

---

MODULE *ZabWithFLEAndSYNC*

---

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

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

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

EXTENDS *FastLeaderElection*

---

The set of requests that can go into history

CONSTANT *Value* \ \* Replaced by *recorder.nClientRequest*

*Value*  $\triangleq$  *Nat*

Zab states

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

Sync modes & message types

CONSTANTS *DIFF*, *TRUNC*

Message types

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

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

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

CONSTANT *Parameters*

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

*MAXEPOCH*  $\triangleq$  10

---

Variables in annotations mean variables defined in *FastLeaderElection*.

Variables that all servers use.

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

containing: *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 record [*sid*, connected].  
 Set of learners leader has received  
*ACKEPOCH* from, namely 'electingFollowers'  
 in Leader. Set of record  
 [*sid*, *peerLastZxid*, *inQuorum*].  
 And *peerLastZxid* =  $\langle -1, -1 \rangle$  means has done  
*syncFollower* with this *sid*.  
*inQuorum* = TRUE means in code it is one  
 element in 'electingFollowers'.

*ackldRecv*, Set of learners leader has received  
*ACK* of *NEWLEADER* from, namely  
 'newLeaderProposal' in Leader.

*forwarding*, Set of record [*sid*, connected].  
 Set of learners that are synced with  
 leader, namely 'forwardingFollowers'  
 in Leader.

*tempMaxEpoch* (*Maximum epoch in FOLLOWERINFO* + 1) that  
 leader has received from learners,  
 namely 'epoch' in Leader.

*leadingVoteSet* \ \* Set of voters that follow leader.

Variables only used for follower.

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

*packetsSync* packets of *PROPOSAL* and *COMMIT* from leader,  
 namely 'packetsNotCommitted' and  
 'packetsCommitted' in *SyncWithLeader*  
 in Learner.

Variables about network channel.

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

*electionMsgs* \ \* Network channel in *FLE* module.

Variables only used in verifying properties.

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

Variables only used for looking.

VARIABLE *currentVote*, \ \* Info of current vote, namely 'currentVote'  
           \ \* in *QuorumPeer*.  
           *logicalClock*, \ \* Election instance, namely 'logicalClock'  
           \ \* in *FastLeaderElection*.  
           *receiveVotes*, \ \* Votes from current *FLE* round, namely  
           \ \* 'recvset' in *FastLeaderElection*.  
           *outOfElection*, \ \* Votes from previous and current *FLE* round,  
           \ \* namely 'outofelection' in *FastLeaderElection*.  
           *recvQueue*, \ \* Queue of received notifications or timeout  
           \ \* signals.  
           *waitNotmsg* \ \* Whether waiting for new *not*. See line 1050  
           \ \* in *FastLeaderElection* for details.

VARIABLE *idTable* \ \* For mapping *Server* to Integers,  
                           to compare ids between servers.  
 Update: we have transformed *idTable* from variable to function.

VARIABLE *clientReuquest* \ \* Start from 0, and increases monotonically  
                                   when *LeaderProcessRequest* performed. To  
                                   avoid existing two requests with same value.  
 Update: Remove it to *recorder.nClientRequest*.

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

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

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

$electionVars \triangleq electionVarsL$

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

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

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

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

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

---


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

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

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

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

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

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

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

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

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

$$ViolationSet \triangleq \{"stateInconsistent", "proposalInconsistent", "commitInconsistent", "ackInconsistent", "messageIllegal"\}$$

$$Connecting \triangleq [sid : Server, connected : BOOLEAN]$$

$$\begin{aligned}
\text{AckldRecv} &\triangleq \text{Connecting} \\
\text{Electing} &\triangleq [\text{sid} : \text{Server}, \\
&\quad \text{peerLastZxid} : \text{Zxid}, \\
&\quad \text{inQuorum} : \text{BOOLEAN} ] \\
\text{Vote} &\triangleq \\
&\quad [\text{proposedLeader} : \text{ServersIncNullPoint}, \\
&\quad \text{proposedZxid} : \text{Zxid}, \\
&\quad \text{proposedEpoch} : \text{Nat}] \\
\text{ElectionVote} &\triangleq \\
&\quad [\text{vote} : \text{Vote}, \text{round} : \text{Nat}, \text{state} : \text{ElectionState}, \text{version} : \text{Nat}] \\
\text{ElectionMsg} &\triangleq \\
&\quad [\text{mtype} : \{\text{NOTIFICATION}\}, \\
&\quad \text{msource} : \text{Server}, \\
&\quad \text{mstate} : \text{ElectionState}, \\
&\quad \text{mround} : \text{Nat}, \\
&\quad \text{mvote} : \text{Vote}] \cup \\
&\quad [\text{mtype} : \{\text{NONE}\}] \\
\text{TypeOK} &\triangleq \\
&\quad \wedge \text{zabState} \in [\text{Server} \rightarrow \text{ZabState}] \\
&\quad \wedge \text{acceptedEpoch} \in [\text{Server} \rightarrow \text{Nat}] \\
&\quad \wedge \text{lastCommitted} \in [\text{Server} \rightarrow \text{LastItem}] \\
&\quad \wedge \text{learners} \in [\text{Server} \rightarrow \text{SUBSET Server}] \\
&\quad \wedge \text{connecting} \in [\text{Server} \rightarrow \text{SUBSET Connecting}] \\
&\quad \wedge \text{electing} \in [\text{Server} \rightarrow \text{SUBSET Electing}] \\
&\quad \wedge \text{ackldRecv} \in [\text{Server} \rightarrow \text{SUBSET AckldRecv}] \\
&\quad \wedge \text{forwarding} \in [\text{Server} \rightarrow \text{SUBSET Server}] \\
&\quad \wedge \text{initialHistory} \in [\text{Server} \rightarrow \text{Seq}(\text{HistoryItem})] \\
&\quad \wedge \text{tempMaxEpoch} \in [\text{Server} \rightarrow \text{Nat}] \\
&\quad \wedge \text{leaderAddr} \in [\text{Server} \rightarrow \text{ServersIncNullPoint}] \\
&\quad \wedge \text{packetsSync} \in [\text{Server} \rightarrow \text{SyncPackets}] \\
&\quad \wedge \text{proposalMsgsLog} \in \text{SUBSET Proposal} \\
&\quad \wedge \text{epochLeader} \in [1 \dots \text{MAXEPOCH} \rightarrow \text{SUBSET Server}] \\
&\quad \wedge \text{violatedInvariants} \in [\text{ViolationSet} \rightarrow \text{BOOLEAN} ] \\
&\quad \wedge \text{msgs} \in [\text{Server} \rightarrow [\text{Server} \rightarrow \text{Seq}(\text{Message})]] \\
&\quad \text{Fast Leader Election} \\
&\quad \wedge \text{electionMsgs} \in [\text{Server} \rightarrow [\text{Server} \rightarrow \text{Seq}(\text{ElectionMsg})]] \\
&\quad \wedge \text{recvQueue} \in [\text{Server} \rightarrow \text{Seq}(\text{ElectionMsg})] \\
&\quad \wedge \text{leadingVoteSet} \in [\text{Server} \rightarrow \text{SUBSET Server}] \\
&\quad \wedge \text{receiveVotes} \in [\text{Server} \rightarrow [\text{Server} \rightarrow \text{ElectionVote}]] \\
&\quad \wedge \text{currentVote} \in [\text{Server} \rightarrow \text{Vote}] \\
&\quad \wedge \text{outOfElection} \in [\text{Server} \rightarrow [\text{Server} \rightarrow \text{ElectionVote}]]
\end{aligned}$$



