

---

MODULE *ZabWithFLE*

---

This is the formal specification for the *Zab* consensus algorithm, which means *Zookeeper Atomic Broadcast*.

Reference:

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

EXTENDS *FastLeaderElection*

---

The set of requests that can go into history

CONSTANT *Value*

Zab states

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

Message types

CONSTANTS *FOLLOWERINFO*, *LEADERINFO*, *ACKEPOCH*, *NEWLEADER*, *ACKLD*, *UPTODATE*, *PR*

NOTE: Additional message types used for recovery in *synchronization*(*TRUNC/DIFF/SNAP*) are not needed since we abstract this part.(see action *RECOVERYSYNC*)

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.

TODO: Consider in the future replacing the magic atomic synchronization *RECOVERYSYNC* with *DIFF* based message passing.

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

CONSTANT *Parameters*

TODO: Here we can add more constraints to decrease space.

*MAXEPOCH*  $\triangleq$  10

---

Variables that all servers need to use.

VARIABLES *zabState*, The current phase of server,  
in {*ELECTION*, *DISCOVERY*, *SYNCHRONIZATION*, *BROADCAST*}.

*acceptedEpoch*, The epoch number of the last *LEADERINFO* packet accepted,  
namely *f.p* in paper.

*history*, The history of servers as the sequence of transactions.

*commitIndex* Maximum index known to be committed.  
Starts from 0, and increases monotonically before restarting.  
Equals to 'lastCommitted' in code.

These transactions whose index \le *commitIndex*[*i*] can be applied to state machine immediately. So if we have a variable *applyIndex*, we can suppose that *applyIndex*[*i*] = *commitIndex*[*i*] when verifying properties. But in phase *SYNC*, follower will apply all queued proposals to state machine when receiving *NEWLEADER*. But follower only serves traffic after receiving *UPTODATE*, so sequential consistency is not violated.

So when we verify properties, we still suppose  $applyIndex[i] = commitIndex[i]$ , because this is an engineering detail.

Variables only used when  $state = LEADING$ .

VARIABLES <i>learners</i> ,	The set of servers which leader $i$ think are connected with $i$ .
<i>epochRecv</i> ,	The set of followers who has successfully sent <i>FOLLOWERINFO</i> to leader. Equals to 'connectingFollowers' in code.
<i>ackRecv</i> ,	The set of followers who has successfully sent <i>ACK-E</i> to leader. Equals to 'electingFollowers' in code.
<i>ackldRecv</i> ,	The set of followers who has successfully sent <i>ACK-LD</i> to leader in leader. Equals to 'newLeaderProposal' in code.
<i>forwarding</i> ,	The set of servers which leader $i$ should broadcast <i>PROPOSAL</i> and <i>COMMIT</i> to. Equals to 'forwardingFollowers' in code.
<i>ackIndex</i> ,	$[i][j]$ : The latest index that leader $i$ has received from follower $j$ via <i>ACK</i> .
<i>currentCounter</i> ,	$[i]$ : The count of transactions that clients have requested leader $i$ .
<i>sendCounter</i> ,	$\setminus * [i]$ : The count of transactions that leader $i$ has broadcast via <i>PROPOSAL</i> .
<i>committedIndex</i> ,	$\setminus * [i]$ : The maximum index of trasactions that leader $i$ has broadcast in <i>COMMIT</i> .
<i>committedCounter</i>	$[i][j]$ : The latest counter of transaction that leader $i$ has confirmed that follower $j$ has committed.

Variables only used when  $state = LEADING$  &  $zabState! = BROADCAST$ .

VARIABLES <i>initialHistory</i> ,	$[i]$ : The initial history of leader $i$ in epoch <i>acceptedEpoch</i> $[i]$ .
<i>tempMaxEpoch</i>	the maximum epoch in <i>CEPOCH</i> the prospective leader received from followers.

Variables only used when  $state = FOLLOWING$ .

VARIABLES <i>epochSent</i> ,	Express whether follower has sent <i>FOLLOWERINFO</i> to leader.
<i>leaderAddr</i> ,	Express whether follower $i$ has connected or lost connection. $[i]$ : The leader id of follower $i$ .
<i>synced</i>	Express whether follower has completed sync with leader.

Variables about network channel.

VARIABLE <i>msgs</i>	Simulates network channel. $msgs[i][j]$ means the input buffer of server $j$ from server $i$ .
----------------------	---

Variables only used in verifying properties.

VARIABLES <i>epochLeader</i> ,	The set of leaders in every epoch.
<i>proposalMsgsLog</i> ,	The set of all broadcast messages.

*inherentViolated*      Check whether there are conditions contrary to the facts.

Variable used for recording data to constrain state space.

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

$serverVarsZ \triangleq \langle state, currentEpoch, lastZxid, zabState, acceptedEpoch, history, commitIndex \rangle$       7 variables

$electionVarsZ \triangleq electionVars$       6 variables

$leaderVarsZ \triangleq \langle leadingVoteSet, learners, cepochRecv, ackeRecv, ackldRecv, forwarding, ackIndex, currentCounter, committedCounter \rangle$       9 variables

$tempVarsZ \triangleq \langle initialHistory, tempMaxEpoch \rangle$       2 variables

$followerVarsZ \triangleq \langle cepochSent, leaderAddr, synced \rangle$       3 variables

$verifyVarsZ \triangleq \langle proposalMsgsLog, epochLeader, inherentViolated \rangle$       3 variables

$msgVarsZ \triangleq \langle msgs, electionMsgs \rangle$       2 variables

$vars \triangleq \langle serverVarsZ, electionVarsZ, leaderVarsZ, tempVarsZ, followerVarsZ, verifyVarsZ, msgVarsZ, idTa$

---

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

$Zxid \triangleq$   
 $Seq(Nat)$   
 $\cup [epoch: Nat, counter: Nat]$

$HistoryItem \triangleq$   
 $[epoch: Nat,$   
 $counter: Nat,$   
 $value: Value]$

$Proposal \triangleq$   
 $[msource: Server, mtype: \{ "RECOVERYSYNC" \}, mepoch: Nat, mproposals: Seq(HistoryItem)] \cup$   
 $[msource: Server, mtype: \{ PROPOSAL \}, mepoch: Nat, mproposal: HistoryItem]$

$Message \triangleq$   
 $[mtype: \{ FOLLOWERINFO \}, mepoch: Nat] \cup$   
 $[mtype: \{ NEWLEADER \}, mepoch: Nat, mlastZxid: Zxid] \cup$   
 $[mtype: \{ ACKLD \}, mepoch: Nat] \cup$   
 $[mtype: \{ LEADERINFO \}, mepoch: Nat] \cup$   
 $[mtype: \{ ACKEPOCH \}, mepoch: Nat, mlastEpoch: Nat, mlastZxid: Zxid] \cup$   
 $[mtype: \{ UPTODATE \}, mepoch: Nat, mcommit: Nat] \cup$   
 $[mtype: \{ PROPOSAL \}, mepoch: Nat, mproposal: HistoryItem] \cup$   
 $[mtype: \{ ACK \}, mepoch: Nat, mzxid: Zxid] \cup$   
 $[mtype: \{ COMMIT \}, mepoch: Nat, mzxid: Zxid]$

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

$\text{Vote} \triangleq$   
 $[\text{proposedLeader} : \text{ServersIncNullPoint},$   
 $\text{proposedZxid} : \text{Zxid},$   
 $\text{proposedEpoch} : \text{Nat}]$

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

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

$\text{TypeOK} \triangleq$   
 $\wedge \text{zabState} \in [\text{Server} \rightarrow \{\text{ELECTION}, \text{DISCOVERY}, \text{SYNCHRONIZATION}, \text{BROADCAST}\}]$   
 $\wedge \text{acceptedEpoch} \in [\text{Server} \rightarrow \text{Nat}]$   
 $\wedge \text{history} \in [\text{Server} \rightarrow \text{Seq}(\text{HistoryItem})]$   
 $\wedge \text{commitIndex} \in [\text{Server} \rightarrow \text{Nat}]$   
 $\wedge \text{learners} \in [\text{Server} \rightarrow \text{SUBSET ServersIncNullPoint}]$   
 $\wedge \text{cepochnRecv} \in [\text{Server} \rightarrow \text{SUBSET ServersIncNullPoint}]$   
 $\wedge \text{ackRecv} \in [\text{Server} \rightarrow \text{SUBSET ServersIncNullPoint}]$   
 $\wedge \text{ackldRecv} \in [\text{Server} \rightarrow \text{SUBSET ServersIncNullPoint}]$   
 $\wedge \text{ackIndex} \in [\text{Server} \rightarrow [\text{Server} \rightarrow \text{Nat}]]$   
 $\wedge \text{currentCounter} \in [\text{Server} \rightarrow \text{Nat}]$   
 $\wedge \text{sendCounter} \in [\text{Server} \rightarrow \text{Nat}]$   
 $\wedge \text{committedIndex} \in [\text{Server} \rightarrow \text{Nat}]$   
 $\wedge \text{committedCounter} \in [\text{Server} \rightarrow [\text{Server} \rightarrow \text{Nat}]]$   
 $\wedge \text{forwarding} \in [\text{Server} \rightarrow \text{SUBSET ServersIncNullPoint}]$   
 $\wedge \text{initialHistory} \in [\text{Server} \rightarrow \text{Seq}(\text{HistoryItem})]$   
 $\wedge \text{tempMaxEpoch} \in [\text{Server} \rightarrow \text{Nat}]$   
 $\wedge \text{cepochnSent} \in [\text{Server} \rightarrow \text{BOOLEAN}]$   
 $\wedge \text{leaderAddr} \in [\text{Server} \rightarrow \text{ServersIncNullPoint}]$   
 $\wedge \text{syncd} \in [\text{Server} \rightarrow \text{BOOLEAN}]$   
 $\wedge \text{proposalMsgsLog} \in \text{SUBSET Proposal}$   
 $\wedge \text{epochLeader} \in [1 \dots \text{MAXEPOCH} \rightarrow \text{SUBSET ServersIncNullPoint}]$   
 $\wedge \text{inherentViolated} \in \text{BOOLEAN}$   
 $\wedge \text{forwarding} \in [\text{Server} \rightarrow \text{SUBSET Server}]$   
 $\wedge \text{msgs} \in [\text{Server} \rightarrow [\text{Server} \rightarrow \text{Seq}(\text{Message})]]$   
 $\text{Fast Leader Election}$   
 $\wedge \text{electionMsgs} \in [\text{Server} \rightarrow [\text{Server} \rightarrow \text{Seq}(\text{ElectionMsg})]]$   
 $\wedge \text{recvQueue} \in [\text{Server} \rightarrow \text{Seq}(\text{ElectionMsg})]$   
 $\wedge \text{leadingVoteSet} \in [\text{Server} \rightarrow \text{SUBSET Server}]$   
 $\wedge \text{receiveVotes} \in [\text{Server} \rightarrow [\text{Server} \rightarrow \text{ElectionVote}]]$   
 $\wedge \text{currentVote} \in [\text{Server} \rightarrow \text{Vote}]$   
 $\wedge \text{outOfElection} \in [\text{Server} \rightarrow [\text{Server} \rightarrow \text{ElectionVote}]]$   
 $\wedge \text{lastZxid} \in [\text{Server} \rightarrow \text{Zxid}]$   
 $\wedge \text{state} \in [\text{Server} \rightarrow \text{ElectionState}]$

