
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* \ * Replaced by *recorder.nClientRequest*

Value \triangleq *Nat*

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

containing: *zxid*, value, *ackSid*, 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 *clientReuquest* * Start from 0, and increases monotonically
 when *LeaderProcessRequest* performed. To
 avoid existing two requests with same value.

Update: Remove it to *recorder.nClientRequest*.

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*, *nClientRequest*]

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

electionVars \triangleq *electionVarsL*

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

followerVars \triangleq $\langle \text{leaderAddr}, \text{packetsSync} \rangle$

verifyVars \triangleq $\langle \text{proposalMsgsLog}, \text{epochLeader}, \text{violatedInvariants} \rangle$

msgVars \triangleq $\langle \text{msgs}, \text{electionMsgs} \rangle$

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

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

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

$[proposedLeader : ServersIncNullPoint,$
 $proposedZxid : Zxid,$
 $proposedEpoch : Nat]$

$ElectionVote \triangleq$
 $[vote : Vote, round : Nat, state : ElectionState, version : Nat]$

$ElectionMsg \triangleq$
 $[mtype : \{NOTIFICATION\},$
 $msource : Server,$
 $mstate : ElectionState,$
 $mround : Nat,$
 $mvote : Vote] \cup$
 $[mtype : \{NONE\}]$

$TypeOK \triangleq$
 $\wedge zabState \in [Server \rightarrow ZabState]$
 $\wedge acceptedEpoch \in [Server \rightarrow Nat]$
 $\wedge lastCommitted \in [Server \rightarrow LastItem]$
 $\wedge learners \in [Server \rightarrow SUBSET Server]$
 $\wedge connecting \in [Server \rightarrow SUBSET ServersIncNullPoint]$
 $\wedge electing \in [Server \rightarrow SUBSET Electing]$
 $\wedge ackldRecv \in [Server \rightarrow SUBSET ServersIncNullPoint]$
 $\wedge forwarding \in [Server \rightarrow SUBSET Server]$
 $\wedge initialHistory \in [Server \rightarrow Seq(HistoryItem)]$
 $\wedge tempMaxEpoch \in [Server \rightarrow Nat]$
 $\wedge leaderAddr \in [Server \rightarrow ServersIncNullPoint]$
 $\wedge packetsSync \in [Server \rightarrow SyncPackets]$
 $\wedge proposalMsgsLog \in SUBSET Proposal$
 $\wedge epochLeader \in [1 .. MAXEPOCH \rightarrow SUBSET Server]$
 $\wedge violatedInvariants \in [ViolationSet \rightarrow BOOLEAN]$
 $\wedge msgs \in [Server \rightarrow [Server \rightarrow Seq(Message)]]$

Fast Leader Election

$\wedge electionMsgs \in [Server \rightarrow [Server \rightarrow Seq(ElectionMsg)]]$
 $\wedge recvQueue \in [Server \rightarrow Seq(ElectionMsg)]$
 $\wedge leadingVoteSet \in [Server \rightarrow SUBSET Server]$
 $\wedge receiveVotes \in [Server \rightarrow [Server \rightarrow ElectionVote]]$
 $\wedge currentVote \in [Server \rightarrow Vote]$
 $\wedge outOfElection \in [Server \rightarrow [Server \rightarrow ElectionVote]]$
 $\wedge lastProcessed \in [Server \rightarrow LastItem]$
 $\wedge history \in [Server \rightarrow Seq(HistoryItem)]$
 $\wedge state \in [Server \rightarrow ElectionState]$
 $\wedge waitNotmsg \in [Server \rightarrow BOOLEAN]$
 $\wedge currentEpoch \in [Server \rightarrow Nat]$
 $\wedge logicalClock \in [Server \rightarrow Nat]$

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

Actions about recorder
 $GetParameter(p) \triangleq$ IF $p \in \text{DOMAIN } Parameters$ THEN $Parameters[p]$ ELSE 0
 $GetRecorder(p) \triangleq$ IF $p \in \text{DOMAIN } recorder$ THEN $recorder[p]$ 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 $pc[1] = \text{"LeaderProcessRequest"}$ THEN
LET $s \triangleq$ CHOOSE $i \in Server :$
 $\forall j \in Server : Len(history'[i]) \geq Len(history'[j])$
IN $Len(history'[s])$
ELSE $recorder["nTransaction"]$)
 $RecorderSetMaxEpoch(pc) \triangleq ("maxEpoch" :>$

$$\begin{aligned}
& \text{IF } pc[1] = \text{"LeaderProcessFOLLOWERINFO"} \text{ THEN} \\
& \quad \text{LET } s \triangleq \text{CHOOSE } i \in \text{Server} : \\
& \quad \quad \forall j \in \text{Server} : \text{acceptedEpoch}'[i] \geq \text{acceptedEpoch}'[j] \\
& \quad \text{IN } \text{acceptedEpoch}'[s] \\
& \quad \text{ELSE } \text{recorder}[\text{"maxEpoch"}]) \\
\text{RecorderSetRequests}(pc) & \triangleq (\text{"nClientRequest"} :> \\
& \quad \text{IF } pc[1] = \text{"LeaderProcessRequest"} \text{ THEN} \\
& \quad \quad \text{recorder}[\text{"nClientRequest"}] + 1 \\
& \quad \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} \\
& \quad \square \quad pc[1] = \text{"LeaderTimeout"} \rightarrow \text{RecorderIncTimeout} \\
& \quad \square \quad pc[1] = \text{"FollowerTimeout"} \rightarrow \text{RecorderIncTimeout} \\
& \quad \square \quad \text{OTHER} \rightarrow \text{RecorderGetTimeout} \\
\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}, \text{"MaxTimeoutFailures"}) \\
\text{CheckTransactionNum} & \triangleq \text{CheckParameterLimit}(\text{recorder.nTransaction}, \text{"MaxTransactionNum"}) \\
\text{CheckEpoch} & \triangleq \text{CheckParameterLimit}(\text{recorder.maxEpoch}, \text{"MaxEpoch"}) \\
\text{CheckStateConstraints} & \triangleq \text{CheckTimeout} \wedge \text{CheckTransactionNum} \wedge \text{CheckEpoch}
\end{aligned}$$

Actions about network

$$\begin{aligned}
\text{PendingFOLLOWERINFO}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{FOLLOWERINFO} \\
\text{PendingLEADERINFO}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{LEADERINFO} \\
\text{PendingACKEPOCH}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{ACKEPOCH} \\
\text{PendingNEWLEADER}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{NEWLEADER} \\
\text{PendingACKLD}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{ACKLD} \\
\text{PendingUPTODATE}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{UPTODATE} \\
\text{PendingPROPOSAL}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{PROPOSAL}
\end{aligned}$$

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

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

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

Leader broadcasts a *message* (*PROPOSAL/COMMIT*) to all other servers in *forwardingFollowers*.
 $Broadcast(i, m) \triangleq msgs' = [msgs \text{ EXCEPT } ![i] = [v \in Server \mapsto \text{IF } \wedge v \in forwarding[i]$
 $\quad \quad \quad \wedge v \neq i$
 $\quad \quad \quad \text{THEN } Append(msgs[i][v], m)$
 $\quad \quad \quad \text{ELSE } msgs[i][v]]]$

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

Leader broadcasts *LEADERINFO* to all other servers in *connectingFollowers*.
 $DiscardAndBroadcastLEADERINFO(i, j, m) \triangleq$
 $msgs' = [msgs \text{ EXCEPT } ![j][i] = Tail(msgs[j][i]),$
 $\quad \quad \quad ![i] = [v \in Server \mapsto \text{IF } \wedge v \in connecting'[i]$
 $\quad \quad \quad \wedge v \in learners[i]$
 $\quad \quad \quad \wedge v \neq i$
 $\quad \quad \quad \text{THEN } Append(msgs[i][v], m)$
 $\quad \quad \quad \text{ELSE } msgs[i][v]]]$

Leader broadcasts *UPTODATE* to all other servers in *newLeaderProposal*.
 $DiscardAndBroadcastUPTODATE(i, j, m) \triangleq$
 $msgs' = [msgs \text{ EXCEPT } ![j][i] = Tail(msgs[j][i]),$
 $\quad \quad \quad ![i] = [v \in Server \mapsto \text{IF } \wedge v \in ackldRecv'[i]$
 $\quad \quad \quad \wedge v \in learners[i]$
 $\quad \quad \quad \wedge v \neq i$
 $\quad \quad \quad \text{THEN } Append(msgs[i][v], m)$
 $\quad \quad \quad \text{ELSE } msgs[i][v]]]$

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

Shuffle input buffer.
 $Clean(i, j) \triangleq msgs' = [msgs \text{ EXCEPT } ![j][i] = \langle \rangle, ![i][j] = \langle \rangle]$

$$\begin{aligned}
CleanInputBuffer(i) &\triangleq msgs' = [s \in Server \mapsto [v \in Server \mapsto \text{IF } v = i \text{ THEN } \langle \rangle \\
&\hspace{15em} \text{ELSE } msgs[s][v]]] \\
CleanInputBufferInCluster(S) &\triangleq msgs' = [s \in Server \mapsto \\
&\hspace{15em} [v \in Server \mapsto \text{IF } v \in S \text{ THEN } \langle \rangle \\
&\hspace{15em} \text{ELSE } msgs[s][v]]]
\end{aligned}$$

