$\overline{\phantom{xxxxxx}}$ MODULE *ZabWithFLEAndSYNC* $\overline{\phantom{xxxxxx}}$

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 |
| | | {*ELECTION*, *DISCOVERY*, *SYNCHRONIZATION*, *BROADCAST*}. |
| | *acceptedEpoch*, | 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, |

1

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* = ⟨ − 1, − 1⟩ 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 FOLLOWEINFO*} + 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    $clientReuqest$ \* 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:
     $[pc, nTransaction, maxEpoch, nTimeout, nClientRequest]$

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

$electionVars \triangleq electionVarsL$

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

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

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

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

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

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

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

$HistoryItem \triangleq$
$\quad [zxid : Zxid,$
$\quad\ value : Value,$
$\quad\ ackSid : \textsc{subset}\ Server,$
$\quad\ epoch : Nat]$

$Proposal \triangleq$
$\quad [source : Server,$
$\quad\ epoch : Nat,$
$\quad\ zxid : Zxid,$
$\quad\ data : Value]$

$LastItem \triangleq$
$\quad [index : Nat, zxid : Zxid]$

$SyncPackets \triangleq$
$\quad [notCommitted : Seq(HistoryItem),$
$\quad\ committed : Seq(Zxid)]$

$Message \triangleq$
$\quad [mtype : \{FOLLOWERINFO\}, mzxid : Zxid] \cup$
$\quad [mtype : \{LEADERINFO\}, mzxid : Zxid] \cup$
$\quad [mtype : \{ACKEPOCH\}, mzxid : Zxid, mepoch : Nat \cup \{-1\}] \cup$
$\quad [mtype : \{DIFF\}, mzxid : Zxid] \cup$
$\quad [mtype : \{TRUNC\}, mtruncZxid : Zxid] \cup$
$\quad [mtype : \{PROPOSAL\}, mzxid : Zxid, mdata : Value] \cup$
$\quad [mtype : \{COMMIT\}, mzxid : Zxid] \cup$
$\quad [mtype : \{NEWLEADER\}, mzxid : Zxid] \cup$
$\quad [mtype : \{ACKLD\}, mzxid : Zxid] \cup$
$\quad [mtype : \{ACK\}, mzxid : Zxid] \cup$
$\quad [mtype : \{UPTODATE\}, mzxid : Zxid]$

$ElectionState \triangleq \{LOOKING, FOLLOWING, LEADING\}$

$ZabState \triangleq \{ELECTION, DISCOVERY, SYNCHRONIZATION, BROADCAST\}$

$ViolationSet \triangleq \{\text{``stateInconsistent''}, \text{``proposalInconsistent''},$
$\qquad\qquad\qquad \text{``commitInconsistent''}, \text{``ackInconsistent''},$
$\qquad\qquad\qquad \text{``messageIllegal''}\}$

$Electing \triangleq [sid : Server,$
$\qquad\qquad\ peerLastZxid : Zxid,$
$\qquad\qquad\ inQuorum : \textsc{boolean}\ ]$

$Vote \triangleq$

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

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

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

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

Fast Leader Election
$$\land electionMsgs \in [Server \rightarrow [Server \rightarrow Seq(ElectionMsg)]]$$
$$\land recvQueue \in [Server \rightarrow Seq(ElectionMsg)]$$
$$\land leadingVoteSet \in [Server \rightarrow \text{SUBSET } Server]$$
$$\land receiveVotes \in [Server \rightarrow [Server \rightarrow ElectionVote]]$$
$$\land currentVote \in [Server \rightarrow Vote]$$
$$\land outOfElection \in [Server \rightarrow [Server \rightarrow ElectionVote]]$$
$$\land lastProcessed \in [Server \rightarrow LastItem]$$
$$\land history \in [Server \rightarrow Seq(HistoryItem)]$$
$$\land state \in [Server \rightarrow ElectionState]$$
$$\land waitNotmsg \in [Server \rightarrow \text{BOOLEAN }]$$
$$\land currentEpoch \in [Server \rightarrow Nat]$$
$$\land logicalClock \quad \in [Server \rightarrow Nat]$$

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

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

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

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

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

$ToZxid(z) \triangleq [epoch \mapsto z[1],\ counter \mapsto z[2]]$

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

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

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

$GetParameter(p) \triangleq$ IF $p \in$ DOMAIN $Parameters$ THEN $Parameters[p]$ ELSE $0$
$GetRecorder(p) \quad\triangleq$ IF $p \in$ DOMAIN $recorder$ $\qquad$ THEN $recorder[p]$ $\qquad$ ELSE $0$

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

$RecorderIncTimeout \triangleq RecorderIncHelper(\text{"nTimeout"})$
$RecorderGetTimeout \triangleq RecorderGetHelper(\text{"nTimeout"})$
$RecorderSetTransactionNum(pc) \triangleq (\text{"nTransaction"} :>$
$\qquad\qquad\qquad\qquad\qquad\qquad$ IF $pc[1] = \text{"LeaderProcessRequest"}$ THEN
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ LET $s \triangleq$ CHOOSE $i \in Server :$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \forall\, j \in Server : Len(history'[i]) \geq Len(history'[j])$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ IN $\ Len(history'[s])$
$\qquad\qquad\qquad\qquad\qquad\qquad$ ELSE $recorder[\text{"nTransaction"}])$
$RecorderSetMaxEpoch(pc) \qquad\triangleq (\text{"maxEpoch"} :>$

$$
\begin{aligned}
&\qquad\qquad\qquad \text{IF } pc[1] = \text{``LeaderProcessFOLLOWERINFO''} \text{ THEN} \\
&\qquad\qquad\qquad\quad \text{LET } s \triangleq \text{CHOOSE } i \in Server : \\
&\qquad\qquad\qquad\qquad\quad \forall\, j \in Server : acceptedEpoch'[i] \geq acceptedEpoch'[j] \\
&\qquad\qquad\qquad\quad \text{IN } \quad acceptedEpoch'[s] \\
&\qquad\qquad\qquad \text{ELSE } \ recorder[\text{``maxEpoch''}])
\end{aligned}
$$

$$
\begin{aligned}
RecorderSetRequests(pc) \quad &\triangleq\ (\text{``nClientRequest''} :> \\
&\quad \text{IF } pc[1] = \text{``LeaderProcessRequest''} \text{ THEN} \\
&\qquad recorder[\text{``nClientRequest''}] + 1 \\
&\quad \text{ELSE } \ recorder[\text{``nClientRequest''}])
\end{aligned}
$$

$$
\begin{aligned}
RecorderSetPc(pc) \quad &\triangleq\ (\text{``pc''} :> pc) \\
RecorderSetFailure(pc) \ &\triangleq\ \text{CASE } pc[1] = \text{``Timeout''} && \to RecorderIncTimeout \\
&\qquad\quad \square \quad pc[1] = \text{``LeaderTimeout''} &\to RecorderIncTimeout \\
&\qquad\quad \square \quad pc[1] = \text{``FollowerTimeout''} &\to RecorderIncTimeout \\
&\qquad\quad \square \quad \text{OTHER} && \to RecorderGetTimeout
\end{aligned}
$$

$$
\begin{aligned}
UpdateRecorder(pc) \ \triangleq\ recorder' = {}&RecorderSetFailure(pc) \qquad @@\ RecorderSetTransactionNum(pc) \\
&@@\ RecorderSetMaxEpoch(pc) \ \ @@\ RecorderSetPc(pc) \\
&@@\ RecorderSetRequests(pc) \ \ @@\ recorder
\end{aligned}
$$

$$
UnchangeRecorder \ \triangleq\ \text{UNCHANGED } recorder
$$

$$
\begin{aligned}
CheckParameterHelper(n, p, Comp(\_,\_)) \ \triangleq\ &\text{IF } p \in \text{DOMAIN } Parameters \\
&\quad \text{THEN } Comp(n,\ Parameters[p]) \\
&\quad \text{ELSE } \ \text{TRUE}
\end{aligned}
$$

$$
CheckParameterLimit(n, p) \ \triangleq\ CheckParameterHelper(n, p, \text{LAMBDA } i, j : i < j)
$$

$$
\begin{aligned}
CheckTimeout \quad &\triangleq\ CheckParameterLimit(recorder.nTimeout, \quad \text{``MaxTimeoutFailures''}) \\
CheckTransactionNum \ &\triangleq\ CheckParameterLimit(recorder.nTransaction, \text{``MaxTransactionNum''}) \\
CheckEpoch \quad &\triangleq\ CheckParameterLimit(recorder.maxEpoch, \quad \text{``MaxEpoch''})
\end{aligned}
$$

$$
CheckStateConstraints \ \triangleq\ CheckTimeout \wedge CheckTransactionNum \wedge CheckEpoch
$$

---

Actions about network

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

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

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

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

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

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

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

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

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

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

$CleanInputBuffer(i) \triangleq msgs' = [s \in Server \mapsto [v \in Server \mapsto \text{IF } v = i \text{ THEN } \langle\rangle$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \text{ELSE } msgs[s][v]]]$

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

---

Define initial values for all variables

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

$InitLeaderVars \triangleq \land InitLeaderVarsL$
$\qquad\qquad\qquad\quad \land learners \qquad\qquad = [s \in Server \mapsto \{\}]$
$\qquad\qquad\qquad\quad \land connecting \qquad\quad\; = [s \in Server \mapsto \{\}]$
$\qquad\qquad\qquad\quad \land electing \qquad\qquad\; = [s \in Server \mapsto \{\}]$
$\qquad\qquad\qquad\quad \land ackldRecv \qquad\quad\;\; = [s \in Server \mapsto \{\}]$
$\qquad\qquad\qquad\quad \land forwarding \qquad\quad\; = [s \in Server \mapsto \{\}]$
$\qquad\qquad\qquad\quad \land tempMaxEpoch \quad\;\; = [s \in Server \mapsto 0]$

$InitElectionVars \triangleq InitElectionVarsL$

$InitFollowerVars \triangleq \land leaderAddr = [s \in Server \mapsto NullPoint]$
$\qquad\qquad\qquad\qquad \land packetsSync = [s \in Server \mapsto$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad [notCommitted \mapsto \langle\rangle,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad committed \qquad \mapsto \langle\rangle]]$