---

$$Maximum(S) \stackrel{\Delta}{=} \text{IF } S = \{\} \text{ THEN } -1$$
$$\textit{Minimum}(S) \stackrel{\Delta}{=} \text{IF } S = \{\} \text{ THEN } -1$$

---

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

---

Actions about network

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

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

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

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

$Broadcast(i, m) \triangleq msgs' = [msgs \text{ EXCEPT } ![i] = [v \in Server \mapsto \text{IF } \wedge v \in forwarding[i] \wedge v \neq i \wedge \wedge m.mtype = PROPOSAL \wedge ackIndex[i][v] < Len(initialHistory[i]) + m.mproposal.counter \vee \wedge m.mtype = COMMIT \wedge committedCounter[i][v] < m.mzxid[2] \text{ THEN } Append(msgs[i][v], m) \text{ ELSE } msgs[i][v]]]$

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

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

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

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

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

---

Define initial values for all variables  
 $InitServerVarsZ \triangleq \wedge InitServerVars \wedge zabState = [s \in Server \mapsto ELECTION]$

$$\begin{aligned}
& \wedge \text{acceptedEpoch} = [s \in \text{Server} \mapsto 0] \\
& \wedge \text{history} = [s \in \text{Server} \mapsto \langle \rangle] \\
& \wedge \text{commitIndex} = [s \in \text{Server} \mapsto 0] \\
\text{InitLeaderVarsZ} & \triangleq \wedge \text{InitLeaderVars} \\
& \wedge \text{learners} = [s \in \text{Server} \mapsto \{\}] \\
& \wedge \text{cepochnRecv} = [s \in \text{Server} \mapsto \{\}] \\
& \wedge \text{ackRecv} = [s \in \text{Server} \mapsto \{\}] \\
& \wedge \text{ackldRecv} = [s \in \text{Server} \mapsto \{\}] \\
& \wedge \text{ackIndex} = [s \in \text{Server} \mapsto [v \in \text{Server} \mapsto 0]] \\
& \wedge \text{currentCounter} = [s \in \text{Server} \mapsto 0] \\
& \wedge \text{sendCounter} = [s \in \text{Server} \mapsto 0] \\
& \wedge \text{committedIndex} = [s \in \text{Server} \mapsto 0] \\
& \wedge \text{committedCounter} = [s \in \text{Server} \mapsto [v \in \text{Server} \mapsto 0]] \\
& \wedge \text{forwarding} = [s \in \text{Server} \mapsto \{\}] \\
\text{InitElectionVarsZ} & \triangleq \text{InitElectionVars} \\
\text{InitTempVarsZ} & \triangleq \wedge \text{initialHistory} = [s \in \text{Server} \mapsto \langle \rangle] \\
& \wedge \text{tempMaxEpoch} = [s \in \text{Server} \mapsto 0] \\
\text{InitFollowerVarsZ} & \triangleq \wedge \text{cepochnSent} = [s \in \text{Server} \mapsto \text{FALSE}] \\
& \wedge \text{leaderAddr} = [s \in \text{Server} \mapsto \text{NullPoint}] \\
& \wedge \text{synced} = [s \in \text{Server} \mapsto \text{FALSE}] \\
\text{InitVerifyVarsZ} & \triangleq \wedge \text{proposalMsgsLog} = \{\} \\
& \wedge \text{epochLeader} = [i \in 1 \dots \text{MAXEPOCH} \mapsto \{\}] \\
& \wedge \text{inherentViolated} = \text{FALSE} \\
\text{InitMsgVarsZ} & \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]] \\
\text{InitRecorder} & \triangleq \text{recorder} = [n\text{Timeout} \mapsto 0, \\
& \quad n\text{Transaction} \mapsto 0, \\
& \quad \text{maxEpoch} \mapsto 0, \\
& \quad pc \mapsto \langle \text{"InitZ"} \rangle] \\
\text{InitZ} & \triangleq \wedge \text{InitServerVarsZ} \\
& \wedge \text{InitLeaderVarsZ} \\
& \wedge \text{InitElectionVarsZ} \\
& \wedge \text{InitTempVarsZ} \\
& \wedge \text{InitFollowerVarsZ} \\
& \wedge \text{InitVerifyVarsZ} \\
& \wedge \text{InitMsgVarsZ} \\
& \wedge \text{idTable} = \text{InitializeIdTable}(\text{Server}) \\
& \wedge \text{InitRecorder}
\end{aligned}$$


---



$$\begin{aligned}
ZabTurnToLeading(i) &\triangleq \\
&\wedge zabState' = [zabState \text{ EXCEPT } ![i] = DISCOVERY] \\
&\wedge learners' = [learners \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge cepochRecv' = [cepochRecv \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge ackeRecv' = [ackeRecv \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge ackldRecv' = [ackldRecv \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge forwarding' = [forwarding \text{ EXCEPT } ![i] = \{\}] \\
&\wedge ackIndex' = [ackIndex \text{ EXCEPT } ![i] = [v \in Server \mapsto \text{IF } v = i \text{ THEN } Len(history[i]) \\
&\hspace{15em} \text{ELSE } 0]] \\
&\wedge currentCounter' = [currentCounter \text{ EXCEPT } ![i] = 0] \\
&\wedge commitIndex' = [commitIndex \text{ EXCEPT } ![i] = 0] \\
&\wedge sendCounter' = [sendCounter \text{ EXCEPT } ![i] = 0] \\
&\wedge committedIndex' = [committedIndex \text{ EXCEPT } ![i] = 0] \\
&\wedge committedCounter' = [committedCounter \text{ EXCEPT } ![i] = [v \in Server \mapsto \text{IF } v = i \text{ THEN } Len(history[i]) \\
&\hspace{15em} \text{ELSE } 0]] \\
&\wedge initialHistory' = [initialHistory \text{ EXCEPT } ![i] = history[i]] \\
&\wedge tempMaxEpoch' = [tempMaxEpoch \text{ EXCEPT } ![i] = acceptedEpoch[i]] \\
\\
ZabTurnToFollowing(i) &\triangleq \\
&\wedge zabState' = [zabState \text{ EXCEPT } ![i] = DISCOVERY] \\
&\wedge cepochSent' = [cepochSent \text{ EXCEPT } ![i] = FALSE] \\
&\wedge synced' = [synced \text{ EXCEPT } ![i] = FALSE] \\
&\wedge commitIndex' = [commitIndex \text{ EXCEPT } ![i] = 0] \\
\\
\text{Fast Leader Election} \\
FLEReceiveNotmsg(i, j) &\triangleq \\
&\wedge ReceiveNotmsg(i, j) \\
&\wedge \text{UNCHANGED } \langle zabState, acceptedEpoch, history, commitIndex, learners, cepochRecv, \\
&\hspace{10em} ackeRecv, ackldRecv, forwarding, ackIndex, currentCounter, \\
&\hspace{10em} committedCounter, tempVarsZ, followerVarsZ, verifyVarsZ, msgs \rangle \\
&\wedge UpdateRecorder(\langle \text{"FLEReceiveNotmsg"}, i, j \rangle) \\
\\
FLENotmsgTimeout(i) &\triangleq \\
&\wedge NotmsgTimeout(i) \\
&\wedge \text{UNCHANGED } \langle zabState, acceptedEpoch, history, commitIndex, learners, cepochRecv, ackeRecv, ackldRecv, \\
&\hspace{10em} forwarding, ackIndex, currentCounter, committedCounter, \\
&\hspace{10em} tempVarsZ, followerVarsZ, verifyVarsZ, msgs \rangle \\
&\wedge UpdateRecorder(\langle \text{"FLENotmsgTimeout"}, i \rangle) \\
\\
FLEHandleNotmsg(i) &\triangleq \\
&\wedge HandleNotmsg(i) \\
&\wedge \text{LET } newState \triangleq state'[i] \\
&\text{IN} \\
&\vee \wedge newState = LEADING \\
&\hspace{2em} \wedge ZabTurnToLeading(i)
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{UNCHANGED } \langle \text{cepocheSent}, \text{synced} \rangle \\
\vee & \wedge \text{newState} = \text{FOLLOWING} \\
& \wedge \text{ZabTurnToFollowing}(i) \\
& \wedge \text{UNCHANGED } \langle \text{learners}, \text{cepocheRecv}, \text{ackeRecv}, \text{ackldRecv}, \text{forwarding}, \text{ackIndex}, \text{currentCounter}, \\
& \quad \text{committedCounter}, \text{tempVarsZ} \rangle \\
\vee & \wedge \text{newState} = \text{LOOKING} \\
& \wedge \text{UNCHANGED } \langle \text{zabState}, \text{learners}, \text{cepocheRecv}, \text{ackeRecv}, \text{ackldRecv}, \text{forwarding}, \text{ackIndex}, \\
& \quad \text{currentCounter}, \text{commitIndex}, \\
& \quad \text{committedCounter}, \text{tempVarsZ}, \text{cepocheSent}, \text{synced} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{history}, \text{leaderAddr}, \text{verifyVarsZ}, \text{msgs} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"FLEHandleNotmsg"}, i \rangle)
\end{aligned}$$