Define initial values for all variables

$$\begin{aligned}
InitServerVars &\triangleq \wedge InitServerVarsL \\
&\wedge zabState = [s \in Server \mapsto ELECTION] \\
&\wedge acceptedEpoch = [s \in Server \mapsto 0] \\
&\wedge lastCommitted = [s \in Server \mapsto [index \mapsto 0, \\
&\hspace{10em} zxid \mapsto \langle 0, 0 \rangle]] \\
&\wedge initialHistory = [s \in Server \mapsto \langle \rangle] \\
InitLeaderVars &\triangleq \wedge InitLeaderVarsL \\
&\wedge learners = [s \in Server \mapsto \{\}] \\
&\wedge connecting = [s \in Server \mapsto \{\}] \\
&\wedge electing = [s \in Server \mapsto \{\}] \\
&\wedge ackldRecv = [s \in Server \mapsto \{\}] \\
&\wedge forwarding = [s \in Server \mapsto \{\}] \\
&\wedge tempMaxEpoch = [s \in Server \mapsto 0] \\
InitElectionVars &\triangleq InitElectionVarsL \\
InitFollowerVars &\triangleq \wedge leaderAddr = [s \in Server \mapsto NullPoint] \\
&\wedge packetsSync = [s \in Server \mapsto \\
&\hspace{10em} [notCommitted \mapsto \langle \rangle, \\
&\hspace{10em} committed \mapsto \langle \rangle]] \\
InitVerifyVars &\triangleq \wedge proposalMsgsLog = \{\} \\
&\wedge epochLeader = [i \in 1 \dots MAXEPOCH \mapsto \{\}] \\
&\wedge violatedInvariants = [stateInconsistent \mapsto FALSE, \\
&\hspace{10em} proposalInconsistent \mapsto FALSE, \\
&\hspace{10em} commitInconsistent \mapsto FALSE, \\
&\hspace{10em} ackInconsistent \mapsto FALSE, \\
&\hspace{10em} messageIllegal \mapsto FALSE] \\
InitMsgVars &\triangleq \wedge msgs = [s \in Server \mapsto [v \in Server \mapsto \langle \rangle]] \\
&\wedge electionMsgs = [s \in Server \mapsto [v \in Server \mapsto \langle \rangle]] \\
InitRecorder &\triangleq recorder = [nTimeout \mapsto 0, \\
&\hspace{2em} nTransaction \mapsto 0, \\
&\hspace{2em} maxEpoch \mapsto 0, \\
&\hspace{2em} pc \mapsto \langle "Init" \rangle, \\
&\hspace{2em} nClientRequest \mapsto 0] \\
Init &\triangleq \wedge InitServerVars
\end{aligned}$$

$\wedge \text{InitLeaderVars}$
 $\wedge \text{InitElectionVars}$
 $\wedge \text{InitFollowerVars}$
 $\wedge \text{InitVerifyVars}$
 $\wedge \text{InitMsgVars}$
 $\wedge \text{InitRecorder}$

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

$ZabTurnToFollowing(i) \triangleq$
 $\wedge zabState' = [zabState \text{ EXCEPT } ![i] = DISCOVERY]$
 $\wedge initialHistory' = [initialHistory \text{ EXCEPT } ![i] = history'[i]]$
 $\wedge packetsSync' = [packetsSync \text{ EXCEPT } ![i].notCommitted = \langle \rangle, \text{![i].committed} = \langle \rangle]$

Fast Leader Election

$FLEReceiveNotmsg(i, j) \triangleq$
 $\wedge \text{ReceiveNotmsg}(i, j)$
 $\wedge \text{UNCHANGED } \langle zabState, acceptedEpoch, lastCommitted, learners, connecting, \text{initialHistory}, electing, ackldRecv, forwarding, tempMaxEpoch, \text{followerVars}, verifyVars, msgs \rangle$
 $\wedge \text{UpdateRecorder}(\langle \text{"FLEReceiveNotmsg"}, i, j \rangle)$

$FLENotmsgTimeout(i) \triangleq$
 $\wedge \text{NotmsgTimeout}(i)$
 $\wedge \text{UNCHANGED } \langle zabState, acceptedEpoch, lastCommitted, learners, connecting, \text{initialHistory}, electing, ackldRecv, forwarding, tempMaxEpoch, \text{followerVars}, verifyVars, msgs \rangle$
 $\wedge \text{UpdateRecorder}(\langle \text{"FLENotmsgTimeout"}, i \rangle)$

$FLEHandleNotmsg(i) \triangleq$
 $\wedge \text{HandleNotmsg}(i)$
 $\wedge \text{LET } newState \triangleq state'[i]$
 IN
 $\vee \wedge newState = LEADING$
 $\wedge ZabTurnToLeading(i)$

$$\begin{aligned}
& \wedge \text{UNCHANGED } \text{packetsSync} \\
\vee & \wedge \text{newState} = \text{FOLLOWING} \\
& \wedge \text{ZabTurnToFollowing}(i) \\
& \wedge \text{UNCHANGED } \langle \text{learners}, \text{connecting}, \text{electing}, \text{ackldRecv}, \\
& \quad \text{forwarding}, \text{tempMaxEpoch} \rangle \\
\vee & \wedge \text{newState} = \text{LOOKING} \\
& \wedge \text{UNCHANGED } \langle \text{zabState}, \text{learners}, \text{connecting}, \text{electing}, \text{ackldRecv}, \\
& \quad \text{forwarding}, \text{tempMaxEpoch}, \text{packetsSync}, \text{initialHistory} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{lastCommitted}, \text{acceptedEpoch}, \text{leaderAddr}, \text{verifyVars}, \text{msgs} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"FLEHandleNotmsg"}, i \rangle)
\end{aligned}$$

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

$$\begin{aligned}
\text{FLEWaitNewNotmsg}(i) & \triangleq \\
& \wedge \text{WaitNewNotmsg}(i) \\
& \wedge \text{UNCHANGED } \langle \text{zabState}, \text{acceptedEpoch}, \text{lastCommitted}, \text{learners}, \text{connecting}, \\
& \quad \text{electing}, \text{ackldRecv}, \text{forwarding}, \text{tempMaxEpoch}, \text{initialHistory}, \\
& \quad \text{followerVars}, \text{verifyVars}, \text{msgs} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"FLEWaitNewNotmsg"}, i \rangle)
\end{aligned}$$

On the premise that $\text{ReceiveVotes.HasQuorums} = \text{TRUE}$,
corresponding to logic in line 1061 – 1066 in *FastLeaderElection*.

$$\begin{aligned}
\text{FLEWaitNewNotmsgEnd}(i) & \triangleq \\
& \wedge \text{WaitNewNotmsgEnd}(i) \\
& \wedge \text{LET } \text{newState} \triangleq \text{state}'[i] \\
& \text{IN} \\
& \vee \wedge \text{newState} = \text{LEADING} \\
& \quad \wedge \text{ZabTurnToLeading}(i) \\
& \quad \wedge \text{UNCHANGED } \text{packetsSync} \\
& \vee \wedge \text{newState} = \text{FOLLOWING} \\
& \quad \wedge \text{ZabTurnToFollowing}(i) \\
& \quad \wedge \text{UNCHANGED } \langle \text{learners}, \text{connecting}, \text{electing}, \text{ackldRecv}, \text{forwarding}, \\
& \quad \quad \text{tempMaxEpoch} \rangle \\
& \vee \wedge \text{newState} = \text{LOOKING} \\
& \quad \wedge \text{PrintT}(\text{"Note: New state is LOOKING in FLEWaitNewNotmsgEnd,"} \circ \\
& \quad \quad \text{" which should not happen."}) \\
& \quad \wedge \text{UNCHANGED } \langle \text{zabState}, \text{learners}, \text{connecting}, \text{electing}, \text{ackldRecv}, \\
& \quad \quad \text{forwarding}, \text{tempMaxEpoch}, \text{initialHistory}, \text{packetsSync} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{lastCommitted}, \text{acceptedEpoch}, \text{leaderAddr}, \text{verifyVars}, \text{msgs} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"FLEWaitNewNotmsgEnd"}, i \rangle)
\end{aligned}$$

$$\begin{aligned}
\text{InitialVotes} & \triangleq [\text{vote} \mapsto \text{InitialVote}, \\
& \quad \text{round} \mapsto 0, \\
& \quad \text{state} \mapsto \text{LOOKING}, \\
& \quad \text{version} \mapsto 0]
\end{aligned}$$

Equals to for every server in S , performing action *ZabTimeout*.

$$\begin{aligned}
& ZabTimeoutInCluster(S) \triangleq \\
& \quad \wedge state' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } LOOKING \text{ ELSE } state[s]] \\
& \quad \wedge lastProcessed' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } InitLastProcessed(s) \\
& \quad \quad \quad \text{ELSE } lastProcessed[s]] \\
& \quad \wedge logicalClock' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } logicalClock[s] + 1 \\
& \quad \quad \quad \text{ELSE } logicalClock[s]] \\
& \quad \wedge currentVote' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN} \\
& \quad \quad \quad [proposedLeader \mapsto s, \\
& \quad \quad \quad \quad proposedZxid \mapsto lastProcessed'[s].zxid, \\
& \quad \quad \quad \quad proposedEpoch \mapsto currentEpoch[s]] \\
& \quad \quad \quad \text{ELSE } currentVote[s]] \\
& \quad \wedge receiveVotes' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } [v \in Server \mapsto InitialVotes] \\
& \quad \quad \quad \text{ELSE } receiveVotes[s]] \\
& \quad \wedge outOfElection' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } [v \in Server \mapsto InitialVotes] \\
& \quad \quad \quad \text{ELSE } outOfElection[s]] \\
& \quad \wedge recvQueue' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } \langle [mtype \mapsto NONE] \rangle \\
& \quad \quad \quad \text{ELSE } recvQueue[s]] \\
& \quad \wedge waitNotmsg' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } FALSE \text{ ELSE } waitNotmsg[s]] \\
& \quad \wedge leadingVoteSet' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } \{ \} \text{ ELSE } leadingVoteSet[s]] \\
& \quad \wedge UNCHANGED \langle electionMsgs, currentEpoch, history \rangle \\
& \quad \wedge zabState' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } ELECTION \text{ ELSE } zabState[s]] \\
& \quad \wedge leaderAddr' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } NullPoint \text{ ELSE } leaderAddr[s]] \\
& \quad \wedge CleanInputBufferInCluster(S)
\end{aligned}$$

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

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

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

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

See *removeLearnerHandler* for details.