$InitVerifyVars \triangleq \land proposalMsgsLog \quad\; = \{\}$
$\qquad\qquad\qquad\quad \land epochLeader \qquad\quad\; = [i \in 1 .. MAXEPOCH \mapsto \{\}]$
$\qquad\qquad\qquad\quad \land violatedInvariants \; = [stateInconsistent \qquad\;\; \mapsto \text{FALSE},$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad proposalInconsistent \mapsto \text{FALSE},$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad commitInconsistent \quad\; \mapsto \text{FALSE},$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad ackInconsistent \qquad\;\; \mapsto \text{FALSE},$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad messageIllegal \qquad\;\; \mapsto \text{FALSE}]$

$InitMsgVars \triangleq \land msgs \qquad\quad\; = [s \in Server \mapsto [v \in Server \mapsto \langle\rangle]]$
$\qquad\qquad\qquad \land electionMsgs = [s \in Server \mapsto [v \in Server \mapsto \langle\rangle]]$

$InitRecorder \triangleq recorder = [nTimeout \qquad\;\; \mapsto 0,$
$\qquad\qquad\qquad\qquad\qquad\qquad nTransaction \quad\; \mapsto 0,$
$\qquad\qquad\qquad\qquad\qquad\qquad maxEpoch \qquad\; \mapsto 0,$
$\qquad\qquad\qquad\qquad\qquad\qquad pc \qquad\qquad\qquad \mapsto \langle \text{"Init"} \rangle,$
$\qquad\qquad\qquad\qquad\qquad\qquad nClientRequest \mapsto 0]$

$Init \triangleq \land InitServerVars$

$\land$ *InitLeaderVars*
$\land$ *InitElectionVars*
$\land$ *InitFollowerVars*
$\land$ *InitVerifyVars*
$\land$ *InitMsgVars*
$\land$ *InitRecorder*

---

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

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

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

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

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

$\land$ UNCHANGED *packetsSync*

$\quad\quad\quad \lor \ \land newState = FOLLOWING$

$\quad\quad\quad\quad \land ZabTurnToFollowing(i)$

$\quad\quad\quad\quad \land$ UNCHANGED $\langle learners,\ connecting,\ electing,\ ackldRecv,$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad forwarding,\ tempMaxEpoch\rangle$

$\quad\quad\quad \lor \ \land newState = LOOKING$

$\quad\quad\quad\quad \land$ UNCHANGED $\langle zabState,\ learners,\ connecting,\ electing,\ ackldRecv,$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad forwarding,\ tempMaxEpoch,\ packetsSync,\ initialHistory\rangle$

$\quad \land$ UNCHANGED $\langle lastCommitted,\ acceptedEpoch,\ leaderAddr,\ verifyVars,\ msgs\rangle$

$\quad \land UpdateRecorder(\langle\text{“FLEHandleNotmsg''},\ i\rangle)$

On the premise that *ReceiveVotes.HasQuorums* = TRUE,
corresponding to logic in line $1050 - 1055$ in *FastLeaderElection*.

$FLEWaitNewNotmsg(i) \ \triangleq$

$\quad\quad \land WaitNewNotmsg(i)$

$\quad\quad \land$ UNCHANGED $\langle zabState,\ acceptedEpoch,\ lastCommitted,\ learners,\ connecting,$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad electing,\ ackldRecv,\ forwarding,\ tempMaxEpoch,\ initialHistory,$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad followerVars,\ verifyVars,\ msgs\rangle$

$\quad\quad \land UpdateRecorder(\langle\text{“FLEWaitNewNotmsg''},\ i\rangle)$

On the premise that *ReceiveVotes.HasQuorums* = TRUE,
corresponding to logic in line $1061 - 1066$ in *FastLeaderElection*.

$FLEWaitNewNotmsgEnd(i) \ \triangleq$

$\quad\quad \land WaitNewNotmsgEnd(i)$

$\quad\quad \land$ LET $newState \ \triangleq\ state'[i]$

$\quad\quad\quad$ IN

$\quad\quad\quad \lor \ \land newState = LEADING$

$\quad\quad\quad\quad \land ZabTurnToLeading(i)$

$\quad\quad\quad\quad \land$ UNCHANGED *packetsSync*

$\quad\quad\quad \lor \ \land newState = FOLLOWING$

$\quad\quad\quad\quad \land ZabTurnToFollowing(i)$

$\quad\quad\quad\quad \land$ UNCHANGED $\langle learners,\ connecting,\ electing,\ ackldRecv,\ forwarding,$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad tempMaxEpoch\rangle$

$\quad\quad\quad \lor \ \land newState = LOOKING$

$\quad\quad\quad\quad \land PrintT(\text{“Note: New state is LOOKING in FLEWaitNewNotmsgEnd,''} \ \circ$

$\quad\quad\quad\quad\quad\quad \text{`` which should not happen.''})$

$\quad\quad\quad\quad \land$ UNCHANGED $\langle zabState,\ learners,\ connecting,\ electing,\ ackldRecv,$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad forwarding,\ tempMaxEpoch,\ initialHistory,\ packetsSync\rangle$

$\quad\quad \land$ UNCHANGED $\langle lastCommitted,\ acceptedEpoch,\ leaderAddr,\ verifyVars,\ msgs\rangle$

$\quad\quad \land UpdateRecorder(\langle\text{“FLEWaitNewNotmsgEnd''},\ i\rangle)$

---

$InitialVotes \ \triangleq\ [vote \quad\ \mapsto InitialVote,$

$\quad\quad\quad\quad\quad\quad\quad round \quad \mapsto 0,$

$\quad\quad\quad\quad\quad\quad\quad state \quad\ \mapsto LOOKING,$

$\quad\quad\quad\quad\quad\quad\quad version \mapsto 0]$

$ZabTimeoutInCluster(S) \triangleq$

$\quad \wedge state' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } LOOKING \text{ ELSE } state[s]]$

$\quad \wedge lastProcessed' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } InitLastProcessed(s)$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{ELSE } lastProcessed[s]]$

$\quad \wedge logicalClock' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } logicalClock[s] + 1$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{ELSE } logicalClock[s]]$

$\quad \wedge currentVote' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN}$
$\qquad\qquad\qquad\qquad\qquad [proposedLeader \mapsto s,$
$\qquad\qquad\qquad\qquad\qquad\ proposedZxid \quad \mapsto lastProcessed'[s].zxid,$
$\qquad\qquad\qquad\qquad\qquad\ proposedEpoch \mapsto currentEpoch[s]]$
$\qquad\qquad\qquad\qquad\qquad \text{ELSE } currentVote[s]]$

$\quad \wedge receiveVotes' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } [v \in Server \mapsto InitialVotes]$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{ELSE } receiveVotes[s]]$

$\quad \wedge outOfElection' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } [v \in Server \mapsto InitialVotes]$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{ELSE } outOfElection[s]]$

$\quad \wedge recvQueue' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } \langle[mtype \mapsto NONE]\rangle$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{ELSE } recvQueue[s]]$

$\quad \wedge waitNotmsg' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN FALSE ELSE } waitNotmsg[s]]$

$\quad \wedge leadingVoteSet' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } \{\} \text{ ELSE } leadingVoteSet[s]]$

$\quad \wedge \text{UNCHANGED } \langle electionMsgs, currentEpoch, history \rangle$

$\quad \wedge zabState' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } ELECTION \text{ ELSE } zabState[s]]$

$\quad \wedge leaderAddr' = [s \in Server \mapsto \text{IF } s \in S \text{ THEN } NullPoint \text{ ELSE } leaderAddr[s]]$

$\quad \wedge CleanInputBufferInCluster(S)$

$FollowerShutdown(i) \triangleq$

$\quad \wedge ZabTimeout(i)$

$\quad \wedge zabState' \quad = [zabState \quad \text{EXCEPT } ![i] = ELECTION]$

$\quad \wedge leaderAddr' = [leaderAddr \text{ EXCEPT } ![i] = NullPoint]$

$\quad \wedge CleanInputBuffer(i)$

$LeaderShutdown(i) \triangleq$

$\quad \wedge \text{LET } cluster \triangleq \{i\} \cup learners[i]$
$\qquad \text{IN} \quad ZabTimeoutInCluster(cluster)$

$\quad \wedge learners' \quad = [learners \quad \text{EXCEPT } ![i] = \{\}]$

$\quad \wedge forwarding' = [forwarding \text{ EXCEPT } ![i] = \{\}]$

$RemoveElecting(set, sid) \triangleq$

$\quad \text{LET } sid\_electing \quad \triangleq \{s.sid : s \in set\}$

$\quad \text{IN} \quad \text{IF } sid \notin sid\_electing \text{ THEN } set$

$\qquad \text{ELSE } \text{LET } info \triangleq \text{CHOOSE } s \in set : s.sid = sid$
$\qquad\qquad\qquad new\_info \triangleq [sid \qquad\qquad \mapsto sid,$
$\qquad\qquad\qquad\qquad\qquad\qquad peerLastZxid \mapsto \langle -1, -1 \rangle,$
$\qquad\qquad\qquad\qquad\qquad\qquad inQuorum \quad \mapsto info.inQuorum]$
$\qquad\qquad \text{IN} \quad (set \setminus \{info\}) \cup \{new\_info\}$

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

---

Follower connecting to leader fails and truns to *LOOKING*.
$FollowerTimeout(i) \triangleq$
$\quad \wedge CheckTimeout$ test restrictions of *timeout_1*
$\quad \wedge IsFollower(i)$
$\quad \wedge HasNoLeader(i)$
$\quad \wedge FollowerShutdown(i)$
$\quad \wedge CleanInputBuffer(i)$
$\quad \wedge \text{UNCHANGED} \ \langle acceptedEpoch, lastCommitted, learners, connecting, electing,$
$\quad\quad\quad\quad\quad\quad\quad\quad ackldRecv, forwarding, tempMaxEpoch, initialHistory,$
$\quad\quad\quad\quad\quad\quad\quad\quad verifyVars, packetsSync \rangle$
$\quad \wedge UpdateRecorder(\langle \text{"FollowerTimeout"}, i \rangle)$

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

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

$Restart(i) \triangleq$
$\quad \wedge \vee \wedge IsLooking(i)$
$\quad\quad\quad \wedge$
$\quad\quad \vee \wedge IsLeader(i)$