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

$$\begin{aligned}
& \text{FLEWaitNewNotmsg}(i) \triangleq \\
& \wedge \text{WaitNewNotmsg}(i) \\
& \wedge \text{UNCHANGED } \langle \text{zabState}, \text{acceptedEpoch}, \text{history}, \text{commitIndex}, \text{learners}, \text{cepocheRecv}, \text{ackeRecv}, \\
& \quad \text{ackldRecv}, \text{forwarding}, \text{ackIndex}, \text{currentCounter}, \\
& \quad \text{committedCounter}, \text{tempVarsZ}, \text{followerVarsZ}, \text{verifyVarsZ}, \text{msgs} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"FLEWaitNewNotmsg"}, i \rangle)
\end{aligned}$$

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

$$\begin{aligned}
& \text{FLEWaitNewNotmsgEnd}(i) \triangleq \\
& \wedge \text{WaitNewNotmsgEnd}(i) \\
& \wedge \text{LET } \text{newState} \triangleq \text{state}'[i] \\
& \text{IN} \\
& \vee \wedge \text{newState} = \text{LEADING} \\
& \quad \wedge \text{ZabTurnToLeading}(i) \\
& \quad \wedge \text{UNCHANGED } \langle \text{cepocheSent}, \text{synced} \rangle \\
& \vee \wedge \text{newState} = \text{FOLLOWING} \\
& \quad \wedge \text{ZabTurnToFollowing}(i) \\
& \quad \wedge \text{UNCHANGED } \langle \text{learners}, \text{cepocheRecv}, \text{ackeRecv}, \text{ackldRecv}, \text{forwarding}, \text{ackIndex}, \text{currentCounter}, \\
& \quad \quad \text{committedCounter}, \text{tempVarsZ} \rangle \\
& \vee \wedge \text{newState} = \text{LOOKING} \\
& \quad \wedge \text{PrintT}(\text{"New state is LOOKING in FLEWaitNewNotmsgEnd, which should not happen."}) \\
& \quad \wedge \text{UNCHANGED } \langle \text{zabState}, \text{learners}, \text{cepocheRecv}, \text{ackeRecv}, \text{ackldRecv}, \text{forwarding}, \text{ackIndex}, \\
& \quad \quad \text{currentCounter}, \text{commitIndex}, \\
& \quad \quad \text{committedCounter}, \text{tempVarsZ}, \text{cepocheSent}, \text{synced} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{history}, \text{leaderAddr}, \text{verifyVarsZ}, \text{msgs} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"FLEWaitNewNotmsgEnd"}, i \rangle)
\end{aligned}$$


---

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

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

$$\begin{aligned}
\text{LeaderShutdown}(i) &\triangleq \\
&\wedge \text{ZabTimeout}(i) \\
&\wedge \text{zabState}' = [\text{zabState} \text{ EXCEPT } ![i] = \text{ELECTION}] \\
&\wedge \text{leaderAddr}' = [s \in \text{Server} \mapsto \text{IF } s \in \text{learners}[i] \text{ THEN } \text{NullPoint} \text{ ELSE } \text{leaderAddr}[s]] \\
&\wedge \text{learners}' = [\text{learners} \text{ EXCEPT } ![i] = \{\}] \\
&\wedge \text{forwarding}' = [\text{forwarding} \text{ EXCEPT } ![i] = \{\}] \\
&\wedge \text{msgs}' = [s \in \text{Server} \mapsto [v \in \text{Server} \mapsto \text{IF } v \in \text{learners}[i] \vee s \in \text{learners}[i] \\
&\quad \text{THEN } \langle \rangle \text{ ELSE } \text{msgs}[s][v]]]
\end{aligned}$$

$$\begin{aligned}
\text{RemoverLearner}(i, j) &\triangleq \\
&\wedge \text{learners}' = [\text{learners} \text{ EXCEPT } ![i] = \text{learners}[i] \setminus \{j\}] \\
&\wedge \text{forwarding}' = [\text{forwarding} \text{ EXCEPT } ![i] = \text{IF } j \in \text{forwarding}[i] \text{ THEN } \text{forwarding}[i] \setminus \{j\} \text{ ELSE } \text{forw}] \\
&\wedge \text{ceepochRecv}' = [\text{ceepochRecv} \text{ EXCEPT } ![i] = \text{IF } j \in \text{ceepochRecv}[i] \text{ THEN } \text{ceepochRecv}[i] \setminus \{j\} \text{ ELSE } \text{cepo}] \\
&\wedge \text{ackeRecv}' = [\text{ackeRecv} \text{ EXCEPT } ![i] = \text{IF } j \in \text{ackeRecv}[i] \text{ THEN } \text{ackeRecv}[i] \setminus \{j\} \text{ ELSE } \text{acke}] \\
&\wedge \text{ackldRecv}' = [\text{ackldRecv} \text{ EXCEPT } ![i] = \text{IF } j \in \text{ackldRecv}[i] \text{ THEN } \text{ackldRecv}[i] \setminus \{j\} \text{ ELSE } \text{ackld}] \\
&\wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i][j] = 0] \\
&\wedge \text{committedCounter}' = [\text{committedCounter} \text{ EXCEPT } ![i][j] = 0]
\end{aligned}$$

$$\begin{aligned}
\text{FollowerTimeout}(i) &\triangleq \\
&\wedge \text{IsFollower}(i) \\
&\wedge \text{HasNoLeader}(i) \\
&\wedge \text{FollowerShutdown}(i) \\
&\wedge \text{msgs}' = [s \in \text{Server} \mapsto [v \in \text{Server} \mapsto \text{IF } v = i \text{ THEN } \langle \rangle \text{ ELSE } \text{msgs}[s][v]]] \\
&\wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{history}, \text{commitIndex}, \text{learners}, \text{ceepochRecv}, \text{ackeRecv}, \text{ackldRecv}, \\
&\quad \text{forwarding}, \text{ackIndex}, \text{currentCounter}, \text{committedCounter}, \\
&\quad \text{tempVarsZ}, \text{ceepochSent}, \text{synced}, \text{verifyVarsZ} \rangle \\
&\wedge \text{UpdateRecorder}(\langle \text{"FollowerTimeout"}, i \rangle)
\end{aligned}$$

$$\begin{aligned}
\text{LeaderTimeout}(i) &\triangleq \\
&\wedge \text{IsLeader}(i) \\
&\wedge \neg \text{IsQuorum}(\text{learners}[i]) \\
&\wedge \text{LeaderShutdown}(i) \\
&\wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{history}, \text{commitIndex}, \text{ceepochRecv}, \text{ackeRecv}, \text{ackldRecv}, \text{ackIndex}, \text{currentCounter}, \\
&\quad \text{committedCounter}, \text{tempVarsZ}, \text{ceepochSent}, \text{synced}, \text{verifyVarsZ} \rangle \\
&\wedge \text{UpdateRecorder}(\langle \text{"LeaderTimeout"}, i \rangle)
\end{aligned}$$


---

Establish connection between leader  $i$  and follower  $j$ . It means  $i$  creates a *learnerHandler* for communicating with  $j$ , and  $j$  finds  $i$ 's address.

$$\begin{aligned}
\text{EstablishConnection}(i, j) &\triangleq \\
&\wedge \text{IsLeader}(i) \\
&\wedge \text{IsFollower}(j) \\
&\wedge \neg \text{IsMyLearner}(i, j) \\
&\wedge \text{HasNoLeader}(i) \\
&\wedge \text{currentVote}[j].\text{proposedLeader} = i \\
&\wedge \text{learners}' = [\text{learners} \text{ EXCEPT } ![i] = \text{learners}[i] \cup \{j\}] \quad \text{Leader: 'addLearnerHandler(peer)'} \\
&\wedge \text{leaderAddr}' = [\text{leaderAddr} \text{ EXCEPT } ![j] = i] \quad \text{Follower: 'connectToLeader(addr, hostname)'}
\end{aligned}$$

$\wedge$  UNCHANGED  $\langle serverVarsZ, electionVarsZ, leadingVoteSet, cepochRecv, ackeRecv, ackldRecv, forwardIndex, currentCounter, committedCounter, tempVarsZ, cepochSent, synced, verifyVarsZ, msgVarsZ, idTable \rangle$   
 $\wedge UpdateRecorder(\langle \text{"EstablishConnection"}, i, j \rangle)$

The leader  $i$  finds timeout and  $TCP$  connection between  $i$  and  $j$  closes.

$Timeout(i, j) \triangleq$   
 $\wedge IsLeader(i)$   
 $\wedge IsFollower(j)$   
 $\wedge IsMyLearner(i, j)$   
 $\wedge IsMyLeader(j, i)$   
 The action of leader  $i$ . (corresponding to function  $'removeLearnerHandler(peer)'$ .)  
 $\wedge RemoverLearner(i, j)$   
 The action of follower  $j$ .  
 $\wedge FollowerShutdown(j)$   
 Clean channel between  $i$  and  $j$ .  
 $\wedge Clean(i, j)$   
 $\wedge$  UNCHANGED  $\langle acceptedEpoch, history, commitIndex, currentCounter, tempVarsZ, cepochSent, synced, verifyVarsZ \rangle$   
 $\wedge UpdateRecorder(\langle \text{"Timeout"}, i, j \rangle)$

Follower sends  $f.p$  to leader via  $FOLLOWERINFO(CEPOCH)$ .