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

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

$$\begin{aligned}
\text{FollowerTimeout}(i) &\triangleq \\
&\wedge \text{CheckTimeout} \quad \text{test restrictions of } \text{timeout_1} \\
&\wedge \text{IsFollower}(i) \\
&\wedge \text{HasNoLeader}(i) \\
&\wedge \text{FollowerShutdown}(i) \\
&\wedge \text{CleanInputBuffer}(i) \\
&\wedge \text{UNCHANGED} \langle \text{acceptedEpoch}, \text{lastCommitted}, \text{learners}, \text{connecting}, \text{electing}, \\
&\quad \quad \quad \text{ackldRecv}, \text{forwarding}, \text{tempMaxEpoch}, \text{initialHistory}, \\
&\quad \quad \quad \text{verifyVars}, \text{packetsSync} \rangle \\
&\wedge \text{UpdateRecorder}(\langle \text{"FollowerTimeout"}, i \rangle)
\end{aligned}$$

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

$$\begin{aligned}
\text{LeaderTimeout}(i) &\triangleq \\
&\wedge \text{CheckTimeout} \quad \text{test restrictions of } \text{timeout_2} \\
&\wedge \text{IsLeader}(i) \\
&\wedge \neg \text{IsQuorum}(\text{learners}[i]) \\
&\wedge \text{LeaderShutdown}(i) \\
&\wedge \text{UNCHANGED} \langle \text{acceptedEpoch}, \text{lastCommitted}, \text{connecting}, \text{electing}, \text{ackldRecv}, \\
&\quad \quad \quad \text{tempMaxEpoch}, \text{initialHistory}, \text{verifyVars}, \text{packetsSync} \rangle \\
&\wedge \text{UpdateRecorder}(\langle \text{"LeaderTimeout"}, i \rangle)
\end{aligned}$$

Timeout between leader and follower.

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

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

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

waitingForNewEpoch in Leader

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

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.

$$\begin{aligned}
& \text{LeaderProcessFOLLOWERINFO}(i, j) \triangleq \\
& \quad \wedge \text{CheckEpoch} \quad \text{test restrictions of max epoch} \\
& \quad \wedge \text{IsLeader}(i) \\
& \quad \wedge \text{PendingFOLLOWERINFO}(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{lastAcceptedEpoch} \triangleq \text{msg.mzxid}[1] \\
& \quad \text{IN} \\
& \quad \wedge \text{infoOk} \\
& \quad \wedge \vee \quad \text{1. has not broadcast LEADERINFO} \\
& \quad \quad \wedge \text{WaitingForNewEpoch}(i) \\
& \quad \quad \wedge \vee \wedge \text{zabState}[i] = \text{DISCOVERY} \\
& \quad \quad \quad \wedge \text{UNCHANGED } \text{violatedInvariants} \\
& \quad \quad \vee \wedge \text{zabState}[i] \neq \text{DISCOVERY} \\
& \quad \quad \quad \wedge \text{PrintT}(\text{"Exception: waitingFotNewEpoch true,"} \circ \\
& \quad \quad \quad \quad \text{" while zabState not DISCOVERY."}) \\
& \quad \quad \quad \wedge \text{violatedInvariants}' = [\text{violatedInvariants} \text{ EXCEPT } !. \text{stateInconsistent} = \text{TRUE}] \\
& \quad \quad \wedge \text{tempMaxEpoch}' = [\text{tempMaxEpoch} \text{ EXCEPT } ![i] = \text{IF } \text{lastAcceptedEpoch} \geq \text{tempMaxEpoch}[i] \\
& \quad \quad \quad \quad \text{THEN } \text{lastAcceptedEpoch} + 1 \\
& \quad \quad \quad \quad \text{ELSE } @] \\
& \quad \quad \wedge \text{connecting}' = [\text{connecting} \text{ EXCEPT } ![i] = @ \cup \{j\}] \\
& \quad \quad \wedge \vee \wedge \text{WaitingForNewEpochTurnToFalse}(i) \\
& \quad \quad \quad \wedge \text{acceptedEpoch}' = [\text{acceptedEpoch} \text{ EXCEPT } ![i] = \text{tempMaxEpoch}'[i]] \\
& \quad \quad \quad \wedge \text{LET } \text{newLeaderZxid} \triangleq \langle \text{acceptedEpoch}'[i], 0 \rangle \\
& \quad \quad \quad \quad m \triangleq [\text{mtype} \mapsto \text{LEADERINFO}, \\
& \quad \quad \quad \quad \quad \text{mzxid} \mapsto \text{newLeaderZxid}]
\end{aligned}$$

$$\begin{aligned}
& \text{IN } \text{DiscardAndBroadcastLEADERINFO}(i, j, m) \\
& \vee \wedge \neg \text{WaitingForNewEpochTurnToFalse}(i) \\
& \quad \wedge \text{Discard}(j, i) \\
& \quad \wedge \text{UNCHANGED } \text{acceptedEpoch} \\
& \vee \quad \text{2. has broadcast LEADERINFO} \\
& \quad \wedge \neg \text{WaitingForNewEpoch}(i) \\
& \quad \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{LEADERINFO}, \\
& \quad \quad \text{mzxid} \mapsto \langle \text{acceptedEpoch}[i], 0 \rangle]) \\
& \quad \wedge \text{UNCHANGED } \langle \text{tempMaxEpoch}, \text{connecting}, \text{acceptedEpoch}, \text{violatedInvariants} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{lastProcessed}, \text{zabState}, \text{history}, \text{lastCommitted}, \\
& \quad \text{followerVars}, \text{electionVars}, \text{initialHistory}, \text{leadingVoteSet}, \text{learners}, \\
& \quad \text{electing}, \text{ackldRecv}, \text{forwarding}, \text{proposalMsgsLog}, \text{epochLeader}, \\
& \quad \text{electionMsgs} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"LeaderProcessFOLLOWERINFO"}, i, j \rangle)
\end{aligned}$$

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

$$\begin{aligned}
& \text{FollowerProcessLEADERINFO}(i, j) \triangleq \\
& \quad \wedge \text{IsFollower}(i) \\
& \quad \wedge \text{PendingLEADERINFO}(i, j) \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \quad \text{newEpoch} \triangleq \text{msg.mzxid}[1] \\
& \quad \quad \text{infoOk} \triangleq \text{IsMyLeader}(i, j) \\
& \quad \quad \text{epochOk} \triangleq \text{newEpoch} \geq \text{acceptedEpoch}[i] \\
& \quad \quad \text{stateOk} \triangleq \text{zabState}[i] = \text{DISCOVERY} \\
& \text{IN } \wedge \text{infoOk} \\
& \quad \wedge \vee \quad \text{1. Normal case} \\
& \quad \quad \wedge \text{epochOk} \\
& \quad \quad \wedge \vee \wedge \text{stateOk} \\
& \quad \quad \quad \wedge \vee \wedge \text{newEpoch} > \text{acceptedEpoch}[i] \\
& \quad \quad \quad \quad \wedge \text{acceptedEpoch}' = [\text{acceptedEpoch} \text{ EXCEPT } ![i] = \text{newEpoch}] \\
& \quad \quad \quad \quad \wedge \text{LET } \text{epochBytes} \triangleq \text{currentEpoch}[i] \\
& \quad \quad \quad \quad \quad m \triangleq [\text{mtype} \mapsto \text{ACKEPOCH}, \\
& \quad \quad \quad \quad \quad \quad \text{mzxid} \mapsto \text{lastProcessed}[i].\text{zxid}, \\
& \quad \quad \quad \quad \quad \quad \text{mepoch} \mapsto \text{epochBytes}] \\
& \quad \quad \quad \quad \text{IN } \text{Reply}(i, j, m) \\
& \quad \quad \vee \wedge \text{newEpoch} = \text{acceptedEpoch}[i] \\
& \quad \quad \quad \wedge \text{LET } m \triangleq [\text{mtype} \mapsto \text{ACKEPOCH}, \\
& \quad \quad \quad \quad \text{mzxid} \mapsto \text{lastProcessed}[i].\text{zxid}, \\
& \quad \quad \quad \quad \text{mepoch} \mapsto -1] \\
& \quad \quad \quad \quad \text{IN } \text{Reply}(i, j, m) \\
& \quad \quad \quad \quad \wedge \text{UNCHANGED } \text{acceptedEpoch} \\
& \quad \quad \wedge \text{zabState}' = [\text{zabState} \text{ EXCEPT } ![i] = \text{SYNCHRONIZATION}] \\
& \quad \quad \wedge \text{UNCHANGED } \text{violatedInvariants}
\end{aligned}$$

```

    ∨ ∧ ¬stateOk
    ∧ PrintT("Exception: Follower receives LEADERINFO," ∘
      " whileZabState not DISCOVERY.")
    ∧ violatedInvariants' = [violatedInvariants EXCEPT !.stateInconsistent = TRUE]
    ∧ Discard(j, i)
    ∧ UNCHANGED ⟨acceptedEpoch, zabState⟩
    ∧ UNCHANGED ⟨varsL, leaderAddr, learners, forwarding, electing⟩
  ∨ 2. Abnormal case - go back to election
    ∧ ¬epochOk
    ∧ FollowerShutdown(i)
    ∧ Clean(i, leaderAddr[i])
    ∧ RemoveLearner(leaderAddr[i], i)
    ∧ UNCHANGED ⟨acceptedEpoch, violatedInvariants⟩
  ∧ UNCHANGED ⟨history, lastCommitted, connecting, ackldRecv, tempMaxEpoch,
    initialHistory, proposalMsgsLog, epochLeader, packetsSync⟩
  ∧ UpdateRecorder(⟨"FollowerProcessLEADERINFO", i, j⟩)