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

$ConnectAndFollowerSendFOLLOWERINFO(i, j) \triangleq$
   ∧ *IsLeader*(*i*) ∧ ¬*IsMyLearner*(*i*, *j*)
   ∧ *IsFollower*(*j*) ∧ *HasNoLeader*(*j*) ∧ *MyVote*(*j*) = *i*
   ∧ *learners'*  = [*learners*   EXCEPT ![*i*] = *learners*[*i*] ∪ {*j*}]
   ∧ *leaderAddr'* = [*leaderAddr* EXCEPT ![*j*] = *i*]
   ∧ *Send*(*j*, *leaderAddr'*[*j*], [*mtype* ↦ *FOLLOWERINFO*,
               *mzxid* ↦ ⟨*acceptedEpoch*[*j*], 0⟩])
   ∧ UNCHANGED ⟨*serverVars*, *electionVars*, *leadingVoteSet*, *connecting*,
           *electing*, *ackldRecv*, *forwarding*, *tempMaxEpoch*,
           *verifyVars*, *electionMsgs*, *packetsSync*⟩
   ∧ *UpdateRecorder*(⟨"ConnectAndFollowerSendFOLLOWERINFO", *i*, *j*⟩)

*waitingForNewEpoch* in Leader
$WaitingForNewEpoch(i) \triangleq (i \in connecting[i] \land IsQuorum(connecting[i])) = \text{FALSE}$
$WaitingForNewEpochTurnToFalse(i) \triangleq \land i \in connecting'[i]$
                    ∧ *IsQuorum*(*connecting'*[*i*])

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

$LeaderProcessFOLLOWERINFO(i, j) \triangleq$
   ∧ *CheckEpoch*    test restrictions of max epoch
   ∧ *IsLeader*(*i*)
   ∧ *PendingFOLLOWERINFO*(*i*, *j*)
   ∧ LET *msg* $\triangleq$ *msgs*[*j*][*i*][1]
     *infoOk* $\triangleq$ *IsMyLearner*(*i*, *j*)
     *lastAcceptedEpoch* $\triangleq$ *msg.mzxid*[1]
   IN
   ∧ *infoOk*
   ∧ ∨   1. has not broadcast *LEADERINFO*
     ∧ *WaitingForNewEpoch*(*i*)
     ∧ ∨ ∧ *zabState*[*i*] = *DISCOVERY*
       ∧ UNCHANGED *violatedInvariants*
      ∨ ∧ *zabState*[*i*] ≠ *DISCOVERY*
       ∧ *PrintT*("Exception: waitingFotNewEpoch true," ∘
        " while zabState not DISCOVERY.")
       ∧ *violatedInvariants'* = [*violatedInvariants* EXCEPT !.*stateInconsistent* = TRUE]
     ∧ *tempMaxEpoch'* = [*tempMaxEpoch* EXCEPT ![*i*] = IF *lastAcceptedEpoch* ≥ *tempMaxEpoch*[*i*]
                       THEN *lastAcceptedEpoch* + 1
                       ELSE @]

14

$$\wedge\ connecting' \qquad = [connecting \quad \text{EXCEPT}\ ![i] \quad = @ \cup \{j\}]$$
$$\wedge\ \vee\ \wedge\ WaitingForNewEpochTurnToFalse(i)$$
$$\qquad \wedge\ acceptedEpoch' = [acceptedEpoch\ \text{EXCEPT}\ ![i] = tempMaxEpoch'[i]]$$
$$\qquad \wedge\ \text{LET}\ newLeaderZxid\ \triangleq\ \langle acceptedEpoch'[i],\, 0\rangle$$
$$\qquad\qquad m\ \triangleq\ [mtype \mapsto LEADERINFO,$$
$$\qquad\qquad\qquad mzxid \mapsto newLeaderZxid]$$
$$\qquad\quad \text{IN}\quad DiscardAndBroadcastLEADERINFO(i,\, j,\, m)$$
$$\quad \vee\ \wedge\ \neg WaitingForNewEpochTurnToFalse(i)$$
$$\qquad \wedge\ Discard(j,\, i)$$
$$\qquad \wedge\ \text{UNCHANGED}\ acceptedEpoch$$
$$\vee\quad \boxed{2.\ \text{has broadcast}\ LEADERINFO}$$
$$\quad \wedge\ \neg WaitingForNewEpoch(i)$$
$$\quad \wedge\ Reply(i,\, j,\, [mtype \mapsto LEADERINFO,$$
$$\qquad\qquad\qquad mzxid \mapsto \langle acceptedEpoch[i],\, 0\rangle])$$
$$\quad \wedge\ \text{UNCHANGED}\ \langle tempMaxEpoch,\, connecting,\, acceptedEpoch,\, violatedInvariants\rangle$$
$$\wedge\ \text{UNCHANGED}\ \langle state,\, currentEpoch,\, lastProcessed,\, zabState,\, history,\, lastCommitted,$$
$$\qquad\qquad followerVars,\, electionVars,\, initialHistory,\, leadingVoteSet,\, learners,$$
$$\qquad\qquad electing,\, ackldRecv,\, forwarding,\, proposalMsgsLog,\, epochLeader,$$
$$\qquad\qquad electionMsgs\rangle$$
$$\wedge\ UpdateRecorder(\langle\text{"LeaderProcessFOLLOWERINFO"},\, i,\, j\rangle)$$

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

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

$$mzxid \mapsto lastProcessed[i].zxid,$$
$$mepoch \mapsto -1]$$

$\quad$ IN $\quad Reply(i, j, m)$
$\quad\quad \wedge$ UNCHANGED $acceptedEpoch$
$\quad\quad \wedge zabState' = [zabState \text{ EXCEPT } ![i] = SYNCHRONIZATION]$
$\quad\quad \wedge$ UNCHANGED $violatedInvariants$
$\quad \vee \ \wedge \neg stateOk$
$\quad\quad \wedge PrintT(\text{"Exception: Follower receives LEADERINFO," } \circ$
$\quad\quad\quad \text{" whileZabState not DISCOVERY."})$
$\quad\quad \wedge violatedInvariants' = [violatedInvariants \text{ EXCEPT } !.stateInconsistent = \text{TRUE}]$
$\quad\quad \wedge Discard(j, i)$
$\quad\quad \wedge$ UNCHANGED $\langle acceptedEpoch, zabState \rangle$
$\quad \wedge$ UNCHANGED $\langle varsL, leaderAddr, learners, forwarding, electing \rangle$
$\vee \ \boxed{\text{2. Abnormal case - go back to election}}$
$\quad \wedge \neg epochOk$
$\quad \wedge FollowerShutdown(i)$
$\quad \wedge Clean(i, leaderAddr[i])$
$\quad \wedge RemoveLearner(leaderAddr[i], i)$
$\quad \wedge$ UNCHANGED $\langle acceptedEpoch, violatedInvariants \rangle$
$\wedge$ UNCHANGED $\langle history, lastCommitted, connecting, ackldRecv, tempMaxEpoch,$
$\quad\quad\quad initialHistory, proposalMsgsLog, epochLeader, packetsSync \rangle$
$\wedge UpdateRecorder(\langle \text{"FollowerProcessLEADERINFO"}, i, j \rangle)$

---

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

$\boxed{\text{There originally existed one bug in } LeaderProcessACK \text{ when}}$
$\boxed{monotonicallyInc = \text{FALSE, and it is we did not add } ackSid \text{ of}}$
$\boxed{\text{history in } SYNC. \text{ So we update } ackSid \text{ in } syncFollower.}$
$UpdateAckSid(his, lastSeenIndex, target) \triangleq$
$\quad$ IF $Len(his) = 0 \vee lastSeenIndex = 0$ THEN $his$
$\quad\quad$ ELSE $\ UpdateAckSidHelper(his, 1, Minimum(\{Len(his), lastSeenIndex\}), target)$

$\boxed{\text{return } -1: \text{ this } zxid \text{ appears at least twice; } Len(his) + 1: \text{ does not exist;}}$
$\boxed{1 \ \neg Len(his): \text{ exists and appears just once.}}$
RECURSIVE $ZxidToIndexHepler(\_, \_, \_, \_)$
$ZxidToIndexHepler(his, zxid, cur, appeared) \triangleq$
$\quad$ IF $cur > Len(his)$ THEN $cur$

16

ELSE IF $TxnZxidEqual(his[cur], zxid)$
    THEN CASE $appeared = \text{TRUE} \to -1$
      □    OTHER             $\to Minimum(\{cur,$
                      $ZxidToIndexHepler(his, zxid, cur + 1, \text{TRUE})\})$
    ELSE $ZxidToIndexHepler(his, zxid, cur + 1, appeared)$

$ZxidToIndex(his, zxid) \triangleq$ IF $ZxidEqual(zxid, \langle 0, 0 \rangle)$ THEN $0$
        ELSE IF $Len(his) = 0$ THEN $1$
            ELSE LET $len \triangleq Len(his)$ IN
                IF $\exists\, idx \in 1\,..\,len : TxnZxidEqual(his[idx], zxid)$
                THEN $ZxidToIndexHepler(his, zxid, 1, \text{FALSE})$
                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 $cur > committed \to$
        LET $m\_proposal \triangleq [mtype \mapsto PROPOSAL,$
                        $mzxid \mapsto his[cur].zxid,$
                        $mdata \mapsto his[cur].value]$
        IN  $queuePackets(Append(queue, m\_proposal), his, cur + 1, committed, end)$
      □  $cur \leq committed \to$
        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 $setPacketsForChecking(\_, \_, \_, \_, \_, \_)$

17

$setPacketsForChecking(set, src, ep, his, cur, end) \triangleq$
  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 $setPacketsForChecking((set \cup \{m\_proposal\}), src, ep, his, cur + 1, end)$

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

$SyncFollower(i, j, peerLastZxid, needRemoveHead) \triangleq$
  LET $IsPeerNewEpochZxid \triangleq peerLastZxid[2] = 0$
     $lastProcessedZxid \triangleq lastProcessed[i].zxid$

$$maxCommittedLog \triangleq \text{ IF } zabState[i] = BROADCAST$$
$$\text{THEN } lastCommitted[i].zxid$$
$$\text{ELSE } \text{ LET } totalLen \triangleq Len(initialHistory[i])$$
$$\text{IN } \text{ IF } totalLen = 0 \text{ THEN } \langle 0, 0 \rangle$$
$$\text{ELSE } history[i][totalLen].zxid$$

Hypothesis: 1. $minCommittedLog$ : $zxid$ of head of history, so no SNAP.