$FollowerSendFOLLOWERINFO(i) \triangleq$   
 $\wedge IsFollower(i)$   
 $\wedge zabState[i] = DISCOVERY$   
 $\wedge HasLeader(i)$   
 $\wedge \neg cepochSent[i]$   
 $\wedge Send(i, leaderAddr[i], [mtype \mapsto FOLLOWERINFO, mepoch \mapsto acceptedEpoch[i]])$   
 $\wedge cepochSent' = [cepochSent \text{ EXCEPT } ![i] = \text{TRUE}]$   
 $\wedge$  UNCHANGED  $\langle serverVarsZ, leaderVarsZ, electionVarsZ, tempVarsZ, leaderAddr, synced, verifyVarsZ, electionMsgs, idTable \rangle$   
 $\wedge UpdateRecorder(\langle \text{"FollowerSendFOLLOWERINFO"}, i \rangle)$

Leader waits for receiving  $FOLLOWERINFO$  from a quorum, and then chooses a new epoch  $e'$  as its own epoch and broadcasts  $LEADERINFO$ .

$LeaderHandleFOLLOWERINFO(i, j) \triangleq$   
 $\wedge IsLeader(i)$   
 $\wedge PendingFOLLOWERINFO(i, j)$   
 $\wedge \text{LET } msg \triangleq msgs[j][i][1]$   
 IN  $\vee$  1. has not broadcast  $LEADERINFO$  – modify  $tempMaxEpoch$   
 $\wedge NullPoint \notin cepochRecv[i]$   
 $\wedge \text{LET } newEpoch \triangleq Maximum(\{tempMaxEpoch[i], msg.mepoch\})$   
 IN  $tempMaxEpoch' = [tempMaxEpoch \text{ EXCEPT } ![i] = newEpoch]$   
 $\wedge Discard(j, i)$

$\vee$  2. has broadcast *LEADERINFO* – no need to handle the *msg*, just send *LEADERINFO* to corresponding  
 $\wedge \text{NullPoint} \in \text{cePOCHRecv}[i]$   
 $\wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{LEADERINFO},$   
 $\quad \text{mepoch} \mapsto \text{acceptedEpoch}[i]])$   
 $\wedge \text{UNCHANGED tempMaxEpoch}$   
 $\wedge \text{cePOCHRecv}' = [\text{cePOCHRecv} \text{ EXCEPT } ![i] = \text{IF } j \in \text{cePOCHRecv}[i] \text{ THEN } \text{cePOCHRecv}[i]$   
 $\quad \text{ELSE } \text{cePOCHRecv}[i] \cup \{j\}]$   
 $\wedge \text{UNCHANGED } \langle \text{serverVarsZ}, \text{followerVarsZ}, \text{electionVarsZ}, \text{initialHistory}, \text{leadingVoteSet}, \text{learners},$   
 $\quad \text{ackRecv}, \text{ackldRecv}, \text{forwarding}, \text{ackIndex}, \text{currentCounter},$   
 $\quad \text{committedCounter}, \text{verifyVarsZ}, \text{electionMsgs}, \text{idTable} \rangle$   
 $\wedge \text{UpdateRecorder}(\langle \text{"LeaderHandleFOLLOWERINFO"}, i, j \rangle)$

$\text{LeaderBroadcastLEADERINFO}(i) \triangleq$   
 $\wedge \text{IsLeader}(i)$   
 $\wedge \text{zabState}[i] = \text{DISCOVERY}$   
 $\wedge \text{IsQuorum}(\text{cePOCHRecv}[i])$   
 $\wedge \text{acceptedEpoch}' = [\text{acceptedEpoch} \text{ EXCEPT } ![i] = \text{tempMaxEpoch}[i] + 1]$   
 $\wedge \text{cePOCHRecv}' = [\text{cePOCHRecv} \text{ EXCEPT } ![i] = \text{cePOCHRecv}[i] \cup \{\text{NullPoint}\}]$   
 $\wedge \text{BroadcastLEADERINFO}(i, [\text{mtype} \mapsto \text{LEADERINFO},$   
 $\quad \text{mepoch} \mapsto \text{acceptedEpoch}'[i]])$   
 $\wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{lastZxid}, \text{zabState}, \text{history}, \text{commitIndex}, \text{electionVarsZ},$   
 $\quad \text{leadingVoteSet}, \text{learners}, \text{ackRecv}, \text{ackldRecv}, \text{forwarding}, \text{ackIndex}, \text{currentCounter},$   
 $\quad \text{committedCounter}, \text{tempVarsZ}, \text{followerVarsZ}, \text{verifyVarsZ},$   
 $\quad \text{electionMsgs}, \text{idTable} \rangle$   
 $\wedge \text{UpdateRecorder}(\langle \text{"LeaderBroadcastLEADERINFO"}, i \rangle)$

In phase *f12*, follower receives *NEWEPOCH*. If  $e' > f.p$ , then follower sends *ACK-E* back, and *ACK-E* contains *f.a* and *lastZxid* to let leader judge whether it is the latest. After handling *NEWEPOCH*, follower's *zabState* turns to *SYNCHRONIZATION*.

$\text{FollowerHandleLEADERINFO}(i, j) \triangleq$   
 $\wedge \text{IsFollower}(i)$   
 $\wedge \text{PendingLEADERINFO}(i, j)$   
 $\wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1]$   
 $\quad \text{infoOk} \triangleq \text{IsMyLeader}(i, j)$   
 $\quad \text{epochOk} \triangleq \wedge \text{infoOk}$   
 $\quad \quad \wedge \text{msg.mepoch} \geq \text{acceptedEpoch}[i]$   
 $\quad \text{correct} \triangleq \wedge \text{epochOk}$   
 $\quad \quad \wedge \text{zabState}[i] = \text{DISCOVERY}$

IN  $\wedge \text{infoOk}$   
 $\wedge \vee$  1. Normal case  
 $\wedge \text{epochOk}$   
 $\wedge \vee \wedge \text{correct}$   
 $\quad \wedge \text{acceptedEpoch}' = [\text{acceptedEpoch} \text{ EXCEPT } ![i] = \text{msg.mepoch}]$   
 $\quad \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{ACKEPOCH},$   
 $\quad \quad \text{mepoch} \mapsto \text{msg.mepoch},$   
 $\quad \quad \text{mlastEpoch} \mapsto \text{currentEpoch}[i],$

$$\begin{aligned}
& mlastZxid \mapsto lastZxid[i]] \\
& \wedge cepochSent' = [ceepochSent \text{ EXCEPT } ![i] = \text{TRUE}] \\
& \wedge \text{UNCHANGED } inherentViolated \\
\vee & \wedge \neg correct \\
& \wedge PrintT(\text{"Exception: Follower receives LEADERINFO while its ZabState is not DISCOVERY"}) \\
& \wedge inherentViolated' = \text{TRUE} \\
& \wedge Discard(j, i) \\
& \wedge \text{UNCHANGED } \langle acceptedEpoch, cepochSent \rangle \\
& \wedge zabState' = [zabState \text{ EXCEPT } ![i] = \text{IF } zabState[i] = \text{DISCOVERY} \text{ THEN } \text{SYNCHRONIZATION} \\
& \hspace{15em} \text{ELSE } zabState[i]] \\
& \wedge \text{UNCHANGED } \langle varsL, leaderAddr \rangle \\
\vee & \text{2. Abnormal case - go back to election} \\
& \wedge \neg epochOk \\
& \wedge FollowerShutdown(i) \\
& \wedge Clean(i, j) \\
& \wedge \text{UNCHANGED } \langle acceptedEpoch, cepochSent, inherentViolated \rangle \\
& \wedge \text{UNCHANGED } \langle history, commitIndex, learners, cepochRecv, ackRecv, ackldRecv, forwarding, ackIndex, \\
& \hspace{10em} currentCounter, committedCounter, tempVarsZ, synced, proposalMsgsLog, epochLead \rangle \\
& \wedge UpdateRecorder(\langle \text{"FollowerHandleLEADERINFO"}, i, j \rangle)
\end{aligned}$$

Abstraction of actions making follower *synced* with leader before leader sending *NEWLEADER*.

$$\begin{aligned}
SubRECOVERYSYNC(i, j) & \triangleq \\
& \text{LET } canSync \triangleq \wedge IsLeader(i) \wedge zabState[i] \neq \text{DISCOVERY} \wedge IsMyLearner(i, j) \\
& \wedge IsFollower(j) \wedge zabState[j] = \text{SYNCHRONIZATION} \wedge IsMyLeader(j, i) \\
& \wedge synced[j] = \text{FALSE} \\
& \text{IN} \\
& \vee \wedge canSync \\
& \wedge history' = [history \text{ EXCEPT } ![j] = history[i]] \\
& \wedge lastZxid' = [lastZxid \text{ EXCEPT } ![j] = lastZxid[i]] \\
& \wedge UpdateProposal(j, leaderAddr[j], lastZxid'[j], currentEpoch[j]) \\
& \wedge commitIndex' = [commitIndex \text{ EXCEPT } ![j] = commitIndex[i]] \\
& \wedge synced' = [synced \text{ EXCEPT } ![j] = \text{TRUE}] \\
& \wedge forwarding' = [forwarding \text{ EXCEPT } ![i] = forwarding[i] \cup \{j\}] \quad j \text{ will join traffic, and receive } PL \\
& \wedge ackIndex' = [ackIndex \text{ EXCEPT } ![i][j] = Len(history[i])] \\
& \wedge committedCounter' = [committedCounter \text{ EXCEPT } ![i][j] = Maximum(\{commitIndex[i] - Len(history[i])\})] \\
& \wedge \text{LET } ms \triangleq [msource \mapsto i, mtype \mapsto \text{"RECOVERYSYNC"}, mepoch \mapsto acceptedEpoch[i], mproposals \mapsto \\
& \hspace{10em} \text{IN } proposalMsgsLog' = \text{IF } ms \in proposalMsgsLog \text{ THEN } proposalMsgsLog \\
& \hspace{15em} \text{ELSE } proposalMsgsLog \cup \{ms\} \\
& \wedge \text{UNCHANGED } inherentViolated \\
& \vee \wedge \neg canSync \\
& \wedge PrintT(\text{"Exception: Leader wants to sync with follower while the condition doesn't allow."}) \\
& \wedge inherentViolated' = \text{TRUE} \\
& \wedge \text{UNCHANGED } \langle history, lastZxid, currentVote, commitIndex, synced, forwarding, ackIndex, committedCounter \rangle
\end{aligned}$$