```

```

RECURSIVE UpdateAckSidHelper(−, −, −, −)
UpdateAckSidHelper(his, cur, end, target) ≜
  IF cur > end THEN his
  ELSE LET curTxn ≜ [xid ↦ his[1].xid,
    value ↦ his[1].value,
    ackSid ↦ IF target ∈ his[1].ackSid THEN his[1].ackSid
      ELSE his[1].ackSid ∪ {target},
    epoch ↦ his[1].epoch]
  IN ⟨curTxn⟩ ∘ UpdateAckSidHelper(Tail(his), cur + 1, end, 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) ≜
  IF Len(his) = 0 ∨ lastSeenIndex = 0 THEN his
  ELSE UpdateAckSidHelper(his, 1, Minimum({Len(his), lastSeenIndex}), target)

```

```

return − 1: this xid appears at least twice; Len(his) + 1: does not exist;
1 ¬Len(his): exists and appears just once.
RECURSIVE XidToIndexHepler(−, −, −, −)
XidToIndexHepler(his, xid, cur, appeared) ≜
  IF cur > Len(his) THEN cur
  ELSE IF TxnXidEqual(his[cur], xid)
    THEN CASE appeared = TRUE → − 1
      □ OTHER → Minimum({cur,
        XidToIndexHepler(his, xid, cur + 1, TRUE)})
    ELSE XidToIndexHepler(his, xid, cur + 1, appeared)

XidToIndex(his, xid) ≜ IF XidEqual(xid, ⟨0, 0⟩) THEN 0

