and *processAck* in Leader for details.

$$LeaderProcessACK(i, j) \triangleq$$
$$\wedge IsLeader(i)$$
$$\wedge PendingACK(i, j)$$
$$\wedge \text{ LET } msg \triangleq msgs[j][i][1]$$
$$infoOk \triangleq IsMyLearner(i, j)$$
$$outstanding \triangleq LastCommitted(i).index < LastProposed(i).index$$

outstandingProposals not null

$$hasCommitted \stackrel{\Delta}{=} \neg ZxidCompare(msg.mzxid, LastCommitted(i).zxid)$$

namely, $lastCommitted \geq zxid$

$$index \triangleq \text{ZxidToIndex}(\text{history}[i], \text{msg.mzxid})$$
$$exist \triangleq index \geq 1 \wedge index \leq LastProposed(i).index$$

$$ackIndex \triangleq \text{the proposal exists in history} \text{ LastAckIndexFromFollower}(i, j)$$

$$monotonicallyInc \triangleq \vee ackIndex = -1$$

$$\vee ackIndex + 1 = index$$

$$IN \quad \text{TCP makes everytime } ackIndex \text{ should just increase by } 1$$

$$IN \quad \wedge infoOk$$

$$\wedge \vee \wedge exist$$

$$\wedge \text{monotonicallyInc}$$

$$\wedge LET \quad txn \triangleq history[i][index]$$

$$txnAfterAddAck \triangleq [zxid \mapsto txn.zxid,$$

$$value \mapsto txn.value,$$

$$ackSid \mapsto txn.ackSid \cup \{j\},$$

$$epoch \mapsto txn.epoch]$$

$$IN \quad p.addAck(sid)$$

$$\wedge history' = [history \text{ EXCEPT } ![i][index] = txnAfterAddAck]$$

$$\wedge \vee \wedge \text{Note: outstanding is 0.}$$

$$/ \text{ proposal has already been committed.}$$

$$\vee \neg outstanding$$

$$\vee hasCommitted$$

$$\wedge Discard(j, i)$$

$$\wedge UNCHANGED \langle violatedInvariants, lastCommitted, lastProcessed \rangle$$

$$\vee \wedge outstanding$$

$$\wedge \neg hasCommitted$$

$$\wedge LeaderTryToCommit(i, index, msg.mzxid, txnAfterAddAck, j)$$

$$\vee \wedge \vee \neg exist$$

$$\vee \neg \text{monotonicallyInc}$$

$$\wedge PrintT(\text{"Exception: No such zxid. " } \circ$$

$$\text{" / ackIndex doesn't inc monotonically."})$$

$$\wedge violatedInvariants' = [violatedInvariants$$

$$\text{EXCEPT } !.ackInconsistent = TRUE]$$

$$\wedge Discard(j, i)$$

$$\wedge UNCHANGED \langle history, lastCommitted, lastProcessed \rangle$$

$$\wedge UNCHANGED \langle state, currentEpoch, zabState, acceptedEpoch, electionVars,$$

$$leaderVars, initialHistory, followerVars, proposalMsgsLog, epochLeader,$$

$$electionMsgs \rangle$$

$$\wedge UpdateRecorder(\langle \text{"LeaderProcessACK"}, i, j \rangle)$$

Follower processes COMMIT in BROADCAST. See processPacket in Follower for details.