Leader waits for receiving *ACKEPOCH* from a quorum, and check whether it has the latest history and epoch from them. If so, leader's *zabState* turns to *SYNCHRONIZATION*.

$$\begin{aligned}
& \text{LeaderHandleACKEPOCH}(i, j) \triangleq \\
& \quad \wedge \text{IsLeader}(i) \\
& \quad \wedge \text{PendingACKEPOCH}(i, j) \\
& \quad \wedge \text{LET } msg \triangleq msgs[j][i][1] \\
& \quad \quad infoOk \triangleq \wedge \text{IsMyLearner}(i, j) \\
& \quad \quad \quad \wedge \text{acceptedEpoch}[i] = msg.mepoch \\
& \quad \quad logOk \triangleq \text{logOk represents whether leader is more up-to-date than follower} \\
& \quad \quad \quad \wedge infoOk \\
& \quad \quad \quad \wedge \vee \text{currentEpoch}[i] > msg.mlastEpoch \\
& \quad \quad \quad \quad \vee \wedge \text{currentEpoch}[i] = msg.mlastEpoch \\
& \quad \quad \quad \quad \quad \wedge \vee \text{lastZxid}[i][1] > msg.mlastZxid[1] \\
& \quad \quad \quad \quad \quad \quad \vee \wedge \text{lastZxid}[i][1] = msg.mlastZxid[1] \\
& \quad \quad \quad \quad \quad \quad \quad \wedge \text{lastZxid}[i][2] \geq msg.mlastZxid[2] \\
& \quad \quad replyOk \triangleq \wedge infoOk \\
& \quad \quad \quad \wedge \text{NullPoint} \in \text{ackeRecv}[i] \\
& \text{IN} \quad \wedge infoOk \\
& \quad \wedge \vee \wedge \text{replyOk} \\
& \quad \quad \wedge \text{SubRECOVERYSYNC}(i, j) \\
& \quad \quad \wedge \text{ackeRecv}' = [\text{ackeRecv} \text{ EXCEPT } ![i] = \text{IF } j \notin \text{ackeRecv}[i] \text{ THEN } \text{ackeRecv}[i] \cup \{j\} \\
& \quad \quad \quad \quad \quad \quad \quad \text{ELSE } \text{ackeRecv}[i]] \\
& \quad \quad \wedge \text{Reply}(i, j, [mtype \mapsto \text{NEWLEADER}, \\
& \quad \quad \quad \quad mepoch \mapsto \text{acceptedEpoch}[i], \\
& \quad \quad \quad \quad mlastZxid \mapsto \text{lastZxid}[i]]) \\
& \quad \quad \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{logicalClock}, \text{receiveVotes}, \text{outOfElection}, \text{recvQueue}, \\
& \quad \quad \quad \text{leadingVoteSet}, \text{electionMsgs}, \text{idTable}, \text{zabState}, \text{leaderAddr}, \text{learners} \rangle \\
& \quad \vee \wedge \neg \text{replyOk} \\
& \quad \quad \wedge \vee \text{normal case} \\
& \quad \quad \quad \wedge logOk \\
& \quad \quad \quad \wedge \text{ackeRecv}' = [\text{ackeRecv} \text{ EXCEPT } ![i] = \text{IF } j \notin \text{ackeRecv}[i] \text{ THEN } \text{ackeRecv}[i] \cup \{j\} \\
& \quad \quad \quad \quad \quad \quad \quad \text{ELSE } \text{ackeRecv}[i]] \\
& \quad \quad \quad \wedge \text{Discard}(j, i) \\
& \quad \quad \quad \wedge \text{UNCHANGED } \langle \text{varsL}, \text{zabState}, \text{leaderAddr}, \text{learners}, \text{forwarding} \rangle \\
& \quad \quad \vee \text{go back to election since there exists follower more up-to-date than leader} \\
& \quad \quad \quad \wedge \neg logOk \\
& \quad \quad \quad \wedge \text{LeaderShutdown}(i) \\
& \quad \quad \quad \wedge \text{UNCHANGED } \text{ackeRecv} \\
& \quad \quad \wedge \text{UNCHANGED } \langle \text{history}, \text{commitIndex}, \text{syncd}, \text{forwarding}, \text{ackIndex}, \\
& \quad \quad \quad \text{committedCounter}, \text{proposalMsgsLog}, \text{inherentViolated} \rangle \\
& \quad \wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{ceepochRecv}, \text{ackldRecv}, \text{currentCounter}, \\
& \quad \quad \text{tempVarsZ}, \text{ceepochSent}, \text{epochLeader} \rangle \\
& \quad \wedge \text{UpdateRecorder}(\langle \text{"LeaderHandleACKEPOCH"}, i, j \rangle) \\
& \text{LeaderTransitionToSynchronization}(i) \triangleq \\
& \quad \wedge \text{IsLeader}(i) \\
& \quad \wedge \text{zabState}[i] = \text{DISCOVERY}
\end{aligned}$$

$\wedge IsQuorum(ackRecv[i])$   
 $\wedge zabState' = [zabState \text{ EXCEPT } ![i] = SYNCHRONIZATION]$   
 $\wedge currentEpoch' = [currentEpoch \text{ EXCEPT } ![i] = acceptedEpoch[i]]$   
 $\wedge initialHistory' = [initialHistory \text{ EXCEPT } ![i] = history[i]]$   
 $\wedge ackRecv' = [ackRecv \text{ EXCEPT } ![i] = ackRecv[i] \cup \{NullPoint\}]$   
 $\wedge ackIndex' = [ackIndex \text{ EXCEPT } ![i][i] = Len(history[i])]$   
 $\wedge UpdateProposal(i, i, lastZxid[i], currentEpoch'[i])$   
 $\wedge \text{LET } epoch \triangleq acceptedEpoch[i]$   
 $\text{IN } epochLeader' = [epochLeader \text{ EXCEPT } ![epoch] = epochLeader[epoch] \cup \{i\}]$   
 $\wedge \text{UNCHANGED } \langle state, lastZxid, acceptedEpoch, history, commitIndex, logicalClock, receiveVotes, outC$   
 $\text{recvQueue, waitNotmsg, leadingVoteSet, learners, cepochRecv, ackldRecv, forwarding}$   
 $\text{committedCounter, tempMaxEpoch, followerVarsZ, proposalMsgsLog,}$   
 $\text{inherentViolated, msgVarsZ, idTable} \rangle$   
 $\wedge UpdateRecorder(\langle \text{"LeaderTransitionToSynchronization"}, i \rangle)$

Note: Set *cepochRecv*, *ackRecv*, *ackldRecv* to  $\{NullPoint\}$  in corresponding three actions to make sure that the prospective leader will not broadcast *NEWEPOCH/NEWLEADER/COMMITLD* twice.

$RECOVERYSYNC(i, j) \triangleq$   
 $\text{LET } canSync \triangleq \wedge IsLeader(i) \wedge zabState[i] \neq DISCOVERY \wedge IsMyLearner(i, j)$   
 $\wedge IsFollower(j) \wedge zabState[j] = SYNCHRONIZATION \wedge IsMyLeader(j, i)$   
 $\text{IN}$   
 $\wedge canSync$   
 $\wedge SubRECOVERYSYNC(i, j)$   
 $\wedge Send(i, j, [mtype \mapsto NEWLEADER,$   
 $\text{mepoch} \mapsto acceptedEpoch[i],$   
 $\text{mlastZxid} \mapsto lastZxid[i]])$   
 $\wedge \text{UNCHANGED } \langle state, zabState, acceptedEpoch, currentEpoch, logicalClock, receiveVotes, outOfElectio$   
 $\text{leadingVoteSet, learners, cepochRecv, ackRecv, ackldRecv, currentCounter,}$   
 $\text{tempVarsZ, cepochSent, leaderAddr, epochLeader, electionMsgs, idTable} \rangle$   
 $\wedge UpdateRecorder(\langle \text{"RECOVERYSYNC"}, i, j \rangle)$

Follower receives *NEWLEADER*. The follower updates its epoch and history, and sends back *ACK-LD* to leader.

$FollowerHandleNEWLEADER(i, j) \triangleq$   
 $\wedge IsFollower(i)$   
 $\wedge PendingNEWLEADER(i, j)$   
 $\wedge \text{LET } msg \triangleq msgs[j][i][1]$   
 $\text{infoOk} \triangleq \wedge IsMyLeader(i, j)$   
 $\wedge acceptedEpoch[i] = msg.mepoch$   
 $\text{correct} \triangleq \wedge infoOk$   
 $\wedge zabState[i] = SYNCHRONIZATION$   
 $\wedge synced[i]$   
 $\wedge ZxidEqual(lastZxid[i], msg.mlastZxid)$   
 $\text{IN } \wedge infoOk$



$$\begin{aligned}
& \wedge \text{currentEpoch}' = [\text{currentEpoch} \text{ EXCEPT } ![i] = \text{msg.mepoch}] \\
& \wedge \text{UpdateProposal}(i, j, \text{lastZxid}[i], \text{currentEpoch}'[i]) \\
& \wedge \vee \wedge \text{correct} \\
& \quad \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{ACKLD}, \\
& \quad \quad \quad \text{mepoch} \mapsto \text{msg.mepoch}]) \\
& \quad \wedge \text{UNCHANGED } \text{inherentViolated} \\
& \vee \wedge \neg \text{correct} \\
& \quad \wedge \text{PrintT}(\text{"Exception: Follower receives NEWLEADER while it has not completed sync with I"} \\
& \quad \wedge \text{inherentViolated}' = \text{TRUE} \\
& \quad \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{lastZxid}, \text{zabState}, \text{acceptedEpoch}, \text{history}, \text{commitIndex}, \text{logicalClock}, \text{receiveV}, \\
& \quad \text{outOfElection}, \text{recvQueue}, \text{waitNotmsg}, \text{leaderVarsZ}, \text{tempVarsZ}, \text{followerVarsZ}, \text{prop}, \\
& \quad \text{epochLeader}, \text{electionMsgs}, \text{idTable} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"FollowerHandleNEWLEADER"}, i, j \rangle)
\end{aligned}$$

Leader receives *ACK-LD* from a quorum of followers, and sends *COMMIT-LD(UPTODATE)* to followers.