```



```

ELSE IF  $Len(his) = 0$  THEN 1
    ELSE LET  $len \triangleq Len(his)$  IN
        IF  $\exists idx \in 1 \dots len : TxnZxidEqual(his[idx], zxid)$ 
            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  $queuePackets(Append(queue, m\_proposal), his, cur + 1, committed, end)$ 
     $\square \ cur \leq committed \rightarrow$ 
        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 $\langle \text{PROPOSAL}, \text{COMMIT} \rangle$. For proposals in *outstandingProposals*, send *PROPOSAL* only.

$\text{StartForwarding}(i, j, \text{lastSeenZxid}, \text{lastSeenIndex}, \text{mode}, \text{needRemoveHead}) \triangleq$
 $\wedge \text{LET } \text{lastCommittedIndex} \triangleq \text{IF } \text{zabState}[i] = \text{BROADCAST}$
 $\quad \text{THEN } \text{lastCommitted}[i].\text{index}$
 $\quad \text{ELSE } \text{Len}(\text{initialHistory}[i])$
 $\text{lastProposedIndex} \triangleq \text{Len}(\text{history}[i])$
 $\text{queue_origin} \triangleq \text{IF } \text{lastSeenIndex} \geq \text{lastProposedIndex}$
 $\quad \text{THEN } \langle \rangle$
 $\quad \text{ELSE } \text{queuePackets}(\langle \rangle, \text{history}[i],$
 $\quad \quad \text{lastSeenIndex} + 1, \text{lastCommittedIndex},$
 $\quad \quad \text{lastProposedIndex})$
 $\text{set_forChecking} \triangleq \text{IF } \text{lastSeenIndex} \geq \text{lastProposedIndex}$
 $\quad \text{THEN } \{ \}$
 $\quad \text{ELSE } \text{setPacketsForChecking}(\{ \}, i,$
 $\quad \quad \text{acceptedEpoch}[i], \text{history}[i],$
 $\quad \quad \text{lastSeenIndex} + 1, \text{lastProposedIndex})$
 $\text{m_trunc} \triangleq [\text{mtype} \mapsto \text{TRUNC}, \text{mtruncZxid} \mapsto \text{lastSeenZxid}]$
 $\text{m_diff} \triangleq [\text{mtype} \mapsto \text{DIFF}, \text{mzxid} \mapsto \text{lastSeenZxid}]$
 $\text{newLeaderZxid} \triangleq \langle \text{acceptedEpoch}[i], 0 \rangle$
 $\text{m_newleader} \triangleq [\text{mtype} \mapsto \text{NEWLEADER},$
 $\quad \text{mzxid} \mapsto \text{newLeaderZxid}]$
 $\text{queue_toSend} \triangleq \text{CASE } \text{mode} = \text{TRUNC} \rightarrow (\langle \text{m_trunc} \rangle \circ \text{queue_origin}) \circ \langle \text{m_newleader} \rangle$
 $\quad \square \quad \text{OTHER} \rightarrow (\langle \text{m_diff} \rangle \circ \text{queue_origin}) \circ \langle \text{m_newleader} \rangle$
 $\text{IN } \wedge \vee \wedge \text{needRemoveHead}$
 $\quad \wedge \text{DiscardAndSendPackets}(i, j, \text{queue_toSend})$
 $\quad \vee \wedge \neg \text{needRemoveHead}$
 $\quad \wedge \text{SendPackets}(i, j, \text{queue_toSend})$
 $\quad \wedge \text{proposalMsgsLog}' = \text{proposalMsgsLog} \cup \text{set_forChecking}$
 $\wedge \text{forwarding}' = [\text{forwarding} \text{ EXCEPT } ![i] = @ \cup \{j\}]$
 $\wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{UpdateAckSid}(@, \text{lastSeenIndex}, j)]$

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

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

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 $\text{case1. } \text{peerLastZxid} = \text{lastProcessedZxid}$
 $\text{DIFF} + \text{StartForwarding}(\text{lastProcessedZxid})$
 $\wedge \text{ZxidEqual}(\text{peerLastZxid}, \text{lastProcessedZxid})$
 $\wedge \text{StartForwarding}(i, j, \text{peerLastZxid}, \text{lastProcessed}[i].\text{index}, \text{DIFF}, \text{needRemoveHead})$
 $\vee \wedge \neg \text{ZxidEqual}(\text{peerLastZxid}, \text{lastProcessedZxid})$
 $\wedge \vee \text{case2. } \text{peerLastZxid} > \text{maxCommittedLog}$
 $\text{TRUNC} + \text{StartForwarding}(\text{maxCommittedLog})$
 $\wedge \text{ZxidCompare}(\text{peerLastZxid}, \text{maxCommittedLog})$
 $\wedge \text{LET } \text{maxCommittedIndex} \triangleq \text{IF } \text{zabState}[i] = \text{BROADCAST}$
 $\text{THEN } \text{lastCommitted}[i].\text{index}$
 $\text{ELSE } \text{Len}(\text{initialHistory}[i])$
IN $\text{StartForwarding}(i, j, \text{maxCommittedLog}, \text{maxCommittedIndex}, \text{TRUNC}, \text{needRemoveHead})$
 $\vee \text{case3. } \text{minCommittedLog} \leq \text{peerLastZxid} \leq \text{maxCommittedLog}$
 $\wedge \neg \text{ZxidCompare}(\text{peerLastZxid}, \text{maxCommittedLog})$
 $\wedge \text{LET } \text{lastSeenIndex} \triangleq \text{ZxidToIndex}(\text{history}[i], \text{peerLastZxid})$
 $\text{exist} \triangleq \wedge \text{lastSeenIndex} \geq 0$
 $\wedge \text{lastSeenIndex} \leq \text{Len}(\text{history}[i])$
 $\text{lastIndex} \triangleq \text{IF } \text{exist} \text{ THEN } \text{lastSeenIndex}$
 $\text{ELSE } \text{IndexOfZxid}(\text{history}[i], \text{peerLastZxid})$
 $\text{Maximum } \text{zxid} \text{ that } < \text{peerLastZxid}$
 $\text{lastZxid} \triangleq \text{IF } \text{exist} \text{ THEN } \text{peerLastZxid}$
 $\text{ELSE IF } \text{lastIndex} = 0 \text{ THEN } \langle 0, 0 \rangle$
 $\text{ELSE } \text{history}[i][\text{lastIndex}].\text{zxid}$
IN
 $\vee \text{case 3.1. } \text{peerLastZxid} \text{ exists in history}$
 $\text{DIFF} + \text{StartForwarding}$
 $\wedge \text{exist}$
 $\wedge \text{StartForwarding}(i, j, \text{peerLastZxid}, \text{lastSeenIndex}, \text{DIFF}, \text{needRemoveHead})$
 $\vee \text{case 3.2. } \text{peerLastZxid} \text{ does not exist in history}$
 $\text{TRUNC} + \text{StartForwarding}$
 $\wedge \neg \text{exist}$
 $\wedge \text{StartForwarding}(i, j, \text{lastZxid}, \text{lastIndex}, \text{TRUNC}, \text{needRemoveHead})$
we will not have case 4 where $\text{peerLastZxid} < \text{minCommittedLog}$, because minCommittedLog default value is 1 in our spec.

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

$$\wedge \text{ZxidCompare}(ss1.\text{lastZxid}, ss2.\text{lastZxid})$$

$$\begin{aligned} \text{electionFinished in Leader} \\ \text{ElectionFinished}(i, \text{set}) &\triangleq \wedge i \in \text{set} \\ &\wedge \text{IsQuorum}(\text{set}) \end{aligned}$$

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)*.

$$\begin{aligned} \text{UpdateElecting}(\text{oldSet}, \text{sid}, \text{peerLastZxid}, \text{inQuorum}) &\triangleq \\ \text{LET } \text{sid_electing} &\triangleq \{s.\text{sid} : s \in \text{oldSet}\} \\ \text{IN IF } \text{sid} \in \text{sid_electing} & \\ \text{THEN LET } \text{old_info} &\triangleq \text{CHOOSE } \text{info} \in \text{oldSet} : \text{info.sid} = \text{sid} \\ \text{follower_info} &\triangleq \\ &\quad [\text{sid} \mapsto \text{sid}, \\ &\quad \text{peerLastZxid} \mapsto \text{peerLastZxid}, \\ &\quad \text{inQuorum} \mapsto (\text{inQuorum} \vee \text{old_info.inQuorum})] \\ \text{IN } (\text{oldSet} \setminus \{\text{old_info}\}) \cup \{\text{follower_info}\} & \\ \text{ELSE LET } \text{follower_info} &\triangleq \\ &\quad [\text{sid} \mapsto \text{sid}, \\ &\quad \text{peerLastZxid} \mapsto \text{peerLastZxid}, \\ &\quad \text{inQuorum} \mapsto \text{inQuorum}] \\ \text{IN } \text{oldSet} \cup \{\text{follower_info}\} & \end{aligned}$$

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

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.

$$\begin{aligned} \text{LeaderProcessACKEPOCH}(i, j) &\triangleq \\ \wedge \text{IsLeader}(i) & \\ \wedge \text{PendingACKEPOCH}(i, j) & \\ \wedge \text{LET } \text{msg} &\triangleq \text{msgs}[j][i][1] \\ \text{infoOk} &\triangleq \text{IsMyLearner}(i, j) \\ \text{leaderStateSummary} &\triangleq [\text{currentEpoch} \mapsto \text{currentEpoch}[i], \\ &\quad \text{lastZxid} \mapsto \text{lastProcessed}[i].\text{zxid}] \\ \text{followerStateSummary} &\triangleq [\text{currentEpoch} \mapsto \text{msg.mepoch}, \\ &\quad \text{lastZxid} \mapsto \text{msg.mzxid}] \\ \text{logOk} &\triangleq \text{whether follower is no more up-to-date than leader} \\ &\quad \neg \text{IsMoreRecentThan}(\text{followerStateSummary}, \text{leaderStateSummary}) \\ \text{electing_quorum} &\triangleq \{e \in \text{electing}[i] : e.\text{inQuorum} = \text{TRUE}\} \\ \text{sid_electing} &\triangleq \{s.\text{sid} : s \in \text{electing_quorum}\} \end{aligned}$$

IN $\wedge \text{infoOk}$

$\wedge \vee$ $\text{electionFinished} = \text{true}$, jump out of waitForEpochAck .
Different from code, here we still need to record info
into electing , to help us perform syncFollower afterwards.
Since electing already meets quorum, it does not break
consistency between code and spec.

$\wedge \text{ElectionFinished}(i, \text{sid_electing})$
 $\wedge \text{electing}' = [\text{electing} \text{ EXCEPT } ![i] = \text{UpdateElecting}(@, j, \text{msg.mzxid}, \text{FALSE})]$
 $\wedge \text{Discard}(j, i)$
 $\wedge \text{UNCHANGED} \langle \text{varsL}, \text{zabState}, \text{forwarding}, \text{leaderAddr},$
 $\text{learners}, \text{epochLeader}, \text{violatedInvariants} \rangle$

$\vee \wedge \neg \text{ElectionFinished}(i, \text{sid_electing})$
 $\wedge \vee \wedge \text{zabState}[i] = \text{DISCOVERY}$
 $\wedge \text{UNCHANGED} \text{violatedInvariants}$
 $\vee \wedge \text{zabState}[i] \neq \text{DISCOVERY}$
 $\wedge \text{PrintT}(\text{"Exception: electionFinished false,"} \circ$
 $\text{" while zabState not DISCOVERY."})$
 $\wedge \text{violatedInvariants}' = [\text{violatedInvariants} \text{ EXCEPT}$
 $!.stateInconsistent = \text{TRUE}]$

$\wedge \vee \wedge \text{followerStateSummary.currentEpoch} = -1$
 $\wedge \text{electing}' = [\text{electing} \text{ EXCEPT } ![i] = \text{UpdateElecting}(@, j,$
 $\text{msg.mzxid}, \text{FALSE})]$
 $\wedge \text{Discard}(j, i)$
 $\wedge \text{UNCHANGED} \langle \text{varsL}, \text{zabState}, \text{forwarding}, \text{leaderAddr},$
 $\text{learners}, \text{epochLeader} \rangle$

$\vee \wedge \text{followerStateSummary.currentEpoch} > -1$
 $\wedge \vee$ normal follower
 $\wedge \text{logOk}$
 $\wedge \text{electing}' = [\text{electing} \text{ EXCEPT } ![i] =$
 $\text{UpdateElecting}(@, j, \text{msg.mzxid}, \text{TRUE})]$
 $\wedge \text{LET } \text{new_electing_quorum} \triangleq \{e \in \text{electing}'[i] : e.inQuorum = \text{TRUE}\}$
 $\text{new_sid_electing} \triangleq \{s.sid : s \in \text{new_electing_quorum}\}$

IN

\vee $\text{electionFinished} = \text{true}$, jump out of waitForEpochAck ,
update currentEpoch and zabState .
 $\wedge \text{ElectionFinished}(i, \text{new_sid_electing})$
 $\wedge \text{LeaderTurnToSynchronization}(i)$
 $\wedge \text{LET } \text{newLeaderEpoch} \triangleq \text{acceptedEpoch}[i]$
IN $\text{epochLeader}' = [\text{epochLeader} \text{ EXCEPT } ![newLeaderEpoch]$
 $= @ \cup \{i\}]$ for checking invariants

\vee $\text{there still exists } \text{electionFinished} = \text{false}$.
 $\wedge \neg \text{ElectionFinished}(i, \text{new_sid_electing})$
 $\wedge \text{UNCHANGED} \langle \text{currentEpoch}, \text{zabState}, \text{epochLeader} \rangle$
 $\wedge \text{Discard}(j, i)$
 $\wedge \text{UNCHANGED} \langle \text{state}, \text{lastProcessed}, \text{electionVars}, \text{leadingVoteSet},$

electionMsgs, leaderAddr, learners, history, forwarding)

\vee Exists follower more recent than leader
 $\wedge \neg \text{logOk}$
 $\wedge \text{LeaderShutdown}(i)$
 $\wedge \text{UNCHANGED } \langle \text{electing}, \text{epochLeader} \rangle$
 $\wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{lastCommitted}, \text{connecting}, \text{ackldRecv},$
 $\text{tempMaxEpoch}, \text{initialHistory}, \text{packetsSync}, \text{proposalMsgsLog} \rangle$
 $\wedge \text{UpdateRecorder}(\langle \text{"LeaderProcessACKEPOCH"}, i, j \rangle)$

Strip *syncFollower* from *LeaderProcessACKEPOCH*.

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

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

$\text{TruncateLog}(\text{his}, \text{index}) \triangleq \text{IF } \text{index} \leq 0 \text{ THEN } \langle \rangle$
 $\text{ELSE } \text{SubSeq}(\text{his}, 1, \text{index})$

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

$\text{FollowerProcessSyncMessage}(i, j) \triangleq$
 $\wedge \text{IsFollower}(i)$
 $\wedge \text{msgs}[j][i] \neq \langle \rangle$
 $\wedge \text{msgs}[j][i][1].\text{mtype} = \text{DIFF} \vee \text{msgs}[j][i][1].\text{mtype} = \text{TRUNC}$
 $\wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1]$

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
 $\wedge nextZxid[2] = 1$
 $\wedge curZxid[1] < nextZxid[1]$
 \vee not first PROPOSAL in this epoch
 $\wedge nextZxid[2] > 1$
 $\wedge curZxid[1] = nextZxid[1]$
 $\wedge curZxid[2] + 1 = nextZxid[2]$

$FollowerProcessPROPOSALInSync(i, j) \triangleq$
 $\wedge IsFollower(i)$
 $\wedge PendingPROPOSAL(i, j)$
 $\wedge zabState[i] = SYNCHRONIZATION$
 $\wedge LET msg \triangleq msgs[j][i][1]$
 $infoOk \triangleq IsMyLeader(i, j)$
 $isNext \triangleq IsNextZxid>LastQueued(i).zxid, msg.mzxid)$
 $newTrn \triangleq [zxid \mapsto msg.mzxid,$
 $value \mapsto msg.mdata,$
 $ackSid \mapsto \{\},$ follower do not consider $ackSid$
 $epoch \mapsto acceptedEpoch[i]]$ epoch of this round
 IN $\wedge infoOk$
 $\wedge \vee \wedge isNext$
 $\wedge packetsSync' = [packetsSync \text{ EXCEPT } ![i].notCommitted$
 $= Append(packetsSync[i].notCommitted, newTrn)]$
 $\vee \wedge \neg isNext$
 $\wedge PrintT(\text{"Warn: Follower receives PROPOSAL,"} \circ$
 $\text{" while } zxid != lastQueued + 1.")$
 $\wedge UNCHANGED packetsSync$
 $logRequest \rightarrow SyncRequestProcessor \rightarrow SendAckRequestProcessor \rightarrow \text{reply ack}$
 So here we do not need to send ack to leader.
 $\wedge Discard(j, i)$
 $\wedge UNCHANGED \langle serverVars, electionVars, leaderVars, leaderAddr,$

$verifyVars, electionMsgs\rangle$
 $\wedge UpdateRecorder(\langle \text{"FollowerProcessPROPOSALInSync"}, i, j \rangle)$

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

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

To simplify specification, we assume $snapshotNeeded = \text{false}$ and $writeToTxnLog = \text{true}$. So here we just call $packetsCommitted.add$.

$FollowerProcessCOMMITInSync(i, j) \triangleq$
 $\wedge IsFollower(i)$
 $\wedge PendingCOMMIT(i, j)$
 $\wedge zabState[i] = SYNCHRONIZATION$
 $\wedge \text{LET } msg \triangleq msgs[j][i][1]$

$$\begin{aligned}
& infoOk \triangleq IsMyLeader(i, j) \\
& committedIndex \triangleq LastCommitted(i).index + 1 \\
& exist \triangleq \wedge committedIndex \leq LastQueued(i).index \\
& \quad \wedge IsNextZxid(LastCommitted(i).zxid, msg.mzxid) \\
& match \triangleq ZxidEqual(msg.mzxid, TrnWithIndex(i, committedIndex).zxid) \\
IN & \wedge infoOk \\
& \wedge \vee \wedge exist \\
& \quad \wedge \vee \wedge match \\
& \quad \quad \wedge packetsSync' = [packetsSync \text{ EXCEPT } ![i].committed \\
& \quad \quad \quad = Append(packetsSync[i].committed, msg.mzxid)] \\
& \quad \quad \wedge UNCHANGED violatedInvariants \\
& \quad \vee \wedge \neg match \\
& \quad \quad \wedge PrintT("Warn: Follower receives COMMIT," \circ \\
& \quad \quad \quad "but zxid not the next committed zxid in COMMIT.") \\
& \quad \quad \wedge violatedInvariants' = [violatedInvariants \text{ EXCEPT } \\
& \quad \quad \quad \quad !.commitInconsistent = TRUE] \\
& \quad \quad \wedge UNCHANGED packetsSync \\
& \vee \wedge \neg exist \\
& \quad \wedge PrintT("Warn: Follower receives COMMIT," \circ \\
& \quad \quad "but no packets with its zxid exists.") \\
& \quad \wedge violatedInvariants' = [violatedInvariants \text{ EXCEPT } \\
& \quad \quad \quad \quad !.commitInconsistent = TRUE] \\
& \quad \wedge UNCHANGED packetsSync \\
& \wedge Discard(j, i) \\
& \wedge UNCHANGED \langle serverVars, electionVars, leaderVars, \\
& \quad \quad leaderAddr, epochLeader, proposalMsgsLog, electionMsgs \rangle \\
& \wedge UpdateRecorder(("FollowerProcessCOMMITInSync", i, j)) \\
\\
RECURSIVE & ACKInBatches(-, -) \\
& ACKInBatches(queue, packets) \triangleq \\
& \quad IF packets = \langle \rangle \text{ THEN } queue \\
& \quad ELSE LET head \triangleq packets[1] \\
& \quad \quad newPackets \triangleq Tail(packets) \\
& \quad \quad m_ack \triangleq [mtype \mapsto ACK, \\
& \quad \quad \quad mzxid \mapsto 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) \triangleq \\
& \wedge IsFollower(i) \\
& \wedge PendingNEWLEADER(i, j) \\
& \wedge LET msg \triangleq msgs[j][i][1] \\
& \quad infoOk \triangleq IsMyLeader(i, j) \\
& \quad packetsInSync \triangleq packetsSync[i].notCommitted
\end{aligned}$$

$m_ackld \triangleq [mtype \mapsto ACKLD,$
 $\quad \quad \quad mzxid \mapsto msg.mzxid]$
 $ms_ack \triangleq ACKInBatches(\langle \rangle, packetsInSync)$
 $queue_toSend \triangleq \langle m_ackld \rangle \circ ms_ack$ send ACK – NEWLEADER first.
IN $\wedge infoOk$
 $\wedge currentEpoch' = [currentEpoch \text{ EXCEPT } ![i] = acceptedEpoch[i]]$
 $\wedge history' = [history \text{ EXCEPT } ![i] = @ \circ packetsInSync]$
 $\wedge packetsSync' = [packetsSync \text{ EXCEPT } ![i].notCommitted = \langle \rangle]$
 $\wedge DiscardAndSendPackets(i, j, queue_toSend)$
 $\wedge \text{UNCHANGED } \langle state, lastProcessed, zabState, acceptedEpoch, lastCommitted,$
 $\quad \quad \quad electionVars, leaderVars, initialHistory, leaderAddr, verifyVars,$
 $\quad \quad \quad electionMsgs \rangle$
 $\wedge UpdateRecorder(\langle \text{"FollowerProcessNEWLEADER"}, i, j \rangle)$

quorumFormed in Leader
 $QuorumFormed(i) \triangleq i \in ackldRecv[i] \wedge IsQuorum(ackldRecv[i])$
 $QuorumFormedTurnToTrue(i) \triangleq i \in ackldRecv'[i] \wedge IsQuorum(ackldRecv'[i])$
 $UpdateElectionVote(i, epoch) \triangleq UpdateProposal(i, currentVote[i].proposedLeader,$
 $\quad \quad \quad currentVote[i].proposedZxid, epoch)$

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.
 $LeaderProcessACKLD(i, j) \triangleq$
 $\wedge IsLeader(i)$
 $\wedge PendingACKLD(i, j)$
 $\wedge \text{LET } msg \triangleq msgs[j][i][1]$
 $\quad infoOk \triangleq IsMyLearner(i, j)$
 $\quad match \triangleq ZxidEqual(msg.mzxid, \langle acceptedEpoch[i], 0 \rangle)$
 $\quad currentZxid \triangleq \langle acceptedEpoch[i], 0 \rangle$
 $\quad m_uptodate \triangleq [mtype \mapsto UPTODATE,$
 $\quad \quad \quad mzxid \mapsto currentZxid]$ not important
IN $\wedge infoOk$
 $\wedge \vee$ just reply UPTODATE.
 $\quad \wedge QuorumFormed(i)$