2. $maxCommittedLog = lastCommitted$, to compress state space.

3. merge $queueCommittedProposals$,$startForwarding$ and sending $NEWLEADER$ into $StartForwarding$.

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

$$\wedge\ StartForwarding(i,\ j,\ lastZxid,\ lastIndex,$$
$$TRUNC,\ needRemoveHead)$$

$IsMoreRecentThan(ss1,\ ss2) \triangleq\ \vee\ ss1.currentEpoch > ss2.currentEpoch$
$$\vee\ \wedge\ ss1.currentEpoch = ss2.currentEpoch$$
$$\wedge\ ZxidCompare(ss1.lastZxid,\ ss2.lastZxid)$$

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

$UpdateElecting(oldSet,\ sid,\ peerLastZxid,\ inQuorum) \triangleq$
  LET $sid\_electing \triangleq \{s.sid : s \in oldSet\}$
  IN  IF $sid \in sid\_electing$
    THEN LET $old\_info \triangleq$ CHOOSE $info \in oldSet : info.sid = sid$
       $follower\_info \triangleq$
        $[sid \qquad\qquad \mapsto sid,$
        $peerLastZxid \mapsto peerLastZxid,$
        $inQuorum \quad\ \mapsto (inQuorum \vee old\_info.inQuorum)]$
      IN $(oldSet \setminus \{old\_info\}) \cup \{follower\_info\}$
    ELSE LET $follower\_info \triangleq$
        $[sid \qquad\qquad \mapsto sid,$
        $peerLastZxid \mapsto peerLastZxid,$
        $inQuorum \quad\ \mapsto inQuorum]$
      IN $oldSet \cup \{follower\_info\}$

$LeaderTurnToSynchronization(i) \triangleq$
  $\wedge\ currentEpoch' = [currentEpoch$ EXCEPT $![i] = acceptedEpoch[i]]$
  $\wedge\ zabState' \qquad = [zabState \qquad$ EXCEPT $![i]\ = SYNCHRONIZATION]$

$LeaderProcessACKEPOCH(i,\ j) \triangleq$
  $\wedge\ IsLeader(i)$
  $\wedge\ PendingACKEPOCH(i,\ j)$
  $\wedge$ LET $msg \triangleq msgs[j][i][1]$
    $infoOk \triangleq IsMyLearner(i,\ j)$
    $leaderStateSummary \qquad \triangleq [currentEpoch \mapsto currentEpoch[i],$

$$followerStateSummary \triangleq [currentEpoch \mapsto msg.mepoch,$$
$$lastZxid \quad \mapsto lastProcessed[i].zxid]$$
$$lastZxid \quad \mapsto msg.mzxid]$$

$logOk \triangleq$ whether follower is no more up-to-date than leader
$\quad\quad \neg IsMoreRecentThan(followerStateSummary, leaderStateSummary)$

$electing\_quorum \triangleq \{e \in electing[i] : e.inQuorum = \text{TRUE}\}$

$sid\_electing \triangleq \{s.sid : s \in electing\_quorum\}$

IN $\quad \wedge infoOk$

$\quad\quad \wedge \vee$ $electionFinished =$ true, jump ouf of $waitForEpochAck$.

$\quad\quad\quad\quad$ Different from code, here we still need to record $info$

$\quad\quad\quad\quad$ into electing, to help us perform $syncFollower$ afterwards.

$\quad\quad\quad\quad$ Since electing already meets quorum, it does not break

$\quad\quad\quad\quad$ consistency between code and spec.

$\quad\quad\quad\quad \wedge ElectionFinished(i, sid\_electing)$

$\quad\quad\quad\quad \wedge electing' = [electing \text{ EXCEPT } ![i] = UpdateElecting(@, j, msg.mzxid, \text{FALSE})]$

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

$\quad\quad\quad\quad \wedge \text{UNCHANGED } \langle varsL, zabState, forwarding, leaderAddr,$
$\quad\quad\quad\quad\quad\quad\quad\quad learners, epochLeader, violatedInvariants \rangle$

$\quad\quad\quad \vee \wedge \neg ElectionFinished(i, sid\_electing)$

$\quad\quad\quad\quad \wedge \vee \wedge zabState[i] = DISCOVERY$

$\quad\quad\quad\quad\quad\quad \wedge \text{UNCHANGED } violatedInvariants$

$\quad\quad\quad\quad\quad \vee \wedge zabState[i] \neq DISCOVERY$

$\quad\quad\quad\quad\quad\quad \wedge PrintT(\text{"Exception: electionFinished false,"} \circ$
$\quad\quad\quad\quad\quad\quad\quad \text{" while zabState not DISCOVERY."})$

$\quad\quad\quad\quad\quad\quad \wedge violatedInvariants' = [violatedInvariants \text{ EXCEPT}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad !.stateInconsistent = \text{TRUE}]$

$\quad\quad\quad\quad \wedge \vee \wedge followerStateSummary.currentEpoch = -1$

$\quad\quad\quad\quad\quad\quad \wedge electing' = [electing \text{ EXCEPT } ![i] \quad = UpdateElecting(@, j,$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad msg.mzxid, \text{FALSE})]$

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

$\quad\quad\quad\quad\quad\quad \wedge \text{UNCHANGED } \langle varsL, zabState, forwarding, leaderAddr,$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad learners, epochLeader \rangle$

$\quad\quad\quad\quad\quad \vee \wedge followerStateSummary.currentEpoch > -1$

$\quad\quad\quad\quad\quad\quad \wedge \vee$ normal follower

$\quad\quad\quad\quad\quad\quad\quad\quad \wedge logOk$

$\quad\quad\quad\quad\quad\quad\quad\quad \wedge electing' = [electing \text{ EXCEPT } ![i] =$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad UpdateElecting(@, j, msg.mzxid, \text{TRUE})]$

$\quad\quad\quad\quad\quad\quad\quad\quad \wedge \text{LET } new\_electing\_quorum \triangleq \{e \in electing'[i] : e.inQuorum = \text{TRUE}\}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad new\_sid\_electing \triangleq \{s.sid : s \in new\_electing\_quorum\}$

$\quad\quad\quad\quad\quad\quad\quad\quad \text{IN}$

$\quad\quad\quad\quad\quad\quad\quad\quad \vee$ $electionFinished =$ true, jump out of $waitForEpochAck$,

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ update $currentEpoch$ and $zabState$.

$\quad\quad\quad\quad\quad\quad\quad\quad\quad \wedge ElectionFinished(i, new\_sid\_electing)$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad \wedge LeaderTurnToSynchronization(i)$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad \wedge \text{LET } newLeaderEpoch \triangleq acceptedEpoch[i]$

$$\text{IN} \quad epochLeader' = [epochLeader \text{ EXCEPT } ![newLeaderEpoch]$$
$$= @ \cup \{i\}] \text{ for checking invariants}$$

$\vee$ there still exists $electionFinished =$ false.
$\quad \wedge \neg ElectionFinished(i, new\_sid\_electing)$
$\quad \wedge \text{UNCHANGED } \langle currentEpoch, zabState, epochLeader \rangle$

$\wedge Discard(j, i)$
$\wedge \text{UNCHANGED } \langle state, lastProcessed, electionVars, leadingVoteSet,$
$\qquad\qquad\qquad\qquad electionMsgs, leaderAddr, learners, history, forwarding \rangle$

$\vee$ Exists follower more recent than leader
$\quad \wedge \neg logOk$
$\quad \wedge LeaderShutdown(i)$
$\quad \wedge \text{UNCHANGED } \langle electing, epochLeader \rangle$

$\wedge \text{UNCHANGED } \langle acceptedEpoch, lastCommitted, connecting, ackldRecv,$
$\qquad\qquad\qquad tempMaxEpoch, initialHistory, packetsSync, proposalMsgsLog \rangle$
$\wedge UpdateRecorder(\langle \text{``LeaderProcessACKEPOCH''}, i, j \rangle)$

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

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

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

$FollowerProcessSyncMessage(i, j) \triangleq$
    $\wedge IsFollower(i)$
    $\wedge msgs[j][i] \neq \langle\rangle$
    $\wedge msgs[j][i][1].mtype = DIFF \vee msgs[j][i][1].mtype = TRUNC$
    $\wedge$ LET $msg \triangleq msgs[j][i][1]$
           $infoOk \triangleq IsMyLeader(i, j)$
           $stateOk \triangleq zabState[i] = SYNCHRONIZATION$
      IN   $\wedge infoOk$
          $\wedge \vee$ Follower should receive packets in $SYNC$.
              $\wedge \neg stateOk$
              $\wedge PrintT($ "Exception: Follower receives DIFF/TRUNC," $\circ$
                    " whileZabState not SYNCHRONIZATION." $)$
              $\wedge violatedInvariants' = [violatedInvariants$ EXCEPT $!.stateInconsistent = $ TRUE$]$
              $\wedge$ UNCHANGED $\langle history, initialHistory, lastProcessed, lastCommitted\rangle$
           $\vee \wedge stateOk$
             $\wedge \vee \wedge msg.mtype = DIFF$
                 $\wedge$ UNCHANGED $\langle history, initialHistory, lastProcessed, lastCommitted,$
                        $violatedInvariants\rangle$
               $\vee \wedge msg.mtype = TRUNC$
                 $\wedge$ LET $truncZxid \triangleq msg.mtruncZxid$
                      $truncIndex \triangleq ZxidToIndex(history[i], truncZxid)$
                  IN
                 $\vee \wedge truncIndex > Len(history[i])$
                   $\wedge PrintT($ "Exception: TRUNC error." $)$
                   $\wedge violatedInvariants' = [violatedInvariants$ EXCEPT
                          $!.proposalInconsistent = $ TRUE$]$
                   $\wedge$ UNCHANGED $\langle history, initialHistory, lastProcessed, lastCommitted\rangle$
                 $\vee \wedge truncIndex \leq Len(history[i])$
                   $\wedge history' = [history$ EXCEPT
                          $![i] = TruncateLog(history[i], truncIndex)]$
                   $\wedge initialHistory' = [initialHistory$ EXCEPT $![i] = history'[i]]$
                   $\wedge lastProcessed' = [lastProcessed$ EXCEPT
                          $![i] = [index \mapsto truncIndex,$
                             $zxid \mapsto truncZxid]]$
                   $\wedge lastCommitted' = [lastCommitted$ EXCEPT
                          $![i] = [index \mapsto truncIndex,$
                             $zxid \mapsto truncZxid]]$
                   $\wedge$ UNCHANGED $violatedInvariants$
         $\wedge Discard(j, i)$
         $\wedge$ UNCHANGED $\langle state, currentEpoch, zabState, acceptedEpoch, electionVars,$
                 $leaderVars, tempMaxEpoch, followerVars,$
                 $proposalMsgsLog, epochLeader, electionMsgs\rangle$