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


---

Actions about network

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

$$\begin{aligned}
& \wedge \text{msgs}[j][i][1].\text{mtype} = \text{NEWLEADER} \\
\text{PendingACKLD}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \wedge \text{msgs}[j][i][1].\text{mtype} = \text{ACKLD} \\
\text{PendingUPTODATE}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \wedge \text{msgs}[j][i][1].\text{mtype} = \text{UPTODATE} \\
\text{PendingPROPOSAL}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \wedge \text{msgs}[j][i][1].\text{mtype} = \text{PROPOSAL} \\
\text{PendingACK}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \wedge \text{msgs}[j][i][1].\text{mtype} = \text{ACK} \\
\text{PendingCOMMIT}(i, j) & \triangleq \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \wedge \text{msgs}[j][i][1].\text{mtype} = \text{COMMIT} \\
\text{Add a message to } \text{msgs} - \text{ add a message } m \text{ to } \text{msgs}. \\
\text{Send}(i, j, m) & \triangleq \text{msgs}' = [\text{msgs} \text{ EXCEPT } ![i][j] = \text{Append}(\text{msgs}[i][j], m)] \\
\text{SendPackets}(i, j, \text{ms}) & \triangleq \text{msgs}' = [\text{msgs} \text{ EXCEPT } ![i][j] = \text{msgs}[i][j] \circ \text{ms}] \\
\text{DiscardAndSendPackets}(i, j, \text{ms}) & \triangleq \text{msgs}' = [\text{msgs} \text{ EXCEPT } ![j][i] = \text{Tail}(\text{msgs}[j][i]), \\
& \quad \quad \quad ![i][j] = \text{msgs}[i][j] \circ \text{ms}] \\
\text{Remove a message from } \text{msgs} - \text{ discard head of } \text{msgs}. \\
\text{Discard}(i, j) & \triangleq \text{msgs}' = \text{IF } \text{msgs}[i][j] \neq \langle \rangle \text{ THEN } [\text{msgs} \text{ EXCEPT } ![i][j] = \text{Tail}(\text{msgs}[i][j])] \\
& \quad \quad \quad \text{ELSE } \text{msgs} \\
\text{Leader broadcasts a message (PROPOSAL/COMMIT) to all other servers in forwardingFollowers.} \\
\text{Broadcast}(i, m) & \triangleq \text{msgs}' = [\text{msgs} \text{ EXCEPT } ![i] = [v \in \text{Server} \mapsto \text{IF } \wedge v \in \text{forwarding}[i] \\
& \quad \quad \quad \wedge v \neq i \\
& \quad \quad \quad \text{THEN } \text{Append}(\text{msgs}[i][v], m) \\
& \quad \quad \quad \text{ELSE } \text{msgs}[i][v]]] \\
\text{DiscardAndBroadcast}(i, j, m) & \triangleq \\
& \text{msgs}' = [\text{msgs} \text{ EXCEPT } ![j][i] = \text{Tail}(\text{msgs}[j][i]), \\
& \quad \quad \quad ![i] = [v \in \text{Server} \mapsto \text{IF } \wedge v \in \text{forwarding}[i] \\
& \quad \quad \quad \wedge v \neq i \\
& \quad \quad \quad \text{THEN } \text{Append}(\text{msgs}[i][v], m) \\
& \quad \quad \quad \text{ELSE } \text{msgs}[i][v]]] \\
\text{Leader broadcasts LEADERINFO to all other servers in connectingFollowers.} \\
\text{DiscardAndBroadcastLEADERINFO}(i, j, m) & \triangleq \\
& \text{LET } \text{new\_connecting\_quorum} \triangleq \{c \in \text{connecting}'[i] : c.\text{connected} = \text{TRUE}\} \\
& \quad \text{new\_sid\_connecting} \triangleq \{c.\text{sid} : c \in \text{new\_connecting\_quorum}\} \\
& \text{IN} \\
& \text{msgs}' = [\text{msgs} \text{ EXCEPT } ![j][i] = \text{Tail}(\text{msgs}[j][i]), \\
& \quad \quad \quad ![i] = [v \in \text{Server} \mapsto \text{IF } \wedge v \in \text{new\_sid\_connecting} \\
& \quad \quad \quad \wedge v \in \text{learners}[i] \\
& \quad \quad \quad \wedge v \neq i \\
& \quad \quad \quad \text{THEN } \text{Append}(\text{msgs}[i][v], m) \\
& \quad \quad \quad \text{ELSE } \text{msgs}[i][v]]] \\
\text{Leader broadcasts UPTODATE to all other servers in newLeaderProposal.} \\
\text{DiscardAndBroadcastUPTODATE}(i, j, m) & \triangleq \\
& \text{LET } \text{new\_ackldRecv\_quorum} \triangleq \{a \in \text{ackldRecv}'[i] : a.\text{connected} = \text{TRUE}\} \\
& \quad \text{new\_sid\_ackldRecv} \triangleq \{a.\text{sid} : a \in \text{new\_ackldRecv\_quorum}\}
\end{aligned}$$



IN

$$\begin{aligned}
\text{msgs}' &= [\text{msgs} \text{ EXCEPT } ![j][i] = \text{Tail}(\text{msgs}[j][i]), \\
&\quad ![i] = [v \in \text{Server} \mapsto \text{IF } \wedge v \in \text{new\_sid\_ackldRecv} \\
&\quad \wedge v \in \text{learners}[i] \\
&\quad \wedge v \neq i \\
&\quad \text{THEN } \text{Append}(\text{msgs}[i][v], m) \\
&\quad \text{ELSE } \text{msgs}[i][v]] \\
\text{Combination of } \text{Send} \text{ and } \text{Discard} &- \text{discard head of } \text{msgs}[j][i] \text{ and add } m \text{ into } \text{msgs}. \\
\text{Reply}(i, j, m) &\triangleq \text{msgs}' = [\text{msgs} \text{ EXCEPT } ![j][i] = \text{Tail}(\text{msgs}[j][i]), \\
&\quad ![i][j] = \text{Append}(\text{msgs}[i][j], m)]
\end{aligned}$$

Shuffle input buffer.