$$\begin{aligned}
& \text{LeaderHandleACKLD}(i, j) \triangleq \\
& \quad \wedge \text{IsLeader}(i) \\
& \quad \wedge \text{PendingACKLD}(i, j) \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \quad \text{infoOk} \triangleq \wedge \text{acceptedEpoch}[i] = \text{msg.mepoch} \\
& \quad \quad \quad \wedge \text{IsMyLearner}(i, j) \\
& \quad \quad \text{replyOk} \triangleq \wedge \text{infoOk} \\
& \quad \quad \quad \wedge \text{NullPoint} \in \text{ackldRecv}[i] \\
& \quad \text{IN } \wedge \text{infoOk} \\
& \quad \quad \wedge \vee \text{leader has broadcast UPTODATE} - \text{ just reply} \\
& \quad \quad \quad \wedge \text{replyOk} \\
& \quad \quad \quad \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{UPTODATE}, \\
& \quad \quad \quad \quad \text{mepoch} \mapsto \text{acceptedEpoch}[i], \\
& \quad \quad \quad \quad \text{mcommit} \mapsto \text{commitIndex}[i]]) \\
& \quad \quad \quad \wedge \text{committedCounter}' = [\text{committedCounter} \text{ EXCEPT } ![i][j] = \text{Maximum}(\{\text{commitIndex}[i] \\
& \quad \vee \text{leader still waits for a quorum's ACKLD to broadcast UPTODATE} \\
& \quad \quad \wedge \neg \text{replyOk} \\
& \quad \quad \wedge \text{Discard}(j, i) \\
& \quad \quad \wedge \text{UNCHANGED } \text{committedCounter} \\
& \quad \quad \wedge \text{ackldRecv}' = [\text{ackldRecv} \text{ EXCEPT } ![i] = \text{IF } j \notin \text{ackldRecv}[i] \text{ THEN } \text{ackldRecv}[i] \cup \{j\} \\
& \quad \quad \quad \text{ELSE } \text{ackldRecv}[i]] \\
& \quad \wedge \text{UNCHANGED } \langle \text{serverVarsZ}, \text{electionVarsZ}, \text{leadingVoteSet}, \text{learners}, \text{ceepochRecv}, \text{ackRecv}, \text{forward}, \\
& \quad \quad \text{ackIndex}, \text{currentCounter}, \text{tempVarsZ}, \text{followerVarsZ}, \text{verifyVarsZ}, \text{electionMsgs}, \text{id} \rangle \\
& \quad \wedge \text{UpdateRecorder}(\langle \text{"LeaderHandleACKLD"}, i, j \rangle)
\end{aligned}$$

$$\begin{aligned}
& \text{LeaderTransitionToBroadcast}(i) \triangleq \\
& \quad \wedge \text{IsLeader}(i) \\
& \quad \wedge \text{zabState}[i] = \text{SYNCHRONIZATION} \\
& \quad \wedge \text{IsQuorum}(\text{ackldRecv}[i])
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{commitIndex}' = [\text{commitIndex} \quad \text{EXCEPT } ![i] = \text{Len}(\text{history}[i])] \\
& \wedge \text{zabState}' = [\text{zabState} \quad \text{EXCEPT } ![i] = \text{BROADCAST}] \\
& \wedge \text{currentCounter}' = [\text{currentCounter} \quad \text{EXCEPT } ![i] = 0] \\
& \wedge \text{sendCounter}' = [\text{sendCounter} \quad \text{EXCEPT } ![i] = 0] \\
& \wedge \text{committedIndex}' = [\text{committedIndex} \quad \text{EXCEPT } ![i] = \text{Len}(\text{history}[i])] \\
& \wedge \text{ackldRecv}' = [\text{ackldRecv} \quad \text{EXCEPT } ![i] = \text{ackldRecv}[i] \cup \{\text{NullPoint}\}] \\
& \wedge \text{BroadcastUPTODATE}(i, [\text{mtype} \mapsto \text{UPTODATE}, \\
& \quad \text{mepoch} \mapsto \text{acceptedEpoch}[i], \\
& \quad \text{mcommit} \mapsto \text{Len}(\text{history}[i])]) \quad \text{In actual UPTODATE doesn't carry this info} \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{lastZxid}, \text{acceptedEpoch}, \text{history}, \text{electionVarsZ}, \text{leadingVoteSet}, \\
& \quad \text{ceepochRecv}, \text{ackRecv}, \text{forwarding}, \text{ackIndex}, \text{committedCounter}, \text{tempVarsZ}, \text{follower}, \\
& \quad \text{electionMsgs}, \text{idTable} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"LeaderTransitionToBroadcast"}, i \rangle) \\
& \text{FollowerHandleUPTODATE}(i, j) \triangleq \\
& \quad \wedge \text{IsFollower}(i) \\
& \quad \wedge \text{PendingUPTODATE}(i, j) \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \quad \text{infoOk} \triangleq \wedge \text{IsMyLeader}(i, j) \\
& \quad \quad \quad \wedge \text{acceptedEpoch}[i] = \text{msg.mepoch} \\
& \quad \quad \text{correct} \triangleq \wedge \text{infoOk} \\
& \quad \quad \quad \wedge \text{zabState}[i] = \text{SYNCHRONIZATION} \\
& \quad \quad \quad \wedge \text{currentEpoch}[i] = \text{msg.mepoch} \\
& \quad \text{IN} \quad \wedge \text{infoOk} \\
& \quad \quad \wedge \vee \wedge \text{correct} \\
& \quad \quad \quad \wedge \text{commitIndex}' = [\text{commitIndex} \quad \text{EXCEPT } ![i] = \text{Maximum}(\{\text{commitIndex}[i], \text{msg.mcommit}\})] \\
& \quad \quad \quad \wedge \text{zabState}' = [\text{zabState} \quad \text{EXCEPT } ![i] = \text{BROADCAST}] \\
& \quad \quad \quad \wedge \text{UNCHANGED } \text{inherentViolated} \\
& \quad \quad \vee \wedge \neg \text{correct} \\
& \quad \quad \quad \wedge \text{PrintT}(\text{"Exception: Follower receives UPTODATE while its ZabState is not SYNCHRONIZATION"}) \\
& \quad \quad \quad \wedge \text{inherentViolated}' = \text{TRUE} \\
& \quad \quad \quad \wedge \text{UNCHANGED } \langle \text{commitIndex}, \text{zabState} \rangle \\
& \quad \wedge \text{Discard}(j, i) \\
& \quad \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{lastZxid}, \text{acceptedEpoch}, \text{history}, \text{electionVarsZ}, \text{leaderVarsZ}, \text{tempVarsZ}, \\
& \quad \quad \text{proposalMsgsLog}, \text{epochLeader}, \text{electionMsgs}, \text{idTable} \rangle \\
& \quad \wedge \text{UpdateRecorder}(\langle \text{"FollowerHandleUPTODATE"}, i, j \rangle)
\end{aligned}$$


---

Leader receives client request and broadcasts *PROPOSAL*.

*Note1:* In production, any server in traffic can receive requests and forward it to leader if necessary. We choose to let leader be the sole one who can receive requests, to simplify spec and keep correctness at the same time.

*Note2:* To compress state space, we now choose to merge action client request and action leader broadcasts proposal into one action.

$$\begin{aligned}
& \text{ClientRequestAndLeaderBroadcastProposal}(i, v) \triangleq \\
& \quad \wedge \text{IsLeader}(i) \\
& \quad \wedge \text{zabState}[i] = \text{BROADCAST}
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{currentCounter}' = [\text{currentCounter} \text{ EXCEPT } ![i] = \text{currentCounter}[i] + 1] \\
& \wedge \text{LET } \text{newTransaction} \triangleq [\text{epoch} \mapsto \text{acceptedEpoch}[i], \\
& \quad \text{counter} \mapsto \text{currentCounter}'[i], \\
& \quad \text{value} \mapsto v] \\
& \text{IN } \wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{Append}(\text{history}[i], \text{newTransaction})] \\
& \wedge \text{lastZxid}' = [\text{lastZxid} \text{ EXCEPT } ![i] = \langle \text{acceptedEpoch}[i], \text{currentCounter}'[i] \rangle] \\
& \wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i][i] = \text{Len}(\text{history}'[i])] \\
& \wedge \text{UpdateProposal}(i, i, \text{lastZxid}'[i], \text{currentEpoch}[i]) \\
& \wedge \text{Broadcast}(i, [\text{mtype} \mapsto \text{PROPOSAL}, \\
& \quad \text{mepoch} \mapsto \text{acceptedEpoch}[i], \\
& \quad \text{mproposal} \mapsto \text{newTransaction}]) \\
& \wedge \text{LET } m \triangleq [\text{msource} \mapsto i, \text{mepoch} \mapsto \text{acceptedEpoch}[i], \text{mtype} \mapsto \text{PROPOSAL}, \text{mproposal} \mapsto \text{newTransaction}] \\
& \text{IN } \text{proposalMsgsLog}' = \text{proposalMsgsLog} \cup \{m\} \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{zabState}, \text{acceptedEpoch}, \text{commitIndex}, \text{logicalClock}, \text{receiveVotes}, \\
& \quad \text{recvQueue}, \text{waitNotmsg}, \text{leadingVoteSet}, \text{learners}, \text{ceepochRecv}, \text{ackRecv}, \text{ackldRecv}, \\
& \quad \text{committedCounter}, \text{tempVarsZ}, \text{followerVarsZ}, \text{epochLeader}, \text{inherentViolated}, \\
& \quad \text{electionMsgs}, \text{idTable} \rangle \\
& \wedge \text{UpdateRecorder}(\langle \text{"ClientRequestAndLeaderBroadcastProposal"}, i, v \rangle) \\
& \text{LeaderBroadcastProposal}(i) \triangleq \\
& \quad \wedge \text{IsLeader}(i) \\
& \quad \wedge \text{zabState}[i] = \text{BROADCAST} \\
& \quad \wedge \text{sendCounter}[i] < \text{currentCounter}[i] \\
& \quad \wedge \text{LET } \text{toBeSentCounter} \triangleq \text{sendCounter}[i] + 1 \\
& \quad \quad \text{toBeSentIndex} \triangleq \text{Len}(\text{initialHistory}[i]) + \text{toBeSentCounter} \\
& \quad \quad \text{toBeSentEntry} \triangleq \text{history}[i][\text{toBeSentIndex}] \\
& \quad \text{IN } \wedge \text{Broadcast}(i, [\text{mtype} \mapsto \text{PROPOSAL}, \\
& \quad \quad \text{mepoch} \mapsto \text{acceptedEpoch}[i], \\
& \quad \quad \text{mproposal} \mapsto \text{toBeSentEntry}]) \\
& \quad \wedge \text{sendCounter}' = [\text{sendCounter} \text{ EXCEPT } ![i] = \text{toBeSentCounter}] \\
& \quad \wedge \text{LET } m \triangleq [\text{msource} \mapsto i, \text{mepoch} \mapsto \text{acceptedEpoch}[i], \text{mtype} \mapsto \text{PROPOSAL}, \\
& \quad \quad \text{mproposal} \mapsto \text{toBeSentEntry}] \\
& \quad \text{IN } \text{proposalMsgsLog}' = \text{proposalMsgsLog} \cup \{m\} \\
& \quad \wedge \text{UpdateRecorder}(\langle \text{"LeaderBroadcastProposal"}, i \rangle) \\
& \quad \wedge \text{UNCHANGED } \langle \text{serverVarsZ}, \text{electionVarsZ}, \text{leadingVoteSet}, \text{learners}, \text{ceepochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{forwarding}, \\
& \quad \quad \text{ackIndex}, \text{currentCounter}, \text{committedCounter}, \text{tempVarsZ}, \text{followerVarsZ}, \text{epochLeader}, \\
& \quad \quad \text{inherentViolated}, \text{electionMsgs}, \text{idTable} \rangle
\end{aligned}$$