$\land$ *UpdateRecorder*($\langle$"FollowerProcessSyncMessage", $i$, $j\rangle$)

$LastProposed(i) \triangleq$ IF $Len(history[i]) = 0$ THEN $[index \mapsto 0,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad zxid \quad \mapsto \langle 0, 0\rangle]$
$\qquad\qquad\qquad\qquad$ ELSE
$\qquad\qquad\qquad\qquad$ LET $lastIndex \triangleq Len(history[i])$
$\qquad\qquad\qquad\qquad\qquad entry \qquad \triangleq history[i][lastIndex]$
$\qquad\qquad\qquad\qquad$ IN $\quad [index \mapsto lastIndex,$
$\qquad\qquad\qquad\qquad\qquad\quad zxid \quad \mapsto entry.zxid]$

$LastQueued(i) \triangleq$ IF $\neg IsFollower(i) \lor zabState[i] \neq SYNCHRONIZATION$
$\qquad\qquad\qquad\quad$ THEN $LastProposed(i)$
$\qquad\qquad\qquad\quad$ ELSE $\quad$ condition: $IsFollower(i) \land zabState = SYNCHRONIZATION$
$\qquad\qquad\qquad\qquad\qquad\quad$ LET $packetsInSync \triangleq packetsSync[i].notCommitted$
$\qquad\qquad\qquad\qquad\qquad\qquad lenSync \quad \triangleq Len(packetsInSync)$
$\qquad\qquad\qquad\qquad\qquad\qquad totalLen \quad \triangleq Len(history[i]) + lenSync$
$\qquad\qquad\qquad\qquad\qquad$ IN $\quad$ IF $lenSync = 0$ THEN $LastProposed(i)$
$\qquad\qquad\qquad\qquad\qquad\qquad$ ELSE $[index \mapsto totalLen,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad zxid \quad \mapsto packetsInSync[lenSync].zxid]$

$IsNextZxid(curZxid, nextZxid) \triangleq$
$\qquad\qquad\quad \lor \quad$ first *PROPOSAL* in this epoch
$\qquad\qquad\qquad\quad \land nextZxid[2] = 1$
$\qquad\qquad\qquad\quad \land curZxid[1] < nextZxid[1]$
$\qquad\qquad\quad \lor \quad$ not first *PROPOSAL* in this epoch
$\qquad\qquad\qquad\quad \land nextZxid[2] > 1$
$\qquad\qquad\qquad\quad \land curZxid[1] = nextZxid[1]$
$\qquad\qquad\qquad\quad \land curZxid[2] + 1 = nextZxid[2]$

$FollowerProcessPROPOSALInSync(i, j) \triangleq$
$\qquad\quad \land IsFollower(i)$
$\qquad\quad \land PendingPROPOSAL(i, j)$
$\qquad\quad \land zabState[i] = SYNCHRONIZATION$
$\qquad\quad \land$ LET $msg \triangleq msgs[j][i][1]$
$\qquad\qquad\qquad infoOk \quad \triangleq IsMyLeader(i, j)$
$\qquad\qquad\qquad isNext \quad \triangleq IsNextZxid(LastQueued(i).zxid, msg.mzxid)$
$\qquad\qquad\qquad newTxn \triangleq [zxid \quad \mapsto msg.mzxid,$
$\qquad\qquad\qquad\qquad\qquad\qquad value \quad \mapsto msg.mdata,$
$\qquad\qquad\qquad\qquad\qquad\qquad ackSid \mapsto \{\}, \qquad$ follower do not consider *ackSid*
$\qquad\qquad\qquad\qquad\qquad\qquad epoch \quad \mapsto acceptedEpoch[i]]$ epoch of this round
$\qquad\quad$ IN $\quad \land infoOk$
$\qquad\qquad\qquad \land \lor \land isNext$
$\qquad\qquad\qquad\qquad\qquad \land packetsSync' = [packetsSync$ EXCEPT $![i].notCommitted$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad = Append(packetsSync[i].notCommitted, newTxn)]$

$\lor \land \neg isNext$
$\quad \land PrintT(\text{``Warn: Follower receives PROPOSAL,''} \circ$
$\qquad \text{`` while zxid != lastQueued + 1.''})$
$\quad \land \text{UNCHANGED } packetsSync$

$\land Discard(j, i)$
$\land \text{UNCHANGED } \langle serverVars, electionVars, leaderVars, leaderAddr,$
$\qquad\qquad\qquad verifyVars, electionMsgs \rangle$
$\land UpdateRecorder(\langle \text{``FollowerProcessPROPOSALInSync''}, i, j \rangle)$

$\text{RECURSIVE } IndexOfFirstTxnWithEpoch(\_, \_, \_, \_)$
$IndexOfFirstTxnWithEpoch(his, epoch, cur, end) \triangleq$
$\qquad \text{IF } cur > end \text{ THEN } cur$
$\qquad\quad \text{ELSE IF } his[cur].epoch = epoch \text{ THEN } cur$
$\qquad\qquad\quad \text{ELSE } IndexOfFirstTxnWithEpoch(his, epoch, cur + 1, end)$

$LastCommitted(i) \triangleq \text{IF } zabState[i] = BROADCAST \text{ THEN } lastCommitted[i]$
$\qquad\qquad\qquad \text{ELSE CASE } IsLeader(i) \quad \to$
$\qquad\qquad\qquad\qquad \text{LET } lastInitialIndex \triangleq Len(initialHistory[i])$
$\qquad\qquad\qquad\qquad \text{IN } \quad \text{IF } lastInitialIndex = 0 \text{ THEN } [index \mapsto 0,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad zxid \quad \mapsto \langle 0, 0 \rangle]$
$\qquad\qquad\qquad\qquad\qquad \text{ELSE } [index \mapsto lastInitialIndex,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad zxid \quad \mapsto history[i][lastInitialIndex].zxid]$
$\qquad\qquad\qquad\quad \square \quad IsFollower(i) \to$
$\qquad\qquad\qquad\qquad \text{LET } completeHis \triangleq history[i] \circ packetsSync[i].notCommitted$
$\qquad\qquad\qquad\qquad\qquad packetsCommitted \triangleq packetsSync[i].committed$
$\qquad\qquad\qquad\qquad\qquad lenCommitted \triangleq Len(packetsCommitted)$
$\qquad\qquad\qquad\qquad \text{IN } \quad \text{IF } lenCommitted = 0 \text{ return last one in initial history}$
$\qquad\qquad\qquad\qquad\qquad \text{THEN LET } lastInitialIndex \triangleq Len(initialHistory[i])$
$\qquad\qquad\qquad\qquad\qquad\qquad \text{IN } \quad \text{IF } lastInitialIndex = 0$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{THEN } [index \mapsto 0,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad zxid \quad \mapsto \langle 0, 0 \rangle]$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{ELSE } [index \mapsto lastInitialIndex,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad zxid \quad \mapsto completeHis[lastInitialIndex].zxid]$
$\qquad\qquad\qquad\qquad\qquad \text{ELSE } \qquad\qquad\qquad\qquad \text{return tail of } packetsCommitted$
$\qquad\qquad\qquad\qquad\qquad\qquad \text{LET } committedIndex \triangleq ZxidToIndex(completeHis,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad packetsCommitted[lenCommitted])$
$\qquad\qquad\qquad\qquad\qquad\qquad \text{IN } \quad [index \mapsto committedIndex,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad zxid \quad \mapsto packetsCommitted[lenCommitted]]$
$\qquad\qquad\qquad\quad \square \quad \text{OTHER } \to lastCommitted[i]$

$TxnWithIndex(i, idx) \triangleq \text{IF } \neg IsFollower(i) \lor zabState[i] \neq SYNCHRONIZATION$
$\qquad\qquad\qquad\qquad \text{THEN } history[i][idx]$
$\qquad\qquad\qquad\qquad \text{ELSE LET } completeHis \triangleq history[i] \circ packetsSync[i].notCommitted$
$\qquad\qquad\qquad\qquad\qquad \text{IN } \quad completeHis[idx]$

$FollowerProcessCOMMITInSync(i, j) \triangleq$

$\qquad \wedge IsFollower(i)$

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

$\qquad \wedge zabState[i] = SYNCHRONIZATION$

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

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

$\qquad\qquad\quad committedIndex \triangleq LastCommitted(i).index + 1$

$\qquad\qquad\quad exist \triangleq \wedge committedIndex \leq LastQueued(i).index$

$\qquad\qquad\qquad\qquad\quad \wedge IsNextZxid(LastCommitted(i).zxid, msg.mzxid)$

$\qquad\qquad\quad match \triangleq ZxidEqual(msg.mzxid, TxnWithIndex(i, committedIndex).zxid)$

$\qquad\quad \text{IN} \quad \wedge infoOk$

$\qquad\qquad\qquad \wedge \vee \wedge exist$

$\qquad\qquad\qquad\qquad \vee \wedge match$

$\qquad\qquad\qquad\qquad\qquad \wedge packetsSync' = [packetsSync \text{ EXCEPT }![i].committed$

$\qquad\qquad\qquad\qquad\qquad\qquad = Append(packetsSync[i].committed, msg.mzxid)]$

$\qquad\qquad\qquad\qquad\qquad \wedge \text{UNCHANGED } violatedInvariants$

$\qquad\qquad\qquad\qquad \vee \wedge \neg match$

$\qquad\qquad\qquad\qquad\qquad \wedge PrintT(\text{"Warn: Follower receives COMMIT,"} \circ$

$\qquad\qquad\qquad\qquad\qquad\qquad \text{" but zxid not the next committed zxid in COMMIT."})$

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

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

$\qquad\qquad\qquad\qquad\qquad \wedge \text{UNCHANGED } packetsSync$

$\qquad\qquad\qquad\quad \vee \wedge \neg exist$

$\qquad\qquad\qquad\qquad\quad \wedge PrintT(\text{"Warn: Follower receives COMMIT,"} \circ$