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


---

Define initial values for all variables

$$\begin{aligned}
\text{InitServerVars} &\triangleq \wedge \text{InitServerVarsL} \\
&\quad \wedge \text{zabState} = [s \in \text{Server} \mapsto \text{ELECTION}] \\
&\quad \wedge \text{acceptedEpoch} = [s \in \text{Server} \mapsto 0] \\
&\quad \wedge \text{lastCommitted} = [s \in \text{Server} \mapsto [\text{index} \mapsto 0, \\
&\quad \quad \quad \text{zxid} \mapsto \langle 0, 0 \rangle]] \\
&\quad \wedge \text{initialHistory} = [s \in \text{Server} \mapsto \langle \rangle] \\
\text{InitLeaderVars} &\triangleq \wedge \text{InitLeaderVarsL} \\
&\quad \wedge \text{learners} = [s \in \text{Server} \mapsto \{\}] \\
&\quad \wedge \text{connecting} = [s \in \text{Server} \mapsto \{\}] \\
&\quad \wedge \text{electing} = [s \in \text{Server} \mapsto \{\}] \\
&\quad \wedge \text{ackldRecv} = [s \in \text{Server} \mapsto \{\}] \\
&\quad \wedge \text{forwarding} = [s \in \text{Server} \mapsto \{\}] \\
&\quad \wedge \text{tempMaxEpoch} = [s \in \text{Server} \mapsto 0] \\
\text{InitElectionVars} &\triangleq \text{InitElectionVarsL} \\
\text{InitFollowerVars} &\triangleq \wedge \text{leaderAddr} = [s \in \text{Server} \mapsto \text{NullPoint}] \\
&\quad \wedge \text{packetsSync} = [s \in \text{Server} \mapsto \\
&\quad \quad \quad [\text{notCommitted} \mapsto \langle \rangle, \\
&\quad \quad \quad \text{committed} \mapsto \langle \rangle]] \\
\text{InitVerifyVars} &\triangleq \wedge \text{proposalMsgsLog} = \{\} \\
&\quad \wedge \text{epochLeader} = [i \in 1 \dots \text{MAXEPOCH} \mapsto \{\}] \\
&\quad \wedge \text{violatedInvariants} = [\text{stateInconsistent} \mapsto \text{FALSE}, \\
&\quad \quad \quad \text{proposalInconsistent} \mapsto \text{FALSE},
\end{aligned}$$

$$\begin{array}{ll}
\text{commitInconsistent} & \mapsto \text{FALSE}, \\
\text{ackInconsistent} & \mapsto \text{FALSE}, \\
\text{messageIllegal} & \mapsto \text{FALSE}
\end{array}$$

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

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

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

---


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

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

#### Fast Leader Election

$$\begin{aligned}
\text{FLEReceiveNotmsg}(i, j) &\triangleq \\
&\wedge \text{ReceiveNotmsg}(i, j) \\
&\wedge \text{UNCHANGED} \langle \text{zabState}, \text{acceptedEpoch}, \text{lastCommitted}, \text{learners}, \text{connecting}, \\
&\quad \text{initialHistory}, \text{electing}, \text{ackldRecv}, \text{forwarding}, \text{tempMaxEpoch},
\end{aligned}$$

$$\begin{aligned}
& \text{followerVars, verifyVars, msgs} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"FLEReceiveNotmsg", } i, j \rangle) \\
\text{FLENotmsgTimeout}(i) & \triangleq \\
& \wedge \text{NotmsgTimeout}(i) \\
& \wedge \text{UNCHANGED} \langle \text{zabState, acceptedEpoch, lastCommitted, learners, connecting,} \\
& \quad \text{initialHistory, electing, ackldRecv, forwarding, tempMaxEpoch,} \\
& \quad \text{followerVars, verifyVars, msgs} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"FLENotmsgTimeout", } i \rangle) \\
\text{FLEHandleNotmsg}(i) & \triangleq \\
& \wedge \text{HandleNotmsg}(i) \\
& \wedge \text{LET } \text{newState} \triangleq \text{state}'[i] \\
& \text{IN} \\
& \vee \wedge \text{newState} = \text{LEADING} \\
& \quad \wedge \text{ZabTurnToLeading}(i) \\
& \quad \wedge \text{UNCHANGED } \text{packetsSync} \\
& \vee \wedge \text{newState} = \text{FOLLOWING} \\
& \quad \wedge \text{ZabTurnToFollowing}(i) \\
& \quad \wedge \text{UNCHANGED} \langle \text{learners, connecting, electing, ackldRecv,} \\
& \quad \quad \text{forwarding, tempMaxEpoch} \rangle \\
& \vee \wedge \text{newState} = \text{LOOKING} \\
& \quad \wedge \text{UNCHANGED} \langle \text{zabState, learners, connecting, electing, ackldRecv,} \\
& \quad \quad \text{forwarding, tempMaxEpoch, packetsSync, initialHistory} \rangle \\
& \wedge \text{UNCHANGED} \langle \text{lastCommitted, acceptedEpoch, leaderAddr, verifyVars, msgs} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"FLEHandleNotmsg", } i \rangle)
\end{aligned}$$

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

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

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

$$\begin{aligned}
\text{FLEWaitNewNotmsgEnd}(i) & \triangleq \\
& \wedge \text{WaitNewNotmsgEnd}(i) \\
& \wedge \text{LET } \text{newState} \triangleq \text{state}'[i] \\
& \text{IN} \\
& \vee \wedge \text{newState} = \text{LEADING} \\
& \quad \wedge \text{ZabTurnToLeading}(i) \\
& \quad \wedge \text{UNCHANGED } \text{packetsSync} \\
& \vee \wedge \text{newState} = \text{FOLLOWING}
\end{aligned}$$

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


---


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

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

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

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

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

$\wedge \text{leaderAddr}' = [\text{leaderAddr} \text{ EXCEPT } ![i] = \text{NullPoint}]$   
 $\wedge \text{CleanInputBuffer}(i)$

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

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

$\text{RemoveConnectingOrAckldRecv}(\text{set}, \text{sid}) \triangleq$   
 $\text{LET } \text{sid\_set} \triangleq \{s.\text{sid} : s \in \text{set}\}$   
 $\text{IN } \text{IF } \text{sid} \notin \text{sid\_set} \text{ THEN } \text{set}$   
 $\text{ELSE } \text{LET } \text{info} \triangleq \text{CHOOSE } s \in \text{set} : s.\text{sid} = \text{sid}$   
 $\text{new\_info} \triangleq [\text{sid} \mapsto \text{sid},$   
 $\text{connected} \mapsto \text{FALSE}]$   
 $\text{IN } (\text{set} \setminus \{\text{info}\}) \cup \{\text{new\_info}\}$

See `removeLearnerHandler` for details.