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

$$\wedge IsFollower(i)$$

$$\wedge PendingCOMMIT(i, j)$$

$$\wedge zabState[i] = BROADCAST$$

$$\wedge LET \quad msg \triangleq msgs[j][i][1]$$

$$infoOk \triangleq IsMyLeader(i, j)$$

$$pendingTxns \triangleq PendingTxns(i)$$

$$\begin{aligned}
& \wedge \vee msg.mtype = FOLLOWERINFO \\
& \vee msg.mtype = ACKEPOCH \\
& \vee msg.mtype = ACKLD \\
& \vee msg.mtype = ACK \\
\vee \wedge IsFollower(i) \\
& \wedge LET infoOk \triangleq IsMyLeader(i, j) \\
& IN \\
& \vee msg.mtype = FOLLOWERINFO \\
& \vee msg.mtype = ACKEPOCH \\
& \vee msg.mtype = ACKLD \\
& \vee msg.mtype = ACK \\
& \vee \wedge \neg infoOk \\
& \wedge \vee msg.mtype = LEADERINFO \\
& \vee msg.mtype = NEWLEADER \\
& \vee msg.mtype = UPTODATE \\
& \vee msg.mtype = PROPOSAL \\
& \vee msg.mtype = COMMIT \\
& \vee IsLooking(i) \\
& \wedge Discard(j, i) \\
& \wedge violatedInvariants' = [violatedInvariants \text{ EXCEPT } !.messageIllegal = \text{TRUE}] \\
& \wedge \text{UNCHANGED } \langle serverVars, electionVars, leaderVars, \\
& \quad followerVars, proposalMsgsLog, epochLeader, electionMsgs \rangle \\
& \wedge UnchangeRecorder
\end{aligned}$$

Defines how the variables may transition.

$Next \triangleq$

FLE module

$$\begin{aligned}
& \vee \exists i, j \in Server : FLEReceiveNotmsg(i, j) \\
& \vee \exists i \in Server : FLENotmsgTimeout(i) \\
& \vee \exists i \in Server : FLEHandleNotmsg(i) \\
& \vee \exists i \in Server : FLEWaitNewNotmsg(i) \\
& \vee \exists i \in Server : FLEWaitNewNotmsgEnd(i)
\end{aligned}$$

Some conditions like failure, network delay

$$\begin{aligned}
& \vee \exists i \in Server : FollowerTimeout(i) \\
& \vee \exists i \in Server : LeaderTimeout(i) \\
& \vee \exists i, j \in Server : Timeout(i, j)
\end{aligned}$$

Zab module - Discovery and Synchronization part

$$\begin{aligned}
& \vee \exists i, j \in Server : ConnectAndFollowerSendFOLLOWERINFO(i, j) \\
& \vee \exists i, j \in Server : LeaderProcessFOLLOWERINFO(i, j) \\
& \vee \exists i, j \in Server : FollowerProcessLEADERINFO(i, j) \\
& \vee \exists i, j \in Server : LeaderProcessACKEPOCH(i, j) \\
& \vee \exists i \in Server : LeaderSyncFollower(i) \\
& \vee \exists i, j \in Server : FollowerProcessSyncMessage(i, j) \\
& \vee \exists i, j \in Server : FollowerProcessPROPOSALInSync(i, j) \\
& \vee \exists i, j \in Server : FollowerProcessCOMMITInSync(i, j)
\end{aligned}$$

$\forall \exists i, j \in \text{Server} : \text{FollowerProcessNEWLEADER}(i, j)$
 $\forall \exists i, j \in \text{Server} : \text{LeaderProcessACKLD}(i, j)$
 $\forall \exists i, j \in \text{Server} : \text{FollowerProcessUPTODATE}(i, j)$

Zab module – Broadcast part

$\forall \exists i \in \text{Server}, v \in \text{Value} : \text{LeaderProcessRequest}(i, v)$
 $\forall \exists i, j \in \text{Server} : \text{FollowerProcessPROPOSAL}(i, j)$
 $\forall \exists i, j \in \text{Server} : \text{LeaderProcessACK}(i, j)$ *Sync + Broadcast*
 $\forall \exists i, j \in \text{Server} : \text{FollowerProcessCOMMIT}(i, j)$

An action used to judge whether there are redundant messages in network

$\forall \exists i \in \text{Server} : \text{FilterNonexistentMessage}(i)$

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

Define safety properties of Zab 1.0 protocol.

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

There is most one established leader for a certain epoch.