$\qquad\qquad\qquad\qquad\qquad \text{" but no packets with its zxid exists."})$

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

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

$\qquad\qquad\qquad\qquad\quad \wedge \text{UNCHANGED } packetsSync$

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

$\qquad \wedge \text{UNCHANGED } \langle serverVars, electionVars, leaderVars,$

$\qquad\qquad\qquad\qquad\qquad leaderAddr, epochLeader, proposalMsgsLog, electionMsgs\rangle$

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

$\text{RECURSIVE } ACKInBatches(\_, \_)$

$ACKInBatches(queue, packets) \triangleq$

$\qquad \text{IF } packets = \langle\rangle \text{ THEN } queue$

$\qquad \text{ELSE LET } head \triangleq packets[1]$

$\qquad\qquad\qquad newPackets \triangleq Tail(packets)$

$\qquad\qquad\qquad m\_ack \triangleq [mtype \mapsto ACK,$

$\qquad\qquad\qquad\qquad\qquad mzxid \mapsto head.zxid]$

$\qquad\quad \text{IN} \quad ACKInBatches(Append(queue, m\_ack), newPackets)$

$FollowerProcessNEWLEADER(i, j) \triangleq$
  $\wedge\ IsFollower(i)$
  $\wedge\ PendingNEWLEADER(i, j)$
  $\wedge\ \text{LET}\ msg \triangleq msgs[j][i][1]$
     $infoOk \triangleq IsMyLeader(i, j)$
     $packetsInSync \triangleq packetsSync[i].notCommitted$
     $m\_ackld \triangleq [mtype \mapsto ACKLD,$
           $mzxid \mapsto msg.mzxid]$
     $ms\_ack \triangleq ACKInBatches(\langle\rangle, packetsInSync)$
     $queue\_toSend \triangleq \langle m\_ackld\rangle \circ ms\_ack$   send $ACK - NEWLEADER$ first.
   $\text{IN}$ $\wedge\ infoOk$
     $\wedge\ currentEpoch' = [currentEpoch\ \text{EXCEPT}\ ![i] = acceptedEpoch[i]]$
     $\wedge\ history' \qquad = [history \qquad\quad \text{EXCEPT}\ ![i]\ = @ \circ packetsInSync]$
     $\wedge\ packetsSync' = [packetsSync\ \text{EXCEPT}\ ![i].notCommitted = \langle\rangle]$
     $\wedge\ DiscardAndSendPackets(i, j, queue\_toSend)$
  $\wedge\ \text{UNCHANGED}\ \langle state, lastProcessed, zabState, acceptedEpoch, lastCommitted,$
          $electionVars, leaderVars, initialHistory, leaderAddr, verifyVars,$
          $electionMsgs\rangle$
  $\wedge\ UpdateRecorder(\langle\text{"FollowerProcessNEWLEADER"}, i, j\rangle)$

<br>

*quorumFormed* in Leader
$QuorumFormed(i) \triangleq i \in ackldRecv[i] \wedge IsQuorum(ackldRecv[i])$
$QuorumFormedTurnToTrue(i) \triangleq i \in ackldRecv'[i] \wedge IsQuorum(ackldRecv'[i])$

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

<br>

See *startZkServer* in Leader for details.
$StartZkServer(i) \triangleq$
  $\text{LET}\ latest \triangleq LastProposed(i)$
  $\text{IN}\quad \wedge\ lastCommitted' = [lastCommitted\ \text{EXCEPT}\ ![i] = latest]$
    $\wedge\ lastProcessed' = [lastProcessed\ \text{EXCEPT}\ ![i]\ = latest]$
    $\wedge\ UpdateElectionVote(i, acceptedEpoch[i])$

<br>

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

<br>

Leader waits for receiving quorum of $ACK$ whose lower bits of *zxid* is 0, and broadcasts
$UPTODATE$. See *waitForNewLeaderAck* for details.
$LeaderProcessACKLD(i, j) \triangleq$
  $\wedge\ IsLeader(i)$
  $\wedge\ PendingACKLD(i, j)$
  $\wedge\ \text{LET}\ msg \quad \triangleq msgs[j][i][1]$
     $infoOk \triangleq IsMyLearner(i, j)$
     $match \triangleq ZxidEqual(msg.mzxid, \langle acceptedEpoch[i], 0\rangle)$

$$currentZxid \triangleq \langle acceptedEpoch[i], 0 \rangle$$
$$m\_uptodate \triangleq [mtype \mapsto UPTODATE,$$
$$mzxid \mapsto currentZxid] \;\boxed{\text{not important}}$$

IN $\quad \wedge infoOk$

$\quad\quad \wedge \vee \;\boxed{\text{just reply } UPTODATE.}$

$\quad\quad\quad\quad \wedge QuorumFormed(i)$

$\quad\quad\quad\quad \wedge Reply(i, j, m\_uptodate)$

$\quad\quad\quad\quad \wedge$ UNCHANGED $\langle ackldRecv, zabState, lastCommitted, lastProcessed,$

$\quad\quad\quad\quad\quad\quad currentVote, violatedInvariants\rangle$

$\quad\quad\quad \vee \wedge \neg QuorumFormed(i)$

$\quad\quad\quad\quad \wedge \vee \wedge match$

$\quad\quad\quad\quad\quad\quad \wedge ackldRecv' = [ackldRecv \text{ EXCEPT } ![i] = @ \cup \{j\}]$

$\quad\quad\quad\quad\quad\quad \wedge \vee \;\boxed{\text{jump out of } waitForNewLeaderAck, \text{ and do } startZkServer,}$

$\quad\quad\quad\quad\quad\quad\quad\quad \boxed{setZabState, \text{ and reply } UPTODATE.}$

$\quad\quad\quad\quad\quad\quad\quad\quad \wedge QuorumFormedTurnToTrue(i)$

$\quad\quad\quad\quad\quad\quad\quad\quad \wedge LeaderTurnToBroadcast(i)$

$\quad\quad\quad\quad\quad\quad\quad\quad \wedge DiscardAndBroadcastUPTODATE(i, j, m\_uptodate)$

$\quad\quad\quad\quad\quad\quad\quad \vee \;\boxed{\text{still wait in } waitForNewLeaderAck.}$

$\quad\quad\quad\quad\quad\quad\quad\quad \wedge \neg QuorumFormedTurnToTrue(i)$

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

$\quad\quad\quad\quad\quad\quad\quad\quad \wedge$ UNCHANGED $\langle zabState, lastCommitted, lastProcessed, currentVote\rangle$

$\quad\quad\quad\quad\quad\quad \wedge$ UNCHANGED $violatedInvariants$

$\quad\quad\quad\quad\quad \vee \wedge \neg match$

$\quad\quad\quad\quad\quad\quad \wedge PrintT(\text{``Exception: NEWLEADER ACK is from a different epoch.''})$

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

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

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

$\quad\quad\quad\quad\quad\quad \wedge$ UNCHANGED $\langle ackldRecv, zabState, lastCommitted,$

$\quad\quad\quad\quad\quad\quad\quad\quad lastProcessed, currentVote\rangle$

$\quad \wedge$ UNCHANGED $\langle state, currentEpoch, acceptedEpoch, history, logicalClock, receiveVotes,$

$\quad\quad\quad outOfElection, recvQueue, waitNotmsg, leadingVoteSet, learners, connecting,$

$\quad\quad\quad electing, forwarding, tempMaxEpoch, initialHistory, followerVars,$

$\quad\quad\quad proposalMsgsLog, epochLeader, electionMsgs\rangle$

$\quad \wedge UpdateRecorder(\langle \text{``LeaderProcessACKLD''}, i, j\rangle)$

$TxnsWithPreviousEpoch(i) \triangleq$

$\quad$ LET $completeHis \triangleq$ IF $\neg IsFollower(i) \vee zabState[i] \neq SYNCHRONIZATION$

$\quad\quad\quad\quad\quad\quad\quad$ THEN $history[i]$

$\quad\quad\quad\quad\quad\quad\quad$ ELSE $history[i] \circ packetsSync[i].notCommitted$

$\quad\quad end \quad \triangleq Len(completeHis)$

$\quad\quad first \quad \triangleq IndexOfFirstTxnWithEpoch(completeHis, acceptedEpoch[i], 1, end)$

$\quad$ IN $\quad$ IF $first > end$ THEN $completeHis$

$\quad\quad\quad$ ELSE $SubSeq(completeHis, 1, first - 1)$

$TxnsRcvWithCurEpoch(i) \triangleq$

$\text{LET } completeHis \triangleq \text{ IF } \neg IsFollower(i) \vee zabState[i] \neq SYNCHRONIZATION$
$\qquad\qquad\qquad\qquad \text{THEN } history[i]$
$\qquad\qquad\qquad\qquad \text{ELSE } history[i] \circ packetsSync[i].notCommitted$
$\qquad end \quad \triangleq Len(completeHis)$
$\qquad first \quad \triangleq IndexOfFirstTxnWithEpoch(completeHis, acceptedEpoch[i], 1, end)$
$\text{IN} \quad \text{IF } first > end \text{ THEN } \langle\rangle$
$\qquad\quad \text{ELSE } SubSeq(completeHis, first, end) \quad \boxed{completeHis[first:end]}$

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

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

$\text{RECURSIVE } TxnsAndCommittedMatch(\_, \_)$
$TxnsAndCommittedMatch(txns, packetsCommitted) \triangleq$
$\qquad \text{LET } len1 \triangleq Len(txns)$
$\qquad\qquad len2 \triangleq Len(packetsCommitted)$
$\qquad \text{IN} \quad \text{IF } len2 = 0 \text{ THEN TRUE}$
$\qquad\qquad \text{ELSE IF } len1 < len2 \text{ THEN FALSE}$
$\qquad\qquad\qquad \text{ELSE} \quad \wedge ZxidEqual(txns[len1].zxid, packetsCommitted[len2])$
$\qquad\qquad\qquad\qquad\qquad \wedge TxnsAndCommittedMatch(SubSeq(txns, 1, len1 - 1),$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad SubSeq(packetsCommitted, 1, len2 - 1))$

$FollowerLogRequestInBatches(i, leader, ms\_ack, packetsNotCommitted) \triangleq$
$\qquad \wedge history' = [history \text{ EXCEPT } ![i] = @ \circ packetsNotCommitted]$
$\qquad \wedge DiscardAndSendPackets(i, leader, ms\_ack)$