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

---

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

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

Leader loses support from a quorum and turns to *LOOKING*.  
 $LeaderTimeout(i) \triangleq$   
 $\wedge CheckTimeout$  test restrictions of *timeout\_2*  
 $\wedge IsLeader(i)$   
 $\wedge \neg IsQuorum(learners[i])$   
 $\wedge LeaderShutdown(i)$   
 $\wedge UNCHANGED \langle acceptedEpoch, lastCommitted, connecting, electing, ackldRecv,$   
 $tempMaxEpoch, initialHistory, verifyVars, packetsSync \rangle$   
 $\wedge UpdateRecorder(\langle "LeaderTimeout", i \rangle)$

Timeout between leader and follower.  
 $Timeout(i, j) \triangleq$   
 $\wedge CheckTimeout$  test restrictions of *timeout\_3*  
 $\wedge IsLeader(i) \wedge IsMyLearner(i, j)$   
 $\wedge IsFollower(j) \wedge IsMyLeader(j, i)$   
The action of leader *i*.  
 $\wedge RemoveLearner(i, j)$   
The action of follower *j*.  
 $\wedge FollowerShutdown(j)$   
 $\wedge Clean(i, j)$   
 $\wedge UNCHANGED \langle acceptedEpoch, lastCommitted, tempMaxEpoch,$   
 $initialHistory, verifyVars, packetsSync \rangle$   
 $\wedge UpdateRecorder(\langle "Timeout", i, j \rangle)$

---

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

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

*waitingForNewEpoch* in Leader  
 $WaitingForNewEpoch(i, set) \triangleq (i \in set \wedge IsQuorum(set)) = FALSE$   
 $WaitingForNewEpochTurnToFalse(i, set) \triangleq i \in set$   
 $\wedge IsQuorum(set)$

There may exists some follower in connecting but shuts down and connects again. So when leader broadcasts *LEADERINFO*, the follower will receive *LEADERINFO*, and receive it again after

sending *FOLLOWERINFO*. So connected makes sure each follower will only receive *LEADERINFO* at most once before timeout.

```

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

```

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

```

LeaderProcessFOLLOWERINFO(i, j)  $\triangleq$ 
   $\wedge$  CheckEpoch test restrictions of max epoch
   $\wedge$  IsLeader(i)
   $\wedge$  PendingFOLLOWERINFO(i, j)
   $\wedge$  LET msg  $\triangleq$  msgs[j][i][1]
      infoOk  $\triangleq$  IsMyLearner(i, j)
      lastAcceptedEpoch  $\triangleq$  msg.mzxid[1]
      sid_connecting  $\triangleq$  {c.sid : c  $\in$  connecting[i]}
  IN
     $\wedge$  infoOk
     $\wedge$   $\vee$  1. has not broadcast LEADERINFO
         $\wedge$  WaitingForNewEpoch(i, sid_connecting)
         $\wedge$   $\vee$   $\wedge$  zabState[i] = DISCOVERY
             $\wedge$  UNCHANGED violatedInvariants
         $\vee$   $\wedge$  zabState[i]  $\neq$  DISCOVERY
             $\wedge$  PrintT("Exception: waitingFotNewEpoch true,"  $\circ$ 
                    " while zabState not DISCOVERY.")
             $\wedge$  violatedInvariants' = [violatedInvariants EXCEPT !.stateInconsistent = TRUE]
         $\wedge$  tempMaxEpoch' = [tempMaxEpoch EXCEPT ![i] = IF lastAcceptedEpoch  $\geq$  tempMaxEpoch[i]
                            THEN lastAcceptedEpoch + 1
                            ELSE @]
     $\wedge$  connecting' = [connecting EXCEPT ![i] = UpdateConnectingOrAckldRecv(@, j)]
     $\wedge$  LET new_sid_connecting  $\triangleq$  {c.sid : c  $\in$  connecting'[i]}
  IN
     $\vee$   $\wedge$  WaitingForNewEpochTurnToFalse(i, new_sid_connecting)
         $\wedge$  acceptedEpoch' = [acceptedEpoch EXCEPT ![i] = tempMaxEpoch'[i]]
         $\wedge$  LET newLeaderZxid  $\triangleq$  <acceptedEpoch'[i], 0>
            m  $\triangleq$  [mtype  $\mapsto$  LEADERINFO,
                  mzxid  $\mapsto$  newLeaderZxid]

```

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

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

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



---

```

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

```

---

```

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

```

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

```

UpdateAckSid(his, lastSeenIndex, target) ≜
  IF Len(his) = 0 ∨ lastSeenIndex = 0 THEN his
  ELSE UpdateAckSidHelper(his, 1, Minimum({Len(his), lastSeenIndex}), target)

```

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

```

RECURSIVE ZxidToIndexHepler(−, −, −, −)
ZxidToIndexHepler(his, zxid, cur, appeared) ≜
  IF cur > Len(his) THEN cur
  ELSE IF TxnZxidEqual(his[cur], zxid)
    THEN CASE appeared = TRUE → − 1
           □   OTHER      → Minimum({cur,
                                     ZxidToIndexHepler(his, zxid, cur + 1, TRUE)})
    ELSE ZxidToIndexHepler(his, zxid, cur + 1, appeared)

```

```

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

```

Find index  $idx$  which meets:

$\text{history}[idx].\text{zxid} \leq \text{zxid} < \text{history}[idx + 1].\text{zxid}$

```

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

IndexOfZxid(his, zxid)  $\triangleq$  IF Len(his) = 0 THEN 0
                        ELSE LET idx  $\triangleq$  ZxidToIndex(his, zxid)
                              len  $\triangleq$  Len(his)
                              IN
                              IF idx  $\leq$  len THEN idx
                              ELSE IndexOfZxidHelper(his, zxid, 1, len)