```

    ∧ Reply(i, j, m_uptodate)
    ∧ UNCHANGED ⟨ackldRecv, zabState, lastCommitted, lastProcessed,
                  currentVote, violatedInvariants⟩
  ∨ ∧ ¬QuorumFormed(i)
    ∧ ∨ ∧ match
      ∧ ackldRecv' = [ackldRecv EXCEPT ![i] = @ ∪ {j}]
      ∧ ∨ jump out of waitForNewLeaderAck, and do startZkServer,
            setZabState, and reply UPTODATE.
      ∧ QuorumFormedTurnToTrue(i)
      ∧ LeaderTurnToBroadcast(i)
      ∧ DiscardAndBroadcastUPTODATE(i, j, m_uptodate)
    ∨ still wait in waitForNewLeaderAck.
      ∧ ¬QuorumFormedTurnToTrue(i)
      ∧ Discard(j, i)
      ∧ UNCHANGED ⟨zabState, lastCommitted, lastProcessed, currentVote⟩
    ∧ 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.

$$\begin{array}{l}
\text{PendingTxns}(i) \triangleq \text{IF } \neg \text{IsFollower}(i) \vee \text{zabState}[i] \neq \text{SYNCHRONIZATION} \\
\quad \text{THEN } \text{SubSeq}(\text{history}[i], \text{lastCommitted}[i].\text{index} + 1, \text{Len}(\text{history}[i])) \\
\quad \text{ELSE LET } \text{packetsCommitted} \triangleq \text{packetsSync}[i].\text{committed} \\
\quad \quad \text{completeHis} \triangleq \text{history}[i] \circ \text{packetsSync}[i].\text{notCommitted} \\
\quad \text{IN IF } \text{Len}(\text{packetsCommitted}) = 0 \\
\quad \quad \text{THEN } \text{SubSeq}(\text{completeHis}, \text{Len}(\text{initialHistory}[i]) + 1, \text{Len}(\text{completeHis})) \\
\quad \quad \text{ELSE } \text{SubSeq}(\text{completeHis}, \text{LastCommitted}(i).\text{index} + 1, \text{Len}(\text{completeHis}))
\end{array}$$
$$\begin{array}{l}
CommittedTrns(i) \triangleq \text{IF } \neg IsFollower(i) \vee zabState[i] \neq SYNCHRONIZATION \\
\quad \text{THEN } SubSeq(history[i], 1, lastCommitted[i].index) \\
\quad \text{ELSE LET } packetsCommitted \triangleq packetsSync[i].committed \\
\quad \quad \quad completeHis \triangleq history[i] \circ packetsSync[i].notCommitted \\
\quad \text{IN IF } Len(packetsCommitted) = 0 \text{ THEN } initialHistory[i] \\
\quad \quad \text{ELSE } SubSeq(completeHis, 1, LastCommitted(i).index)
\end{array}$$

Each *xxid* of *packetsCommitted* equals to *xxid* of

corresponding txn in $txns$.

```

RECURSIVE TrnsAndCommittedMatch(-, -)
TrnsAndCommittedMatch(trns, packetsCommitted)  $\triangleq$ 
  LET len1  $\triangleq$  Len(trns)
      len2  $\triangleq$  Len(packetsCommitted)
  IN  IF len2 = 0 THEN TRUE
      ELSE IF len1 < len2 THEN FALSE
          ELSE  $\wedge$  ZxidEqual(trns[len1].zxid, packetsCommitted[len2])
               $\wedge$  TrnsAndCommittedMatch(SubSeq(trns, 1, len1 - 1),
                                      SubSeq(packetsCommitted, 1, len2 - 1))

```

$$\begin{aligned} & \text{FollowerLogRequestInBatches}(i, \text{leader}, \text{ms_ack}, \text{packetsNotCommitted}) \triangleq \\ & \quad \wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = @ \circ \text{packetsNotCommitted}] \\ & \quad \wedge \text{DiscardAndSendPackets}(i, \text{leader}, \text{ms_ack}) \end{aligned}$$

Since `commit` will call `commitProcessor.commit`, which will finally

update *lastProcessed*, we update it here atomically.

$$\begin{array}{l}
\text{FollowerCommitInBatches}(i) \triangleq \\
\quad \text{LET } committedTxns \triangleq CommittedTxns(i) \\
\quad \quad packetsCommitted \triangleq packetsSync[i].committed \\
\quad \quad match \triangleq TxnsAndCommittedMatch(committedTxns, packetsCommitted) \\
\text{IN} \\
\quad \vee \wedge match \\
\quad \wedge lastCommitted' = [lastCommitted \text{ EXCEPT } ![i] = LastCommitted(i)] \\
\quad \wedge lastProcessed' = [lastProcessed \text{ EXCEPT } ![i] = lastCommitted'[i]] \\
\quad \wedge \text{UNCHANGED } violatedInvariants
\end{array}$$

$\vee \wedge \neg match$
 $\wedge PrintT(\text{"Warn: Committing zxid without see txn. /"} \circ$
 $\quad \text{"Committing zxid != pending txn zxid."})$
 $\wedge violatedInvariants' = [violatedInvariants \text{ EXCEPT}$
 $\quad \quad !.commitInconsistent = TRUE]$
 $\wedge UNCHANGED \langle lastCommitted, lastProcessed \rangle$

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.