Follower accepts proposal and append it to history.

$$\begin{aligned}
& \text{FollowerHandlePROPOSAL}(i, j) \triangleq \\
& \quad \wedge \text{IsFollower}(i) \\
& \quad \wedge \text{PendingPROPOSAL}(i, j) \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \quad \text{infoOk} \triangleq \wedge \text{IsMyLeader}(i, j) \\
& \quad \quad \quad \wedge \text{acceptedEpoch}[i] = \text{msg.mepoch} \\
& \quad \quad \text{correct} \triangleq \wedge \text{infoOk}
\end{aligned}$$

$\wedge \text{zabState}[i] \neq \text{DISCOVERY}$   
 $\wedge \text{synced}[i]$   
 $\text{logOk} \triangleq \begin{array}{l} \text{the first PROPOSAL in this epoch} \\ \vee \wedge \text{msg.mproposal.counter} = 1 \\ \quad \wedge \vee \text{Len}(\text{history}[i]) = 0 \\ \quad \vee \wedge \text{Len}(\text{history}[i]) > 0 \\ \quad \quad \wedge \text{history}[i][\text{Len}(\text{history}[i])].\text{epoch} < \text{msg.mepoch} \\ \text{not the first PROPOSAL in this epoch} \\ \vee \wedge \text{msg.mproposal.counter} > 1 \\ \quad \wedge \text{Len}(\text{history}[i]) > 0 \\ \quad \wedge \text{history}[i][\text{Len}(\text{history}[i])].\text{epoch} = \text{msg.mepoch} \\ \quad \wedge \text{history}[i][\text{Len}(\text{history}[i])].\text{counter} = \text{msg.mproposal.counter} - 1 \end{array}$   
IN  $\wedge \text{infoOk}$   
 $\wedge \vee \wedge \text{correct}$   
 $\quad \wedge \vee \wedge \text{logOk}$   
 $\quad \quad \wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{Append}(\text{history}[i], \text{msg.mproposal})]$   
 $\quad \quad \wedge \text{lastZxid}' = [\text{lastZxid} \text{ EXCEPT } ![i] = \langle \text{msg.mepoch}, \text{msg.mproposal.counter} \rangle]$   
 $\quad \quad \wedge \text{UpdateProposal}(i, j, \text{lastZxid}'[i], \text{currentEpoch}[i])$   
 $\quad \quad \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{ACK},$   
 $\quad \quad \quad \text{mepoch} \mapsto \text{acceptedEpoch}[i],$   
 $\quad \quad \quad \text{mzxid} \mapsto \langle \text{msg.mepoch}, \text{msg.mproposal.counter} \rangle])$   
 $\quad \quad \wedge \text{UNCHANGED } \text{inherentViolated}$   
 $\quad \vee \wedge \neg \text{logOk}$   
 $\quad \quad \wedge \text{PrintT}(\text{"Exception: Follower receives PROPOSAL while the transaction is not the next"})$   
 $\quad \quad \wedge \text{inherentViolated}' = \text{TRUE}$   
 $\quad \quad \wedge \text{Discard}(j, i)$   
 $\quad \quad \wedge \text{UNCHANGED } \langle \text{history}, \text{lastZxid}, \text{currentVote} \rangle$   
 $\quad \vee \wedge \neg \text{correct}$   
 $\quad \quad \wedge \text{PrintT}(\text{"Exception: Follower receives PROPOSAL while it has not completed sync with leader"})$   
 $\quad \quad \wedge \text{inherentViolated}' = \text{TRUE}$   
 $\quad \quad \wedge \text{Discard}(j, i)$   
 $\quad \quad \wedge \text{UNCHANGED } \langle \text{history}, \text{lastZxid}, \text{currentVote} \rangle$   
 $\wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{zabState}, \text{acceptedEpoch}, \text{commitIndex}, \text{logicalClock}, \text{receiveVotes},$   
 $\quad \text{outOfElection}, \text{recvQueue}, \text{waitNotmsg}, \text{leaderVarsZ}, \text{tempVarsZ}, \text{followerVarsZ}, \text{propose},$   
 $\quad \text{epochLeader}, \text{electionMsgs}, \text{idTable} \rangle$   
 $\wedge \text{UpdateRecorder}(\langle \text{"FollowerHandlePROPOSAL"}, i, j \rangle)$   
 $\quad \text{Create a commit packet and send it to all the members of the quorum.}$   
 $\text{LeaderCommit}(s, \text{source}, \text{index}, \text{zxid}) \triangleq$   
 $\quad \wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![s] = \text{index}]$   
 $\quad \wedge \text{DiscardAndBroadcast}(s, \text{source}, [\text{mtype} \mapsto \text{COMMIT},$   
 $\quad \quad \text{mepoch} \mapsto \text{acceptedEpoch}[s],$   
 $\quad \quad \text{mzxid} \mapsto \langle \text{zxid.epoch}, \text{zxid.counter} \rangle])$   
 $\quad \text{return true if committed, otherwise false.}$

$LeaderTryToCommit(s, zxid, follower) \triangleq$   
 LET  $pindex \triangleq Len(initialHistory[s]) + zxid.counter$   
 Only when all proposals before  $zxid$  all committed, this proposal can be permitted to be committed.  
 $allProposalsBeforeCommitted \triangleq \vee zxid.counter = 1$   
 $\vee \wedge zxid.counter > 1$   
 $\wedge commitIndex[s] \geq pindex - 1$   
 In order to be committed, a proposal must be accepted by a quorum.  
 $agreeSet \triangleq \{s\} \cup \{k \in (Server \setminus \{s\}) : ackIndex'[s][k] \geq pindex\}$   
 $hasAllQuorums \triangleq IsQuorum(agreeSet)$   
 Commit proposals in order.  
 $ordered \triangleq commitIndex[s] + 1 = pindex$   
 IN  $\vee \wedge \vee \neg allProposalsBeforeCommitted$   
 $\vee \neg hasAllQuorums$   
 $\wedge Discard(follower, s)$   
 $\wedge UNCHANGED \langle inherentViolated, commitIndex \rangle$   
 $\vee \wedge allProposalsBeforeCommitted$   
 $\wedge hasAllQuorums$   
 $\wedge \vee \wedge \neg ordered$   
 $\wedge PrintT(\text{"Exception: Committing } zxid \text{ " } \circ zxid \circ \text{"not first."})$   
 $\wedge inherentViolated' = TRUE$   
 $\vee \wedge ordered$   
 $\wedge UNCHANGED inherentViolated$   
 $\wedge LeaderCommit(s, follower, pindex, zxid)$   
 Keep a count of acks that are received by the leader for a particular proposal, and commit the proposal.  
 $LeaderProcessACK(i, j) \triangleq$   
 $\wedge IsLeader(i)$   
 $\wedge PendingACK(i, j)$   
 $\wedge LET msg \triangleq msgs[j][i][1]$   
 $infoOk \triangleq \wedge j \in forwarding[i]$   
 $\wedge acceptedEpoch[i] = msg.mepoch$   
 $correct \triangleq \wedge infoOk$   
 $\wedge zabState[i] = BROADCAST$   
 $\wedge currentCounter[i] \geq msg.mzxid[2]$   
 $noOutstanding \triangleq commitIndex[i] = currentCounter[i]$   
 $outstandingProposals: \text{proposals in } history[commitIndex + 1 : currentCounter]$   
 $hasCommitted \triangleq commitIndex[i] \geq Len(initialHistory[i]) + msg.mzxid[2]$   
 $\text{namely, } lastCommitted \geq zxid$   
 $logOk \triangleq \wedge infoOk$   
 $\wedge ackIndex[i][j] + 1 = Len(initialHistory[i]) + msg.mzxid[2]$   
 $\text{everytime, } ackIndex \text{ should just increase by 1}$   
 IN  $\wedge infoOk$   
 $\wedge \vee \wedge correct$   
 $\wedge logOk$

```

    ∧ ∨ ∧ ∨ noOutstanding
      ∨ hasCommitted
      ∧ PrintT("Note: outstanding is 0 / proposal has already been committed.")
      ∧ Discard(j, i)
      ∧ UNCHANGED ⟨ackIndex, commitIndex, inherentViolated⟩
    ∨ ∧ ¬noOutstanding
      ∧ ¬hasCommitted
      ∧ ackIndex' = [ackIndex EXCEPT ![i][j] = Len(initialHistory[i] + msg.mzxid[2])]
      ∧ LeaderTryToCommit(i, ToZxid(msg.mzxid), j)
  ∨ ∧ ∨ ¬correct
    ∨ ¬logOk
    ∧ PrintT("Exception: ackIndex doesn't increase monotonically.")
    ∧ inherentViolated' = TRUE
    ∧ Discard(j, i)
    ∧ UNCHANGED ⟨ackIndex, commitIndex⟩
  ∧ UNCHANGED ⟨state, currentEpoch, lastZxid, zabState, acceptedEpoch, history, electionVarsZ,
    leadingVoteSet, learners, ceepochRecv, ackRecv, ackldRecv,
    forwarding, currentCounter, committedCounter, tempVarsZ,
    followerVarsZ, proposalMsgsLog, epochLeader, electionMsgs, idTable⟩
  ∧ UpdateRecorder(("LeaderProcessACK", i, j))

LeaderAdvanceCommit(i) ≜
  ∧ IsLeader(i)
  ∧ zabState[i] = BROADCAST
  ∧ commitIndex[i] < Len(history[i])
  ∧ LET Agree(index) ≜ {i} ∪ {k ∈ (Server \ {i}) : ackIndex[i][k] ≥ index}
    agreeIndexes ≜ {index ∈ (commitIndex[i] + 1) .. Len(history[i]) :
      Agree(index) ∈ Quorums}
    newCommitIndex ≜ IF agreeIndexes ≠ {} THEN Maximum(agreeIndexes)
      ELSE commitIndex[i]
  IN commitIndex' = [commitIndex EXCEPT ![i] = newCommitIndex]
  ∧ UpdateRecorder(("LeaderAdvanceCommit", i))
  ∧ UNCHANGED ⟨state, currentEpoch, lastZxid, zabState, acceptedEpoch, history, electionVarsZ, leaderVarsZ,
    tempVarsZ, followerVarsZ, verifyVarsZ, msgVarsZ, idTable⟩

LeaderBroadcastCommit(i) ≜
  ∧ IsLeader(i)
  ∧ zabState[i] = BROADCAST
  ∧ committedIndex[i] < commitIndex[i]
  ∧ Len(initialHistory[i] + sendCounter[i] > committedIndex[i]
  ∧ LET newCommittedIndex ≜ committedIndex[i] + 1
  IN ∧ Broadcast(i, [mtype ↦ COMMIT,
    mepoch ↦ acceptedEpoch[i],
    mzxid ↦ ⟨history[i][newCommittedIndex].epoch, history[i][newCommittedIndex].counter⟩])

  ∧ committedIndex' = [committedIndex EXCEPT ![i] = committedIndex[i] + 1]
  ∧ committedCounter' = [committedCounter EXCEPT ![i] = [v ∈ Server ↦ IF ∧
    v ∈ forwarding[i]