```

```

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

```

```

RECURSIVE setPacketsForChecking(-, -, -, -, -, -)
setPacketsForChecking(set, src, ep, his, cur, end)  $\triangleq$ 
  IF cur > end THEN set
  ELSE LET m_proposal  $\triangleq$  [source  $\mapsto$  src,
                          epoch  $\mapsto$  ep,
                          zxid  $\mapsto$  his[cur].zxid,
                          data  $\mapsto$  his[cur].value]

```

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

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

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

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

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

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

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

we will not have case 4 where  $peerLastZxid < minCommittedLog$ , because  
 $minCommittedLog$  default value is 1 in our spec.

compare state summary of two servers

$$\begin{aligned}
IsMoreRecentThan(ss1, ss2) \triangleq & \vee ss1.currentEpoch > ss2.currentEpoch \\
& \vee \wedge ss1.currentEpoch = ss2.currentEpoch \\
& \wedge ZxidCompare(ss1.lastZxid, ss2.lastZxid)
\end{aligned}$$

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

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

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

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

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

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

$electing\_quorum \triangleq \{e \in electing[i] : e.inQuorum = \text{TRUE}\}$   
 $sid\_electing \triangleq \{s.sid : s \in electing\_quorum\}$   
IN  $\wedge infoOk$   
 $\wedge \vee$   $electionFinished = \text{true}$ , jump out of *waitForEpochAck*.  
Different from code, here we still need to record *info*  
into *electing*, to help us perform *syncFollower* afterwards.  
Since *electing* already meets quorum, it does not break  
consistency between code and spec.  
 $\wedge ElectionFinished(i, sid\_electing)$   
 $\wedge electing' = [electing \text{ EXCEPT } ![i] = UpdateElecting(@, j, msg.mzxid, \text{FALSE})]$   
 $\wedge Discard(j, i)$   
 $\wedge \text{UNCHANGED } \langle varsL, zabState, forwarding, leaderAddr,$   
 $\quad learners, epochLeader, violatedInvariants \rangle$   
 $\vee \wedge \neg ElectionFinished(i, sid\_electing)$   
 $\wedge \vee \wedge zabState[i] = DISCOVERY$   
 $\quad \wedge \text{UNCHANGED } violatedInvariants$   
 $\quad \vee \wedge zabState[i] \neq DISCOVERY$   
 $\quad \wedge PrintT(\text{"Exception: electionFinished false,"} \circ$   
 $\quad \quad \text{"while zabState not DISCOVERY."})$   
 $\quad \wedge violatedInvariants' = [violatedInvariants \text{ EXCEPT}$   
 $\quad \quad \quad !.stateInconsistent = \text{TRUE}]$   
 $\wedge \vee \wedge followerStateSummary.currentEpoch = -1$   
 $\quad \wedge electing' = [electing \text{ EXCEPT } ![i] = UpdateElecting(@, j,$   
 $\quad \quad \quad msg.mzxid, \text{FALSE})]$   
 $\quad \wedge Discard(j, i)$   
 $\quad \wedge \text{UNCHANGED } \langle varsL, zabState, forwarding, leaderAddr,$   
 $\quad \quad learners, epochLeader \rangle$   
 $\vee \wedge followerStateSummary.currentEpoch > -1$   
 $\quad \wedge \vee$  *normal follower*  
 $\quad \wedge logOk$   
 $\quad \wedge electing' = [electing \text{ EXCEPT } ![i] =$   
 $\quad \quad \quad UpdateElecting(@, j, msg.mzxid, \text{TRUE})]$   
 $\quad \wedge \text{LET } new\_electing\_quorum \triangleq \{e \in electing'[i] : e.inQuorum = \text{TRUE}\}$   
 $\quad \quad new\_sid\_electing \triangleq \{s.sid : s \in new\_electing\_quorum\}$   
IN  
 $\vee$   $electionFinished = \text{true}$ , jump out of *waitForEpochAck*,  
update *currentEpoch* and *zabState*.  
 $\wedge ElectionFinished(i, new\_sid\_electing)$   
 $\wedge LeaderTurnToSynchronization(i)$   
 $\wedge \text{LET } newLeaderEpoch \triangleq acceptedEpoch[i]$   
IN  $epochLeader' = [epochLeader \text{ EXCEPT } ![newLeaderEpoch]$   
 $\quad = @ \cup \{i\}]$  for checking invariants  
 $\vee$  there still exists  $electionFinished = \text{false}$ .  
 $\wedge \neg ElectionFinished(i, new\_sid\_electing)$   
 $\wedge \text{UNCHANGED } \langle currentEpoch, zabState, epochLeader \rangle$

$\wedge \text{Discard}(j, i)$   
 $\wedge \text{UNCHANGED } \langle \text{state}, \text{lastProcessed}, \text{electionVars}, \text{leadingVoteSet},$   
 $\text{electionMsgs}, \text{leaderAddr}, \text{learners}, \text{history}, \text{forwarding} \rangle$   
 $\vee$  Exists follower more recent than leader  
 $\wedge \neg \text{logOk}$   
 $\wedge \text{LeaderShutdown}(i)$   
 $\wedge \text{UNCHANGED } \langle \text{electing}, \text{epochLeader} \rangle$   
 $\wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{lastCommitted}, \text{connecting}, \text{ackldRecv},$   
 $\text{tempMaxEpoch}, \text{initialHistory}, \text{packetsSync}, \text{proposalMsgsLog} \rangle$   
 $\wedge \text{UpdateRecorder}(\langle \text{"LeaderProcessACKEPOCH"}, i, j \rangle)$

Strip *syncFollower* from *LeaderProcessACKEPOCH*.