$FollowerCommitInBatches(i) \triangleq$
$\qquad \text{LET } committedTxns \triangleq CommittedTxns(i)$
$\qquad\qquad\quad packetsCommitted \triangleq packetsSync[i].committed$

$$match \triangleq TxnsAndCommittedMatch(committedTxns, packetsCommitted)$$

IN
$$\lor \land match$$
$$\land lastCommitted' = [lastCommitted \text{ EXCEPT } ![i] = LastCommitted(i)]$$
$$\land lastProcessed' \quad = [lastProcessed \text{ EXCEPT } ![i] \quad = lastCommitted'[i]]$$
$$\land \text{UNCHANGED } violatedInvariants$$
$$\lor \land \neg match$$
$$\land PrintT(\text{"Warn: Committing zxid withou see txn. /"} \circ$$
$$\text{"Committing zxid != pending txn zxid."})$$
$$\land violatedInvariants' = [violatedInvariants \text{ EXCEPT }$$
$$!.commitInconsistent = \text{TRUE}]$$
$$\land \text{UNCHANGED } \langle lastCommitted, lastProcessed \rangle$$

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

$$FollowerProcessUPTODATE(i, j) \triangleq$$
$$\land IsFollower(i)$$
$$\land PendingUPTODATE(i, j)$$
$$\land \text{LET } msg \triangleq msgs[j][i][1]$$
$$infoOk \triangleq IsMyLeader(i, j)$$
$$packetsNotCommitted \triangleq packetsSync[i].notCommitted$$
$$ms\_ack \triangleq ACKInBatches(\langle\rangle, packetsNotCommitted)$$
$$\text{IN} \quad \land infoOk$$

Here we ignore ack of *UPTODATE*.

$$\land UpdateElectionVote(i, acceptedEpoch[i])$$
$$\land FollowerLogRequestInBatches(i, j, ms\_ack, packetsNotCommitted)$$
$$\land FollowerCommitInBatches(i)$$
$$\land packetsSync' = [packetsSync \text{ EXCEPT } ![i].notCommitted = \langle\rangle,$$
$$![i].committed = \langle\rangle]$$
$$\land zabState' = [zabState \text{ EXCEPT } ![i] = BROADCAST]$$
$$\land \text{UNCHANGED } \langle state, currentEpoch, acceptedEpoch, logicalClock,$$
$$receiveVotes, outOfElection, recvQueue, waitNotmsg, leaderVars,$$
$$initialHistory, leaderAddr, epochLeader, proposalMsgsLog, electionMsgs \rangle$$
$$\land UpdateRecorder(\langle\text{"FollowerProcessUPTODATE"}, i, j\rangle)$$

---

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

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

$LeaderProcessRequest(i) \triangleq$
  $\wedge\ CheckTransactionNum$ <span style="background:#ccc">test restrictions of transaction num</span>
  $\wedge\ IsLeader(i)$
  $\wedge\ zabState[i] = BROADCAST$
  $\wedge\ \text{LET } request\_value\ \triangleq\ GetRecorder(\text{"nClientRequest"})$ <span style="background:#ccc">unique value</span>
    $newTxn\ \triangleq\ [zxid\quad \mapsto IncZxid(i, LastProposed(i).zxid),$
          $value\ \mapsto request\_value,$
          $ackSid \mapsto \{i\},$ <span style="background:#ccc">assume we have done $leader.processAck$</span>
          $epoch\ \ \mapsto acceptedEpoch[i]]$
    $m\_proposal\ \triangleq\ [mtype \mapsto PROPOSAL,$
           $mzxid \mapsto newTxn.zxid,$
           $mdata \mapsto request\_value]$
    $m\_proposal\_for\_checking\ \triangleq\ [source \mapsto i,$
                $epoch\ \ \mapsto acceptedEpoch[i],$
                $zxid\quad \mapsto newTxn.zxid,$
                $data\quad \mapsto request\_value]$
  $\text{IN}\quad \wedge\ history' = [history\ \text{EXCEPT } ![i] = Append(@, newTxn)]$
    $\wedge\ Broadcast(i, m\_proposal)$
    $\wedge\ proposalMsgsLog' = proposalMsgsLog \cup \{m\_proposal\_for\_checking\}$
  $\wedge\ \text{UNCHANGED } \langle state, currentEpoch, lastProcessed, zabState, acceptedEpoch,$
    $lastCommitted, electionVars, leaderVars, followerVars, initialHistory,$
    $epochLeader, violatedInvariants, electionMsgs\rangle$
  $\wedge\ UpdateRecorder(\langle\text{"LeaderProcessRequest"}, i\rangle)$

<span style="background:#ccc">Follower processes $PROPOSAL$ in $BROADCAST$. See $processPacket$ in Follower for details.</span>
$FollowerProcessPROPOSAL(i, j)\ \triangleq$
  $\wedge\ IsFollower(i)$
  $\wedge\ PendingPROPOSAL(i, j)$
  $\wedge\ zabState[i] = BROADCAST$
  $\wedge\ \text{LET } msg\ \triangleq\ msgs[j][i][1]$
    $infoOk\quad \triangleq\ IsMyLeader(i, j)$
    $isNext\quad \triangleq\ IsNextZxid(LastQueued(i).zxid, msg.mzxid)$
    $newTxn\ \triangleq\ [zxid\quad \mapsto msg.mzxid,$
          $value\ \ \mapsto msg.mdata,$
          $ackSid \mapsto \{\},$
          $epoch\ \ \mapsto acceptedEpoch[i]]$
    $m\_ack\quad \triangleq\ [mtype \mapsto ACK,$
         $mzxid \mapsto msg.mzxid]$
  $\text{IN}\quad \wedge\ infoOk$
    $\wedge\ \vee\ \wedge\ isNext$
        $\wedge\ \text{UNCHANGED } violatedInvariants$
      $\vee\ \wedge\ \neg isNext$
        $\wedge\ PrintT(\text{"Exception: Follower receives PROPOSAL, while"}\ \circ$
         $\text{" the transaction is not the next."})$
        $\wedge\ violatedInvariants' = [violatedInvariants\ \text{EXCEPT}$

$$!.proposalInconsistent = \text{TRUE}]$$
$$\land\ history' = [history\ \text{EXCEPT}\ ![i] = Append(@,\ newTxn)]$$
$$\land\ Reply(i,\ j,\ m\_ack)$$
$$\land\ \text{UNCHANGED}\ \langle state,\ currentEpoch,\ lastProcessed,\ zabState,\ acceptedEpoch,$$
$$lastCommitted,\ electionVars,\ leaderVars,\ followerVars,\ initialHistory,$$
$$epochLeader,\ proposalMsgsLog,\ electionMsgs\rangle$$
$$\land\ UpdateRecorder(\langle\text{``FollowerProcessPROPOSAL''},\ i,\ j\rangle)$$

See *outstandingProposals* in Leader
$$OutstandingProposals(i)\ \triangleq\ \text{IF}\ zabState[i] \neq BROADCAST\ \text{THEN}\ \langle\rangle$$
$$\text{ELSE}\ \ SubSeq(history[i],\ lastCommitted[i].index + 1,$$
$$Len(history[i]))$$

$$LastAckIndexFromFollower(i,\ j)\ \triangleq$$
$$\text{LET}\ set\_index\ \triangleq\ \{idx \in 1\mathinner{\ldotp\ldotp} Len(history[i]) : j \in history[i][idx].ackSid\}$$
$$\text{IN}\ \ \ Maximum(set\_index)$$

See commit in Leader for details.
$$LeaderCommit(s,\ follower,\ index,\ zxid)\ \triangleq$$
$$\land\ lastCommitted' = [lastCommitted\ \text{EXCEPT}\ ![s] = [index \mapsto index,$$
$$zxid\ \ \ \mapsto zxid]]$$
$$\land\ \text{LET}\ m\_commit\ \triangleq\ [mtype \mapsto COMMIT,$$
$$mzxid \mapsto zxid]$$
$$\text{IN}\ \ \ DiscardAndBroadcast(s,\ follower,\ m\_commit)$$

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.
$$LeaderTryToCommit(s,\ index,\ zxid,\ newTxn,\ follower)\ \triangleq$$
$$\text{LET}\ allTxnsBeforeCommitted\ \triangleq\ lastCommitted[s].index \geq index - 1$$
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.
$$\text{IN}\ \ \ \lor\ \land\ \text{Current conditions do not satisfy committing the proposal.}$$
$$\lor\ \neg allTxnsBeforeCommitted$$
$$\lor\ \neg hasAllQuorums$$
$$\land\ Discard(follower,\ s)$$
$$\land\ \text{UNCHANGED}\ \langle violatedInvariants,\ lastCommitted,\ lastProcessed\rangle$$
$$\lor\ \land\ allTxnsBeforeCommitted$$
$$\land\ hasAllQuorums$$

32

$\wedge \; \vee \; \wedge \neg ordered$
$\qquad \wedge PrintT(\text{``Warn: Committing zxid ''} \circ ToString(zxid) \circ \text{`` not first.''})$
$\qquad \wedge violatedInvariants' = [violatedInvariants \;\text{EXCEPT}$
$\qquad\qquad\qquad !.commitInconsistent = \text{TRUE}]$
$\quad\; \vee \; \wedge ordered$
$\qquad \wedge \text{UNCHANGED} \; violatedInvariants$
$\wedge LeaderCommit(s, follower, index, zxid)$
$\wedge lastProcessed' = [lastProcessed \;\text{EXCEPT}\; ![s] = [index \mapsto index,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad zxid \;\; \mapsto zxid]]$

<div style="background:#cccccc">Leader Keeps a count of acks for a particular proposal, and try to commit the proposal. See case <em>Leader.ACK</em> in <em>LearnerHandler</em>, <em>processRequest</em> in <em>AckRequestProcessor</em>, and <em>processAck</em> in Leader for details.</div>