$FollowerProcessUPTODATE(i, j) \triangleq$
 $\wedge IsFollower(i)$
 $\wedge PendingUPTODATE(i, j)$
 $\wedge LET \ msg \triangleq \ msgs[j][i][1]$
 $\quad infoOk \triangleq IsMyLeader(i, j)$
 $\quad packetsNotCommitted \triangleq packetsSync[i].notCommitted$
 $\quad ms_ack \triangleq ACKInBatches(\langle \rangle, packetsNotCommitted)$
 IN $\wedge infoOk$
 $\quad \text{Here we ignore ack of UPTODATE.}$
 $\wedge UpdateElectionVote(i, acceptedEpoch[i])$
 $\wedge FollowerLogRequestInBatches(i, j, ms_ack, packetsNotCommitted)$
 $\wedge FollowerCommitInBatches(i)$
 $\wedge packetsSync' = [packetsSync \text{ EXCEPT } ![i].notCommitted = \langle \rangle,$
 $\quad \quad \quad ![i].committed = \langle \rangle]$
 $\wedge zabState' = [zabState \text{ EXCEPT } ![i] = BROADCAST]$
 $\wedge UNCHANGED \langle state, currentEpoch, acceptedEpoch, logicalClock,$
 $\quad receiveVotes, outOfElection, recvQueue, waitNotmsg, leaderVars,$
 $\quad initialHistory, leaderAddr, epochLeader, proposalMsgsLog, electionMsgs \rangle$
 $\wedge UpdateRecorder(\langle \text{"FollowerProcessUPTODATE"}, i, j \rangle)$

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

Leader receives client propose and broadcasts *PROPOSAL*. See *processRequest* in
ProposalRequestProcessor and propose in Leader for details. Since
 $proposalProcessor.processRequest \rightarrow syncProcessor.processRequest \rightarrow$
 $ackProcessor.processRequest \rightarrow leader.processAck$, we initially set $txn.ackSid = \{i\}$, assuming
 we have done *leader.processAck*. 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.

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

$$\begin{aligned}
newTrn &\triangleq [xid \mapsto IncZxid(i, LastProposed(i).xid), \\
&\quad value \mapsto request_value, \\
&\quad ackSid \mapsto \{i\}, \quad \text{assume we have done } leader.processAck \\
&\quad epoch \mapsto acceptedEpoch[i]] \\
m_proposal &\triangleq [mtype \mapsto PROPOSAL, \\
&\quad mzxid \mapsto newTrn.xid, \\
&\quad mdata \mapsto request_value] \\
m_proposal_for_checking &\triangleq [source \mapsto i, \\
&\quad epoch \mapsto acceptedEpoch[i], \\
&\quad xid \mapsto newTrn.xid, \\
&\quad data \mapsto request_value] \\
IN \quad &\wedge history' = [history \text{ EXCEPT } ![i] = Append(@, newTrn)] \\
&\wedge Broadcast(i, m_proposal) \\
&\wedge proposalMsgsLog' = proposalMsgsLog \cup \{m_proposal_for_checking\} \\
&\wedge UNCHANGED \langle state, currentEpoch, lastProcessed, zabState, acceptedEpoch, \\
&\quad lastCommitted, electionVars, leaderVars, followerVars, initialHistory, \\
&\quad epochLeader, violatedInvariants, electionMsgs \rangle \\
&\wedge UpdateRecorder(("LeaderProcessRequest", i))
\end{aligned}$$

Follower processes *PROPOSAL* in *BROADCAST*. See *processPacket* in Follower for details.

$$\begin{aligned}
FollowerProcessPROPOSAL(i, j) &\triangleq \\
&\wedge IsFollower(i) \\
&\wedge PendingPROPOSAL(i, j) \\
&\wedge zabState[i] = BROADCAST \\
&\wedge LET \ msg \triangleq \ msgs[j][i][1] \\
&\quad infoOk \triangleq IsMyLeader(i, j) \\
&\quad isNext \triangleq IsNextZxid(LastQueued(i).xid, msg.mzxid) \\
&\quad newTrn \triangleq [xid \mapsto msg.mzxid, \\
&\quad \quad value \mapsto msg.mdata, \\
&\quad \quad ackSid \mapsto \{\}, \\
&\quad \quad epoch \mapsto acceptedEpoch[i]] \\
&\quad m_ack \triangleq [mtype \mapsto ACK, \\
&\quad \quad mzxid \mapsto msg.mzxid] \\
IN \quad &\wedge infoOk \\
&\wedge \vee \wedge isNext \\
&\quad \wedge UNCHANGED \ violatedInvariants \\
&\quad \vee \wedge \neg isNext \\
&\quad \wedge PrintT("Exception: Follower receives PROPOSAL, while" \circ \\
&\quad \quad "the transaction is not the next.") \\
&\quad \wedge violatedInvariants' = [violatedInvariants \text{ EXCEPT } \\
&\quad \quad \quad !.proposalInconsistent = TRUE] \\
&\quad \wedge history' = [history \text{ EXCEPT } ![i] = Append(@, newTrn)] \\
&\quad \wedge Reply(i, j, m_ack) \\
&\wedge UNCHANGED \langle state, currentEpoch, lastProcessed, zabState, acceptedEpoch, \\
&\quad lastCommitted, electionVars, leaderVars, followerVars, initialHistory,
\end{aligned}$$

$epochLeader, proposalMsgsLog, electionMsgs\rangle$
 $\wedge UpdateRecorder(\langle \text{"FollowerProcessPROPOSAL"}, i, j \rangle)$

See *outstandingProposals* in Leader

$OutstandingProposals(i) \triangleq$ IF $zabState[i] \neq BROADCAST$ THEN $\langle \rangle$
ELSE $SubSeq(history[i], lastCommitted[i].index + 1,$
 $Len(history[i]))$

$LastAckIndexFromFollower(i, j) \triangleq$
LET $set_index \triangleq \{idx \in 1 \dots Len(history[i]) : j \in history[i][idx].ackSid\}$
IN $Maximum(set_index)$

See *commit* in Leader for details.

$LeaderCommit(s, follower, index, zxid) \triangleq$
 $\wedge lastCommitted' = [lastCommitted \text{ EXCEPT } ![s] = [index \mapsto index,$
 $zxid \mapsto zxid]]$
 \wedge LET $m_commit \triangleq [mtype \mapsto COMMIT,$
 $mzxid \mapsto zxid]$
IN $DiscardAndBroadcast(s, follower, m_commit)$

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.

$LeaderTryToCommit(s, index, zxid, newTxn, follower) \triangleq$
LET $allTxnsBeforeCommitted \triangleq lastCommitted[s].index \geq index - 1$
Only when all proposals before $zxid$ 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.
 $ordered \triangleq lastCommitted[s].index + 1 = index$
Commit proposals in order.
IN $\vee \wedge$ Current conditions do not satisfy committing the proposal.
 $\vee \neg allTxnsBeforeCommitted$
 $\vee \neg hasAllQuorums$
 $\wedge Discard(follower, s)$
 $\wedge UNCHANGED \langle violatedInvariants, lastCommitted, lastProcessed \rangle$
 $\vee \wedge allTxnsBeforeCommitted$
 $\wedge hasAllQuorums$
 $\wedge \vee \wedge \neg ordered$
 $\wedge PrintT(\text{"Warn: Committing zxid " } \circ ToString(zxid) \circ \text{" not first."})$
 $\wedge violatedInvariants' = [violatedInvariants \text{ EXCEPT }$
 $!.commitInconsistent = TRUE]$
 $\vee \wedge ordered$

$$\begin{aligned}
& \wedge \text{UNCHANGED } \text{violatedInvariants} \\
& \wedge \text{LeaderCommit}(s, \text{follower}, \text{index}, \text{zxid}) \\
& \wedge \text{lastProcessed}' = [\text{lastProcessed} \text{ EXCEPT } ![s] = [\text{index} \mapsto \text{index}, \\
& \hspace{15em} \text{zxid} \mapsto \text{zxid}]]
\end{aligned}$$

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.