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

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

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

$\text{FollowerProcessSyncMessage}(i, j) \triangleq$   
 $\wedge \text{IsFollower}(i)$   
 $\wedge \text{msgs}[j][i] \neq \langle \rangle$

```

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

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

```



```

LET  $lastIndex \triangleq Len(history[i])$ 
    $entry \triangleq history[i][lastIndex]$ 
IN   $[index \mapsto lastIndex,$ 
      $zxid \mapsto entry.zxid]$ 

```

See *lastQueued* in *Learner* for details.

```

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

```

```

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

```

```

 $FollowerProcessPROPOSALInSync(i, j) \triangleq$ 
   $\wedge IsFollower(i)$ 
   $\wedge PendingPROPOSAL(i, j)$ 
   $\wedge zabState[i] = SYNCHRONIZATION$ 
   $\wedge LET msg \triangleq msgs[j][i][1]$ 
     $infoOk \triangleq IsMyLeader(i, j)$ 
     $isNext \triangleq IsNextZxid>LastQueued(i).zxid, msg.mzxid)$ 
     $newTrn \triangleq [zxid \mapsto msg.mzxid,$ 
                    $value \mapsto msg.mdata,$ 
                    $ackSid \mapsto \{\},$ 
                    $epoch \mapsto acceptedEpoch[i]]$ 
  IN   $\wedge infoOk$ 
       $\wedge \vee \wedge isNext$ 
         $\wedge packetsSync' = [packetsSync \text{ EXCEPT } ![i].notCommitted$ 
                            $= Append(packetsSync[i].notCommitted, newTrn)]$ 
       $\vee \wedge \neg isNext$ 
         $\wedge PrintT("Warn: Follower receives PROPOSAL," \circ$ 
                    $" \text{ while } zxid \neq lastQueued + 1."$ )
         $\wedge UNCHANGED packetsSync$ 

```

$logRequest \rightarrow SyncRequestProcessor \rightarrow SendAckRequestProcessor \rightarrow \text{reply ack}$

So here we do not need to send ack to leader.

```

 $\wedge \text{Discard}(j, i)$ 
 $\wedge \text{UNCHANGED } \langle \text{serverVars}, \text{electionVars}, \text{leaderVars}, \text{leaderAddr},$ 
 $\text{verifyVars}, \text{electionMsgs} \rangle$ 
 $\wedge \text{UpdateRecorder}(\langle \text{"FollowerProcessPROPOSALInSync"}, i, j \rangle)$ 

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

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

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

```

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

```

 $\text{FollowerProcessCOMMITInSync}(i, j) \triangleq$ 
 $\wedge \text{IsFollower}(i)$ 
 $\wedge \text{PendingCOMMIT}(i, j)$ 

```

```

 $\wedge zabState[i] = SYNCHRONIZATION$ 
 $\wedge \text{LET } msg \triangleq msgs[j][i][1]$ 
 $\quad infoOk \triangleq IsMyLeader(i, j)$ 
 $\quad committedIndex \triangleq LastCommitted(i).index + 1$ 
 $\quad exist \triangleq \wedge committedIndex \leq LastQueued(i).index$ 
 $\quad \quad \wedge IsNextZxid(LastCommitted(i).zxid, msg.mzxid)$ 
 $\quad match \triangleq ZxidEqual(msg.mzxid, TrnWithIndex(i, committedIndex).zxid)$ 
IN  $\wedge infoOk$ 
 $\quad \wedge \vee \wedge exist$ 
 $\quad \quad \wedge \vee \wedge match$ 
 $\quad \quad \quad \wedge packetsSync' = [packetsSync \text{ EXCEPT } ![i].committed$ 
 $\quad \quad \quad \quad = Append(packetsSync[i].committed, msg.mzxid)]$ 
 $\quad \quad \quad \wedge \text{UNCHANGED } violatedInvariants$ 
 $\quad \quad \vee \wedge \neg match$ 
 $\quad \quad \quad \wedge PrintT(\text{"Warn: Follower receives COMMIT,"} \circ$ 
 $\quad \quad \quad \quad \text{" but zxid not the next committed zxid in COMMIT."})$ 
 $\quad \quad \quad \wedge violatedInvariants' = [violatedInvariants \text{ EXCEPT}$ 
 $\quad \quad \quad \quad \quad \quad !.commitInconsistent = TRUE]$ 
 $\quad \quad \quad \wedge \text{UNCHANGED } packetsSync$ 
 $\quad \vee \wedge \neg exist$ 
 $\quad \quad \wedge PrintT(\text{"Warn: Follower receives COMMIT,"} \circ$ 
 $\quad \quad \quad \quad \text{" but no packets with its zxid exists."})$ 
 $\quad \quad \quad \wedge violatedInvariants' = [violatedInvariants \text{ EXCEPT}$ 
 $\quad \quad \quad \quad \quad \quad !.commitInconsistent = TRUE]$ 
 $\quad \quad \quad \wedge \text{UNCHANGED } packetsSync$ 
 $\wedge Discard(j, i)$ 
 $\wedge \text{UNCHANGED } \langle serverVars, electionVars, leaderVars,$ 
 $\quad \quad leaderAddr, epochLeader, proposalMsgsLog, electionMsgs \rangle$ 
 $\wedge UpdateRecorder(\langle \text{"FollowerProcessCOMMITInSync"}, i, j \rangle)$ 

RECURSIVE  $ACKInBatches(-, -)$ 
 $ACKInBatches(queue, packets) \triangleq$ 
IF  $packets = \langle \rangle$  THEN  $queue$ 
ELSE LET  $head \triangleq packets[1]$ 
 $\quad newPackets \triangleq Tail(packets)$ 
 $\quad m\_ack \triangleq [mtype \mapsto ACK,$ 
 $\quad \quad \quad mzxid \mapsto head.zxid]$ 
IN  $ACKInBatches(Append(queue, m\_ack), newPackets)$ 

```

Update *currentEpoch*, and *logRequest* every packets in *packetsNotCommitted* and clear it. As *syncProcessor* will be called in *logRequest*, we have to reply acks here.

```

FollowerProcessNEWLEADER( $i, j$ )  $\triangleq$ 
 $\wedge IsFollower(i)$ 
 $\wedge PendingNEWLEADER(i, j)$ 
 $\wedge \text{LET } msg \triangleq msgs[j][i][1]$ 

```





```

ELSE history[i]  $\circ$  packetsSync[i].notCommitted
end  $\triangleq$  Len(completeHis)
first  $\triangleq$  IndexOfFirstTxnWithEpoch(completeHis, acceptedEpoch[i], 1, end)
IN IF first > end THEN  $\langle \rangle$ 
ELSE SubSeq(completeHis, first, end) completeHis[first : end]

```

Txns received in current epoch but not committed.