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

33

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

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

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

$$
\begin{aligned}
&\quad\quad\quad\quad\quad ![i] = [index \mapsto lastCommitted[i].index + 1, \\
&\quad\quad\quad\quad\quad\quad\quad\quad zxid \;\;\mapsto firstElementZxid]] \\
&\quad\quad\quad \land \textsc{unchanged}\ violatedInvariants \\
&\quad \land Discard(j,\, i) \\
&\quad \land \textsc{unchanged}\ \langle state,\, currentEpoch,\, zabState,\, acceptedEpoch,\, history, \\
&\quad\quad\quad\quad\quad electionVars,\, leaderVars,\, initialHistory,\, followerVars, \\
&\quad\quad\quad\quad\quad proposalMsgsLog,\, epochLeader,\, electionMsgs\rangle \\
&\quad \land UpdateRecorder(\langle\text{``FollowerProcessCOMMIT''},\, i,\, j\rangle)
\end{aligned}
$$

---

<div style="background:#d3d3d3">Used to discard some messages which should not exist in network channel. This action should not be triggered.</div>

$FilterNonexistentMessage(i) \;\triangleq$
$\quad \land \exists j \in Server \setminus \{i\} : \;\land msgs[j][i] \neq \langle\rangle$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \land \textsc{let}\ msg \;\triangleq\; msgs[j][i][1]$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\ \ \textsc{in}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ \land IsLeader(i)$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\ \land \textsc{let}\ infoOk \;\triangleq\; IsMyLearner(i,\, j)$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \textsc{in}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = LEADERINFO$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = NEWLEADER$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = UPTODATE$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = PROPOSAL$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = COMMIT$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ \land \neg infoOk$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\ \land\ \lor\ msg.mtype = FOLLOWERINFO$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = ACKEPOCH$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = ACKLD$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = ACK$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ \land IsFollower(i)$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\ \land \textsc{let}\ infoOk \;\triangleq\; IsMyLeader(i,\, j)$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \textsc{in}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = FOLLOWERINFO$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = ACKEPOCH$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = ACKLD$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = ACK$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ \land \neg infoOk$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\ \land\ \lor\ msg.mtype = LEADERINFO$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = NEWLEADER$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = UPTODATE$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = PROPOSAL$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ msg.mtype = COMMIT$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \lor\ IsLooking(i)$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\ \land Discard(j,\, i)$
$\quad \land violatedInvariants' = [violatedInvariants\ \textsc{except}\ !.messageIllegal = \textsc{true}]$
$\quad \land \textsc{unchanged}\ \langle serverVars,\, electionVars,\, leaderVars,$

$$\langle followerVars,\ proposalMsgsLog,\ epochLeader,\ electionMsgs\rangle$$
$$\land\ UnchangeRecorder$$

---

Defines how the variables may transition.
$Next\ \triangleq$

        $FLE$ modlue
          $\lor\ \exists\,i,j\in Server : FLEReceiveNotmsg(i,j)$
          $\lor\ \exists\,i\in Server :\quad FLENotmsgTimeout(i)$
          $\lor\ \exists\,i\in Server :\quad FLEHandleNotmsg(i)$
          $\lor\ \exists\,i\in Server :\quad FLEWaitNewNotmsg(i)$
          $\lor\ \exists\,i\in Server :\quad FLEWaitNewNotmsgEnd(i)$
        Some conditions like failure, network delay
          $\lor\ \exists\,i\in Server :\quad FollowerTimeout(i)$
          $\lor\ \exists\,i\in Server :\quad LeaderTimeout(i)$
          $\lor\ \exists\,i,j\in Server : Timeout(i,j)$
        Zab module - Discovery and Synchronization part
          $\lor\ \exists\,i,j\in Server : ConnectAndFollowerSendFOLLOWERINFO(i,j)$
          $\lor\ \exists\,i,j\in Server : LeaderProcessFOLLOWERINFO(i,j)$
          $\lor\ \exists\,i,j\in Server : FollowerProcessLEADERINFO(i,j)$
          $\lor\ \exists\,i,j\in Server : LeaderProcessACKEPOCH(i,j)$
          $\lor\ \exists\,i\in Server :\quad LeaderSyncFollower(i)$
          $\lor\ \exists\,i,j\in Server : FollowerProcessSyncMessage(i,j)$
          $\lor\ \exists\,i,j\in Server : FollowerProcessPROPOSALInSync(i,j)$
          $\lor\ \exists\,i,j\in Server : FollowerProcessCOMMITInSync(i,j)$
          $\lor\ \exists\,i,j\in Server : FollowerProcessNEWLEADER(i,j)$
          $\lor\ \exists\,i,j\in Server : LeaderProcessACKLD(i,j)$
          $\lor\ \exists\,i,j\in Server : FollowerProcessUPTODATE(i,j)$
        Zab module $-\ Broadcast$ part
          $\lor\ \exists\,i\in Server :\quad LeaderProcessRequest(i)$
          $\lor\ \exists\,i,j\in Server : FollowerProcessPROPOSAL(i,j)$
          $\lor\ \exists\,i,j\in Server : LeaderProcessACK(i,j)$  $Sync+Broadcast$
          $\lor\ \exists\,i,j\in Server : FollowerProcessCOMMIT(i,j)$
        An action used to judge whether there are redundant messages in network
          $\lor\ \exists\,i\in Server :\quad FilterNonexistentMessage(i)$

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

---

Define safety properties of $Zab$ 1.0 protocol.

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

There is most one established leader for a certain epoch.
$Leadership1\ \triangleq\ \forall\,i,j\in Server :$
            $\land\ IsLeader(i)\land zabState[i]\in\{SYNCHRONIZATION,\ BROADCAST\}$
            $\land\ IsLeader(j)\land zabState[j]\in\{SYNCHRONIZATION,\ BROADCAST\}$
            $\land\ acceptedEpoch[i]=acceptedEpoch[j]$

$$\Rightarrow i = j$$

$Leadership2 \;\triangleq\; \forall\, epoch \in 1\mathrel{..} MAXEPOCH : Cardinality(epochLeader[epoch]) \leq 1$

$PrefixConsistency$: The prefix that have been committed
in history in any process is the same.
$PrefixConsistency \;\triangleq\; \forall\, i,\, j \in Server :$
$\qquad\qquad\qquad$ LET $smaller \;\triangleq\; Minimum(\{lastCommitted[i].index,\; lastCommitted[j].index\})$
$\qquad\qquad\qquad$ IN $\quad \vee\, smaller = 0$
$\qquad\qquad\qquad\qquad\quad \vee\, \wedge\, smaller > 0$
$\qquad\qquad\qquad\qquad\qquad \wedge\, \forall\, index \in 1\mathrel{..} smaller :$
$\qquad\qquad\qquad\qquad\qquad\qquad TxnEqual(history[i][index],\; history[j][index])$

Integrity: If some follower delivers one transaction, then some primary has broadcast it.
$Integrity \;\triangleq\; \forall\, i \in Server :$
$\qquad\qquad\quad \wedge\, IsFollower(i)$
$\qquad\qquad\quad \wedge\, lastCommitted[i].index > 0$
$\qquad\qquad\quad \Rightarrow \forall\, idx \in 1\mathrel{..} lastCommitted[i].index : \exists\, proposal \in proposalMsgsLog :$
$\qquad\qquad\qquad$ LET $txn\_proposal \;\triangleq\; [zxid \;\mapsto\; proposal.zxid,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad value \mapsto proposal.data]$
$\qquad\qquad\qquad$ IN $\quad TxnEqual(history[i][idx],\; txn\_proposal)$

Agreement: If some follower $f$ delivers transaction a and some follower $f'$ delivers transaction b,
$\qquad\quad$ then $f'$ delivers a or $f$ delivers b.
$Agreement \;\triangleq\; \forall\, i,\, j \in Server :$
$\qquad\qquad\quad \wedge\, IsFollower(i) \wedge lastCommitted[i].index > 0$
$\qquad\qquad\quad \wedge\, IsFollower(j) \wedge lastCommitted[j].index > 0$
$\qquad\qquad\quad \Rightarrow$
$\qquad\qquad\quad \forall\, idx1 \in 1\mathrel{..} lastCommitted[i].index,\; idx2 \in 1\mathrel{..} lastCommitted[j].index :$
$\qquad\qquad\qquad \vee\, \exists\, idx\_j \in 1\mathrel{..} lastCommitted[j].index :$
$\qquad\qquad\qquad\qquad TxnEqual(history[j][idx\_j],\; history[i][idx1])$
$\qquad\qquad\qquad \vee\, \exists\, idx\_i \in 1\mathrel{..} lastCommitted[i].index :$
$\qquad\qquad\qquad\qquad TxnEqual(history[i][idx\_i],\; history[j][idx2])$

Total order: If some follower delivers a before b, then any process that delivers b
$\qquad\qquad$ must also deliver a and deliver a before b.
$TotalOrder \;\triangleq\; \forall\, i,\, j \in Server :$
$\qquad\qquad\qquad$ LET $committed1 \;\triangleq\; lastCommitted[i].index$
$\qquad\qquad\qquad\qquad committed2 \;\triangleq\; lastCommitted[j].index$
$\qquad\qquad\qquad$ IN $\quad committed1 \geq 2 \wedge committed2 \geq 2$
$\qquad\qquad\qquad\qquad \Rightarrow \forall\, idx\_i1 \in 1\mathrel{..} (committed1 - 1) : \forall\, idx\_i2 \in (idx\_i1 + 1)\mathrel{..} committed1 :$
$\qquad\qquad\qquad\qquad$ LET $logOk \;\triangleq\; \exists\, idx \in 1\mathrel{..} committed2 :$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad TxnEqual(history[i][idx\_i2],\; history[j][idx])$
$\qquad\qquad\qquad\qquad$ IN $\quad \vee\, \neg logOk$
$\qquad\qquad\qquad\qquad\qquad\quad \vee\, \wedge\, logOk$
$\qquad\qquad\qquad\qquad\qquad\qquad\quad \wedge\, \exists\, idx\_j2 \in 1\mathrel{..} committed2 :$

$$\wedge\, TxnEqual(history[i][idx\_i2],\ history[j][idx\_j2])$$
$$\wedge\, \exists\, idx\_j1 \in 1\,..\,(idx\_j2 - 1):$$
$$TxnEqual(history[i][idx\_i1],\ history[j][idx\_j1])$$

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

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.

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

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.

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

$$TxnEqual(history[i][idx\_i],\ history[j][idx\_j])$$

\ * Modification History
\ * Last modified *Mon Nov* 22 21:49:29 *CST* 2021 by Dell
\ * Created Sat *Oct* 23 16:05:04 *CST* 2021 by Dell