$$\begin{aligned}
& \text{LeaderProcessACK}(i, j) \triangleq \\
& \quad \wedge \text{IsLeader}(i) \\
& \quad \wedge \text{PendingACK}(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{outstanding} \triangleq \text{LastCommitted}(i).\text{index} < \text{LastProposed}(i).\text{index} \\
& \quad \quad \quad \text{outstandingProposals not null} \\
& \quad \quad \text{hasCommitted} \triangleq \neg \text{ZxidCompare}(\text{msg.mzxid}, \text{LastCommitted}(i).\text{zxid}) \\
& \quad \quad \quad \text{namely, lastCommitted} \geq \text{zxid} \\
& \quad \quad \text{index} \triangleq \text{ZxidToIndex}(\text{history}[i], \text{msg.mzxid}) \\
& \quad \quad \text{exist} \triangleq \text{index} \geq 1 \wedge \text{index} \leq \text{LastProposed}(i).\text{index} \\
& \quad \quad \quad \text{the proposal exists in history} \\
& \quad \quad \text{ackIndex} \triangleq \text{LastAckIndexFromFollower}(i, j) \\
& \quad \quad \text{monotonicallyInc} \triangleq \vee \text{ackIndex} = -1 \\
& \quad \quad \quad \vee \text{ackIndex} + 1 = \text{index} \\
& \quad \quad \quad \text{TCP makes everytime } \text{ackIndex} \text{ should just increase by 1} \\
& \text{IN } \quad \wedge \text{infoOk} \\
& \quad \wedge \vee \wedge \text{exist} \\
& \quad \quad \wedge \text{monotonicallyInc} \\
& \quad \quad \wedge \text{LET } \text{txn} \triangleq \text{history}[i][\text{index}] \\
& \quad \quad \quad \text{txnAfterAddAck} \triangleq [\text{zxid} \mapsto \text{txn.zxid}, \\
& \quad \quad \quad \quad \text{value} \mapsto \text{txn.value}, \\
& \quad \quad \quad \quad \text{ackSid} \mapsto \text{txn.ackSid} \cup \{j\}, \\
& \quad \quad \quad \quad \text{epoch} \mapsto \text{txn.epoch}] \\
& \quad \text{IN } \quad \text{p.addAck}(\text{sid}) \\
& \quad \wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i][\text{index}] = \text{txnAfterAddAck}] \\
& \quad \wedge \quad \vee \wedge \quad \text{Note: outstanding is 0.} \\
& \quad \quad \quad \text{/ proposal has already been committed.} \\
& \quad \quad \vee \neg \text{outstanding} \\
& \quad \quad \vee \text{hasCommitted} \\
& \quad \quad \wedge \text{Discard}(j, i) \\
& \quad \quad \wedge \text{UNCHANGED } \langle \text{violatedInvariants}, \text{lastCommitted}, \text{lastProcessed} \rangle \\
& \quad \quad \vee \wedge \text{outstanding} \\
& \quad \quad \wedge \neg \text{hasCommitted} \\
& \quad \quad \wedge \text{LeaderTryToCommit}(i, \text{index}, \text{msg.mzxid}, \text{txnAfterAddAck}, j) \\
& \vee \wedge \vee \neg \text{exist} \\
& \quad \vee \neg \text{monotonicallyInc}
\end{aligned}$$

$\wedge \text{PrintT}(\text{"Exception: No such zxid. " } \circ$
 $\quad \text{" / ackIndex doesn't inc monotonically."})$
 $\wedge \text{violatedInvariants}' = [\text{violatedInvariants}$
 $\quad \text{EXCEPT !.ackInconsistent} = \text{TRUE}]$
 $\wedge \text{Discard}(j, i)$
 $\wedge \text{UNCHANGED} \langle \text{history}, \text{lastCommitted}, \text{lastProcessed} \rangle$
 $\wedge \text{UNCHANGED} \langle \text{state}, \text{currentEpoch}, \text{zabState}, \text{acceptedEpoch}, \text{electionVars},$
 $\quad \text{leaderVars}, \text{initialHistory}, \text{followerVars}, \text{proposalMsgsLog}, \text{epochLeader},$
 $\quad \text{electionMsgs} \rangle$
 $\wedge \text{UpdateRecorder}(\langle \text{"LeaderProcessACK"}, i, j \rangle)$

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

$\text{FollowerProcessCOMMIT}(i, j) \triangleq$
 $\wedge \text{IsFollower}(i)$
 $\wedge \text{PendingCOMMIT}(i, j)$
 $\wedge \text{zabState}[i] = \text{BROADCAST}$
 $\wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1]$
 $\quad \text{infoOk} \triangleq \text{IsMyLeader}(i, j)$
 $\quad \text{pendingTxns} \triangleq \text{PendingTxns}(i)$
 $\quad \text{noPending} \triangleq \text{Len}(\text{pendingTxns}) = 0$
 IN
 $\wedge \text{infoOk}$
 $\wedge \vee \wedge \text{noPending}$
 $\quad \wedge \text{PrintT}(\text{"Warn: Committing zxid without seeing txn."})$
 $\quad \wedge \text{UNCHANGED} \langle \text{lastCommitted}, \text{lastProcessed}, \text{violatedInvariants} \rangle$
 $\vee \wedge \neg \text{noPending}$
 $\quad \wedge \text{LET } \text{firstElementZxid} \triangleq \text{pendingTxns}[1].\text{zxid}$
 $\quad \quad \text{match} \triangleq \text{ZxidEqual}(\text{firstElementZxid}, \text{msg.mzxid})$
 IN
 $\vee \wedge \neg \text{match}$
 $\quad \wedge \text{PrintT}(\text{"Exception: Committing zxid not equals" } \circ$
 $\quad \quad \text{" next pending txn zxid."})$
 $\quad \wedge \text{violatedInvariants}' = [\text{violatedInvariants} \text{ EXCEPT}$
 $\quad \quad \text{!.commitInconsistent} = \text{TRUE}]$
 $\quad \wedge \text{UNCHANGED} \langle \text{lastCommitted}, \text{lastProcessed} \rangle$
 $\vee \wedge \text{match}$
 $\quad \wedge \text{lastCommitted}' = [\text{lastCommitted} \text{ EXCEPT}$
 $\quad \quad \text{!}[i] = [\text{index} \mapsto \text{lastCommitted}[i].\text{index} + 1,$
 $\quad \quad \quad \text{zxid} \mapsto \text{firstElementZxid}]$
 $\quad \wedge \text{lastProcessed}' = [\text{lastProcessed} \text{ EXCEPT}$
 $\quad \quad \text{!}[i] = [\text{index} \mapsto \text{lastCommitted}[i].\text{index} + 1,$
 $\quad \quad \quad \text{zxid} \mapsto \text{firstElementZxid}]$
 $\quad \wedge \text{UNCHANGED } \text{violatedInvariants}$
 $\wedge \text{Discard}(j, i)$
 $\wedge \text{UNCHANGED} \langle \text{state}, \text{currentEpoch}, \text{zabState}, \text{acceptedEpoch}, \text{history},$

$electionVars, leaderVars, initialHistory, followerVars,$
 $proposalMsgsLog, epochLeader, electionMsgs\}$
 $\wedge UpdateRecorder(\langle \text{"FollowerProcessCOMMIT"}, i, j \rangle)$

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

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

Defines how the variables may transition.

$Next \triangleq$

FLE module

$\forall \exists i, j \in \text{Server} : \text{FLEReceiveNotmsg}(i, j)$
 $\forall \exists i \in \text{Server} : \text{FLENotmsgTimeout}(i)$
 $\forall \exists i \in \text{Server} : \text{FLEHandleNotmsg}(i)$
 $\forall \exists i \in \text{Server} : \text{FLEWaitNewNotmsg}(i)$
 $\forall \exists i \in \text{Server} : \text{FLEWaitNewNotmsgEnd}(i)$

Some conditions like failure, network delay

$\forall \exists i \in \text{Server} : \text{FollowerTimeout}(i)$
 $\forall \exists i \in \text{Server} : \text{LeaderTimeout}(i)$
 $\forall \exists i, j \in \text{Server} : \text{Timeout}(i, j)$

Zab module - Discovery and Synchronization part

$\forall \exists i, j \in \text{Server} : \text{ConnectAndFollowerSendFOLLOWERINFO}(i, j)$
 $\forall \exists i, j \in \text{Server} : \text{LeaderProcessFOLLOWERINFO}(i, j)$
 $\forall \exists i, j \in \text{Server} : \text{FollowerProcessLEADERINFO}(i, j)$
 $\forall \exists i, j \in \text{Server} : \text{LeaderProcessACKEPOCH}(i, j)$
 $\forall \exists i \in \text{Server} : \text{LeaderSyncFollower}(i)$
 $\forall \exists i, j \in \text{Server} : \text{FollowerProcessSyncMessage}(i, j)$
 $\forall \exists i, j \in \text{Server} : \text{FollowerProcessPROPOSALInSync}(i, j)$
 $\forall \exists i, j \in \text{Server} : \text{FollowerProcessCOMMITInSync}(i, j)$
 $\forall \exists i, j \in \text{Server} : \text{FollowerProcessNEWLEADER}(i, j)$
 $\forall \exists i, j \in \text{Server} : \text{LeaderProcessACKLD}(i, j)$
 $\forall \exists i, j \in \text{Server} : \text{FollowerProcessUPTODATE}(i, j)$

Zab module - Broadcast part

$\forall \exists i \in \text{Server} : \text{LeaderProcessRequest}(i)$
 $\forall \exists i, j \in \text{Server} : \text{FollowerProcessPROPOSAL}(i, j)$
 $\forall \exists i, j \in \text{Server} : \text{LeaderProcessACK}(i, j)$ *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.

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

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

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

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

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

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

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

Local primary order: If a primary broadcasts a before it broadcasts b , then a follower that

delivers b must also deliver a before b.
 $LocalPrimaryOrder \triangleq LET\ p_set(i, e) \triangleq \{p \in proposalMsgsLog : \wedge p.source = i$
 $\wedge p.epoch = e\}$

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

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.

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

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.

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

\ * Modification History
 \ * Last modified Mon Nov 22 22:23:20 CST 2021 by Dell
 \ * Created Sat Oct 23 16:05:04 CST 2021 by Dell