See *pendingTxns* in *FollowerZooKeeper* for details.

```

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

```

```

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

```

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

```

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

```

```

FollowerLogRequestInBatches(i, leader, ms_ack, packetsNotCommitted)  $\triangleq$ 
   $\wedge$  history' = [history EXCEPT ![i] = @  $\circ$  packetsNotCommitted]
   $\wedge$  DiscardAndSendPackets(i, leader, ms_ack)

```

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

```

FollowerCommitInBatches(i)  $\triangleq$ 
  LET committedTxns  $\triangleq$  CommittedTxns(i)
      packetsCommitted  $\triangleq$  packetsSync[i].committed
      match  $\triangleq$  TxnsAndCommittedMatch(committedTxns, packetsCommitted)
IN

```



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

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

$\text{FollowerProcessPROPOSAL}(i, j) \triangleq$   
 $\wedge \text{IsFollower}(i)$   
 $\wedge \text{PendingPROPOSAL}(i, j)$   
 $\wedge \text{zabState}[i] = \text{BROADCAST}$   
 $\wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1]$   
 $\text{infoOk} \triangleq \text{IsMyLeader}(i, j)$   
 $\text{isNext} \triangleq \text{IsNextZxid}(\text{LastQueued}(i).\text{zxid}, \text{msg.mzxid})$   
 $\text{newTrn} \triangleq [\text{zxid} \mapsto \text{msg.mzxid},$   
 $\text{value} \mapsto \text{msg.mdata},$   
 $\text{ackSid} \mapsto \{i\},$   
 $\text{epoch} \mapsto \text{acceptedEpoch}[i]]$   
 $\text{m\_ack} \triangleq [\text{mtype} \mapsto \text{ACK},$   
 $\text{mzxid} \mapsto \text{msg.mzxid}]$   
IN  $\wedge \text{infoOk}$   
 $\wedge \vee \wedge \text{isNext}$   
 $\wedge \text{UNCHANGED } \text{violatedInvariants}$   
 $\vee \wedge \neg \text{isNext}$   
 $\wedge \text{PrintT}(\text{"Exception: Follower receives PROPOSAL, while" } \circ$   
 $\text{" the transaction is not the next."})$   
 $\wedge \text{violatedInvariants}' = [\text{violatedInvariants} \text{ EXCEPT }$   
 $!. \text{proposalInconsistent} = \text{TRUE}]$



$$\begin{aligned}
& \wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{Append}(@, \text{newTxn})] \\
& \wedge \text{Reply}(i, j, m\_ack) \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{lastProcessed}, \text{zabState}, \text{acceptedEpoch}, \\
& \quad \text{lastCommitted}, \text{electionVars}, \text{leaderVars}, \text{followerVars}, \text{initialHistory}, \\
& \quad \text{epochLeader}, \text{proposalMsgsLog}, \text{electionMsgs} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"FollowerProcessPROPOSAL"}, i, j \rangle)
\end{aligned}$$

See *outstandingProposals* in Leader

$$\begin{aligned}
\text{OutstandingProposals}(i) & \triangleq \text{IF } \text{zabState}[i] \neq \text{BROADCAST} \text{ THEN } \langle \rangle \\
& \quad \text{ELSE } \text{SubSeq}(\text{history}[i], \text{lastCommitted}[i].\text{index} + 1, \\
& \quad \quad \text{Len}(\text{history}[i]))
\end{aligned}$$

$$\begin{aligned}
\text{LastAckIndexFromFollower}(i, j) & \triangleq \\
& \text{LET } \text{set\_index} \triangleq \{idx \in 1 \dots \text{Len}(\text{history}[i]) : j \in \text{history}[i][idx].\text{ackSid}\} \\
& \text{IN } \text{Maximum}(\text{set\_index})
\end{aligned}$$

See *commit* in Leader for details.

$$\begin{aligned}
\text{LeaderCommit}(s, \text{follower}, \text{index}, \text{zxid}) & \triangleq \\
& \wedge \text{lastCommitted}' = [\text{lastCommitted} \text{ EXCEPT } ![s] = [\text{index} \mapsto \text{index}, \\
& \quad \quad \quad \text{zxid} \mapsto \text{zxid}]] \\
& \wedge \text{LET } m\_commit \triangleq [\text{mtype} \mapsto \text{COMMIT}, \\
& \quad \quad \quad mzxid \mapsto \text{zxid}] \\
& \text{IN } \text{DiscardAndBroadcast}(s, \text{follower}, m\_commit)
\end{aligned}$$

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.

$$\begin{aligned}
\text{LeaderTryToCommit}(s, \text{index}, \text{zxid}, \text{newTxn}, \text{follower}) & \triangleq \\
& \text{LET } \text{allTxnsBeforeCommitted} \triangleq \text{lastCommitted}[s].\text{index} \geq \text{index} - 1 \\
& \quad \text{Only when all proposals before } \text{zxid} \text{ has been committed,} \\
& \quad \text{this proposal can be permitted to be committed.} \\
& \text{hasAllQuorums} \triangleq \text{IsQuorum}(\text{newTxn.ackSid}) \\
& \quad \text{In order to be committed, a proposal must be accepted} \\
& \quad \text{by a quorum.} \\
& \text{ordered} \triangleq \text{lastCommitted}[s].\text{index} + 1 = \text{index} \\
& \quad \text{Commit proposals in order.} \\
& \text{IN } \vee \wedge \text{Current conditions do not satisfy committing the proposal.} \\
& \quad \vee \neg \text{allTxnsBeforeCommitted} \\
& \quad \vee \neg \text{hasAllQuorums} \\
& \quad \wedge \text{Discard}(\text{follower}, s) \\
& \quad \wedge \text{UNCHANGED } \langle \text{violatedInvariants}, \text{lastCommitted}, \text{lastProcessed} \rangle \\
& \quad \vee \wedge \text{allTxnsBeforeCommitted} \\
& \quad \wedge \text{hasAllQuorums} \\
& \quad \wedge \vee \wedge \neg \text{ordered}
\end{aligned}$$

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

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

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



$$\begin{aligned}
& \text{zxid} \mapsto \text{firstElementZxid}] \\
& \wedge \text{UNCHANGED } \text{violatedInvariants} \\
& \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{zabState}, \text{acceptedEpoch}, \text{history}, \\
& \quad \text{electionVars}, \text{leaderVars}, \text{initialHistory}, \text{followerVars}, \\
& \quad \text{proposalMsgsLog}, \text{epochLeader}, \text{electionMsgs} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"FollowerProcessCOMMIT"}, i, j \rangle)
\end{aligned}$$

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

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

$\wedge \text{UnchangeRecorder}$

Defines how the variables may transition.