$\text{Leadership1} \triangleq \forall i, j \in \text{Server} :$

$\wedge \text{IsLeader}(i) \wedge \text{zabState}[i] \in \{\text{SYNCHRONIZATION}, \text{BROADCAST}\}$
 $\wedge \text{IsLeader}(j) \wedge \text{zabState}[j] \in \{\text{SYNCHRONIZATION}, \text{BROADCAST}\}$
 $\wedge \text{acceptedEpoch}[i] = \text{acceptedEpoch}[j]$
 $\Rightarrow i = j$

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

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

$\text{PrefixConsistency} \triangleq \forall i, j \in \text{Server} :$

$\text{LET } \text{smaller} \triangleq \text{Minimum}(\{\text{lastCommitted}[i].\text{index}, \text{lastCommitted}[j].\text{index}\})$
 $\text{IN } \quad \vee \text{smaller} = 0$
 $\quad \vee \wedge \text{smaller} > 0$
 $\quad \wedge \forall \text{index} \in 1 \dots \text{smaller} :$
 $\quad \text{TxnEqual}(\text{history}[i][\text{index}], \text{history}[j][\text{index}])$

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

$\text{Integrity} \triangleq \forall i \in \text{Server} :$

$\wedge \text{IsFollower}(i)$
 $\wedge \text{lastCommitted}[i].\text{index} > 0$
 $\Rightarrow \forall \text{idx} \in 1 \dots \text{lastCommitted}[i].\text{index} : \exists \text{proposal} \in \text{proposalMsgsLog} :$
 $\quad \text{ZxidEqual}(\text{history}[i][\text{idx}].\text{zxid}, \text{proposal}.\text{zxid})$

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

$\text{Agreement} \triangleq \forall i, j \in \text{Server} :$

$\wedge \text{IsFollower}(i) \wedge \text{lastCommitted}[i].\text{index} > 0$
 $\wedge \text{IsFollower}(j) \wedge \text{lastCommitted}[j].\text{index} > 0$

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

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

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

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

$$\begin{aligned}
\text{LocalPrimaryOrder} &\triangleq \text{LET } p_set(i, e) \triangleq \{p \in \text{proposalMsgsLog} : \wedge p.\text{source} = i \\
&\quad \wedge p.\text{epoch} = e\} \\
&\quad \text{zxid_set}(i, e) \triangleq \{p.\text{zxid} : p \in p_set(i, e)\} \\
\text{IN } \forall i \in \text{Server} : \forall e \in 1 \dots \text{currentEpoch}[i] : \\
&\quad \vee \text{Cardinality}(\text{zxid_set}(i, e)) < 2 \\
&\quad \vee \wedge \text{Cardinality}(\text{zxid_set}(i, e)) \geq 2 \\
&\quad \wedge \exists \text{zxid1}, \text{zxid2} \in \text{zxid_set}(i, e) : \\
&\quad \vee \text{ZxidEqual}(\text{zxid1}, \text{zxid2}) \\
&\quad \vee \wedge \neg \text{ZxidEqual}(\text{zxid1}, \text{zxid2}) \\
&\quad \wedge \text{LET } \text{zxidPre} \triangleq \text{IF } \text{ZxidCompare}(\text{zxid1}, \text{zxid2}) \text{ THEN } \text{zxid2} \text{ ELSE } \text{zxid1} \\
&\quad \quad \text{zxidNext} \triangleq \text{IF } \text{ZxidCompare}(\text{zxid1}, \text{zxid2}) \text{ THEN } \text{zxid1} \text{ ELSE } \text{zxid2} \\
&\quad \text{IN } \forall j \in \text{Server} : \wedge \text{lastCommitted}[j].\text{index} \geq 2 \\
&\quad \quad \wedge \exists \text{idx} \in 1 \dots \text{lastCommitted}[j].\text{index} : \\
&\quad \quad \quad \text{ZxidEqual}(\text{history}[j][\text{idx}].\text{zxid}, \text{zxidNext}) \\
&\quad \Rightarrow \exists \text{idx2} \in 1 \dots \text{lastCommitted}[j].\text{index} : \\
&\quad \quad \wedge \text{ZxidEqual}(\text{history}[j][\text{idx2}].\text{zxid}, \text{zxidNext}) \\
&\quad \quad \wedge \text{idx2} > 1 \\
&\quad \quad \wedge \exists \text{idx1} \in 1 \dots (\text{idx2} - 1) : \\
&\quad \quad \quad \text{ZxidEqual}(\text{history}[j][\text{idx1}].\text{zxid}, \text{zxidPre})
\end{aligned}$$

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.

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

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.

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

\ * Modification History
\ * Last modified Sun Nov 21 21:41:22 CST 2021 by Dell
\ * Created Sat Oct 23 16:05:04 CST 2021 by Dell