$\text{Next} \triangleq$

**FLE module**

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

**Some conditions like failure, network delay**

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

**Zab module - Discovery and Synchronization part**

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

**Zab module - Broadcast part**

$\vee \exists i \in \text{Server} : \text{LeaderProcessRequest}(i)$   
 $\vee \exists i, j \in \text{Server} : \text{FollowerProcessPROPOSAL}(i, j)$   
 $\vee \exists i, j \in \text{Server} : \text{LeaderProcessACK}(i, j) \quad \text{Sync} + \text{Broadcast}$   
 $\vee \exists i, j \in \text{Server} : \text{FollowerProcessCOMMIT}(i, j)$

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

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

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

Define safety properties of Zab 1.0 protocol.

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

There is most one established leader for a certain epoch.

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

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

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

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

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

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

$Integrity \triangleq \forall i \in Server :$   
 $\quad \wedge IsFollower(i)$   
 $\quad \wedge lastCommitted[i].index > 0$   
 $\quad \Rightarrow \forall idx \in 1 \dots lastCommitted[i].index : \exists proposal \in proposalMsgsLog :$   
 $\quad \quad LET \ txn\_proposal \triangleq [zxid \mapsto proposal.zxid,$   
 $\quad \quad \quad value \mapsto proposal.data]$   
 $\quad \quad IN \quad TsnEqual(history[i][idx], txn\_proposal)$

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

$Agreement \triangleq \forall i, j \in Server :$   
 $\quad \wedge IsFollower(i) \wedge lastCommitted[i].index > 0$   
 $\quad \wedge IsFollower(j) \wedge lastCommitted[j].index > 0$   
 $\quad \Rightarrow$   
 $\quad \forall idx1 \in 1 \dots lastCommitted[i].index, idx2 \in 1 \dots lastCommitted[j].index :$   
 $\quad \quad \vee \exists idx\_j \in 1 \dots lastCommitted[j].index :$   
 $\quad \quad \quad TsnEqual(history[j][idx\_j], history[i][idx1])$   
 $\quad \quad \vee \exists idx\_i \in 1 \dots lastCommitted[i].index :$   
 $\quad \quad \quad TsnEqual(history[i][idx\_i], history[j][idx2])$

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

$TotalOrder \triangleq \forall i, j \in Server :$   
 $\quad LET \ committed1 \triangleq lastCommitted[i].index$   
 $\quad \quad committed2 \triangleq lastCommitted[j].index$   
 $\quad IN \quad committed1 \geq 2 \wedge committed2 \geq 2$   
 $\quad \quad \Rightarrow \forall idx\_i1 \in 1 \dots (committed1 - 1) : \forall idx\_i2 \in (idx\_i1 + 1) \dots committed1 :$   
 $\quad \quad \quad LET \ logOk \triangleq \exists idx \in 1 \dots committed2 :$   
 $\quad \quad \quad \quad TsnEqual(history[i][idx\_i2], history[j][idx])$   
 $\quad \quad IN \quad \vee \neg logOk$   
 $\quad \quad \quad \vee \wedge logOk$   
 $\quad \quad \quad \wedge \exists idx\_j2 \in 1 \dots committed2 :$   
 $\quad \quad \quad \quad \wedge TsnEqual(history[i][idx\_i2], history[j][idx\_j2])$   
 $\quad \quad \quad \quad \wedge \exists idx\_j1 \in 1 \dots (idx\_j2 - 1) :$

$$TxnEqual(history[i][idx\_i1], history[j][idx\_j1])$$

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

$$LocalPrimaryOrder \triangleq LET \ p\_set(i, e) \triangleq \{p \in proposalMsgsLog : \wedge p.source = i \\ \wedge p.epoch = e\}$$

$$txn\_set(i, e) \triangleq \{[zxid \mapsto p.zxid, \\ value \mapsto p.data] : p \in p\_set(i, e)\}$$

$$IN \quad \forall i \in Server : \forall e \in 1 \dots currentEpoch[i] :$$

$$\quad \vee Cardinality(txn\_set(i, e)) < 2$$

$$\quad \vee \wedge Cardinality(txn\_set(i, e)) \geq 2$$

$$\quad \wedge \exists txn1, txn2 \in txn\_set(i, e) :$$

$$\quad \vee TxnEqual(txn1, txn2)$$

$$\quad \vee \wedge \neg TxnEqual(txn1, txn2)$$

$$\quad \wedge LET \ TxnPre \triangleq IF \ ZxidCompare(txn1.zxid, txn2.zxid) \ THEN \ txn2 \ ELSE$$

$$\quad \quad TxnNext \triangleq IF \ ZxidCompare(txn1.zxid, txn2.zxid) \ THEN \ txn1 \ ELSE$$

$$\quad IN \quad \forall j \in Server : \wedge lastCommitted[j].index \geq 2$$

$$\quad \quad \wedge \exists idx \in 1 \dots lastCommitted[j].index :$$

$$\quad \quad \quad TxnEqual(history[j][idx], TxnNext)$$

$$\quad \Rightarrow \exists idx2 \in 1 \dots lastCommitted[j].index :$$

$$\quad \quad \wedge TxnEqual(history[j][idx2], TxnNext)$$

$$\quad \quad \wedge idx2 > 1$$

$$\quad \quad \wedge \exists idx1 \in 1 \dots (idx2 - 1) :$$

$$\quad \quad \quad TxnEqual(history[j][idx1], TxnPre)$$

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

$$GlobalPrimaryOrder \triangleq \forall i \in Server : lastCommitted[i].index \geq 2$$

$$\Rightarrow \forall idx1, idx2 \in 1 \dots lastCommitted[i].index :$$

$$\quad \vee \neg EpochPrecedeInTxn(history[i][idx1], history[i][idx2])$$

$$\quad \vee \wedge EpochPrecedeInTxn(history[i][idx1], history[i][idx2])$$

$$\quad \wedge idx1 < idx2$$

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.

$$PrimaryIntegrity \triangleq \forall i, j \in Server : \wedge IsLeader(i) \quad \wedge IsMyLearner(i, j)$$

$$\quad \wedge IsFollower(j) \wedge IsMyLeader(j, i)$$

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

$$\quad \wedge zabState[j] = BROADCAST$$

$$\quad \wedge lastCommitted[j].index \geq 1$$

$$\Rightarrow \forall idx\_j \in 1 \dots lastCommitted[j].index :$$

$$\quad \vee history[j][idx\_j].zxid[1] \geq currentEpoch[i]$$

$$\quad \vee \wedge history[j][idx\_j].zxid[1] < currentEpoch[i]$$

$$\quad \wedge \exists idx\_i \in 1 \dots lastCommitted[i].index :$$

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

\\* Modification History  
\\* Last modified *Mon Nov 22 22:25:23 CST 2021* by Dell  
\\* Created Sat *Oct 23 16:05:04 CST 2021* by Dell