```

```

         $\wedge$  committedCounter[i][v] <
        history[i][newCommittedIndex].counter

        THEN history[i][newCommittedIndex].counter

        ELSE committedCounter[i][v]]

 $\wedge$  UpdateRecorder( $\langle$ “LeaderBroadcastCommit”, i $\rangle$ )
 $\wedge$ UNCHANGED  $\langle$ serverVarsZ, electionVarsZ, leadingVoteSet, learners, cepochRecv, ackeRecv, ackldRecv,

        forwarding, ackIndex, currentCounter, sendCounter, tempVarsZ, followerVarsZ, verifyVarsZ,

        electionMsgs, idTable $\rangle$ 

```

Follower receives COMMIT and commits transaction.

```

FollowerHandleCOMMIT(i, j)  $\triangleq$ 
     $\wedge$  IsFollower(i)
     $\wedge$  PendingCOMMIT(i, j)
     $\wedge$  LET msg  $\triangleq$  msgs[j][i][1]
        infoOk  $\triangleq$   $\wedge$  IsMyLeader(i, j)
             $\wedge$  acceptedEpoch[i] = msg.mepoch
        correct  $\triangleq$   $\wedge$  infoOk
             $\wedge$  zabState[i]  $\neq$  DISCOVERY
             $\wedge$  synced[i]
        mindex  $\triangleq$  IF Len(history[i]) = 0 THEN - 1
            ELSE IF  $\exists$  idx  $\in$  1 .. Len(history[i]) : PZxidEqual(history[i][idx], msg.mzxid)
                THEN CHOOSE idx  $\in$  1 .. Len(history[i]) : PZxidEqual(history[i][idx], msg.mzxid)
                ELSE - 1
        logOk  $\triangleq$  mindex > 0
        latest  $\triangleq$  commitIndex[i] + 1 = mindex
    IN  $\wedge$  infoOk
         $\wedge$   $\vee$   $\wedge$  correct
             $\wedge$   $\vee$   $\wedge$  logOk
                 $\wedge$   $\vee$   $\wedge$  latest
                     $\wedge$  commitIndex' = [commitIndex EXCEPT !i] = commitIndex[i] + 1
                     $\wedge$  UNCHANGED inherentViolated
                 $\vee$   $\wedge$   $\neg$ latest
                     $\wedge$  PrintT(“Note: Follower receives COMMIT while the index is not the next comm
                     $\wedge$  inherentViolated' = TRUE
                     $\wedge$  UNCHANGED commitIndex
             $\vee$   $\wedge$   $\neg$ logOk
                 $\wedge$  PrintT(“Exception: Follower receives COMMIT while the transaction has not been sa
                 $\wedge$  inherentViolated' = TRUE
                 $\wedge$  UNCHANGED commitIndex
         $\vee$   $\wedge$   $\neg$ correct
             $\wedge$  PrintT(“Exception: Follower receives COMMIT while it has not completed sync with leade
             $\wedge$  inherentViolated' = TRUE
             $\wedge$  UNCHANGED commitIndex

```

$\wedge \text{Discard}(j, i)$   
 $\wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{lastZxid}, \text{zabState}, \text{acceptedEpoch}, \text{history}, \text{electionVarsZ}, \text{leader}$   
 $\text{tempVarsZ}, \text{followerVarsZ}, \text{proposalMsgsLog}, \text{epochLeader}, \text{electionMsgs}, \text{idTable} \rangle$   
 $\wedge \text{UpdateRecorder}(\langle \text{"FollowerHandleCOMMIT"}, i, j \rangle)$

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

$\text{FilterNonexistentMessage}(i) \triangleq$   
 $\wedge \exists j \in \text{Server} \setminus \{i\} : \wedge \text{msgs}[j][i] \neq \langle \rangle$   
 $\wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1]$   
 $\text{IN}$   
 $\vee \wedge \text{IsLeader}(i)$   
 $\wedge \text{LET } \text{infoOk} \triangleq \wedge j \in \text{learners}[i]$   
 $\wedge \text{acceptedEpoch}[i] = \text{msg.mepoch}$   
 $\text{IN}$   
 $\vee \text{msg.mtype} = \text{LEADERINFO}$   
 $\vee \text{msg.mtype} = \text{NEWLEADER}$   
 $\vee \text{msg.mtype} = \text{UPTODATE}$   
 $\vee \text{msg.mtype} = \text{PROPOSAL}$   
 $\vee \text{msg.mtype} = \text{COMMIT}$   
 $\vee \wedge j \notin \text{learners}[i]$   
 $\wedge \text{msg.mtype} = \text{FOLLOWERINFO}$   
 $\vee \wedge \neg \text{infoOk}$   
 $\wedge \vee \text{msg.mtype} = \text{ACKEPOCH}$   
 $\vee \text{msg.mtype} = \text{ACKLD}$   
 $\vee \text{msg.mtype} = \text{ACK}$   
 $\vee \wedge \text{IsFollower}(i)$   
 $\wedge \text{LET } \text{infoOk} \triangleq \wedge j = \text{leaderAddr}[i]$   
 $\wedge \text{acceptedEpoch}[i] = \text{msg.mepoch}$   
 $\text{IN}$   
 $\vee \text{msg.mtype} = \text{FOLLOWERINFO}$   
 $\vee \text{msg.mtype} = \text{ACKEPOCH}$   
 $\vee \text{msg.mtype} = \text{ACKLD}$   
 $\vee \text{msg.mtype} = \text{ACK}$   
 $\vee \wedge j \neq \text{leaderAddr}[i]$   
 $\wedge \text{msg.mtype} = \text{LEADERINFO}$   
 $\vee \wedge \neg \text{infoOk}$   
 $\wedge \vee \text{msg.mtype} = \text{NEWLEADER}$   
 $\vee \text{msg.mtype} = \text{UPTODATE}$   
 $\vee \text{msg.mtype} = \text{PROPOSAL}$   
 $\vee \text{msg.mtype} = \text{COMMIT}$   
 $\vee \text{IsLooking}(i)$   
 $\wedge \text{Discard}(j, i)$   
 $\wedge \text{inherentViolated}' = \text{TRUE}$   
 $\wedge \text{UnchangeRecorder}$



$\wedge \text{UNCHANGED } \langle \text{serverVarsZ}, \text{electionVarsZ}, \text{leaderVarsZ}, \text{tempVarsZ}, \text{followerVarsZ}, \text{proposalMsgsL}, \text{epochLeader}, \text{electionMsgs}, \text{idTable} \rangle$

---

Defines how the variables may transition.

$\text{NextZ} \triangleq$

*FLE* modlue

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

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

$\vee \exists i, j \in \text{Server} : \text{EstablishConnection}(i, j)$   
 $\vee \exists i \in \text{Server} : \text{FollowerSendFOLLOWERINFO}(i)$   
 $\vee \exists i, j \in \text{Server} : \text{LeaderHandleFOLLOWERINFO}(i, j)$   
 $\vee \exists i \in \text{Server} : \text{LeaderBroadcastLEADERINFO}(i)$   
 $\vee \exists i, j \in \text{Server} : \text{FollowerHandleLEADERINFO}(i, j)$   
 $\vee \exists i, j \in \text{Server} : \text{LeaderHandleACKEPOCH}(i, j)$   
 $\vee \exists i \in \text{Server} : \text{LeaderTransitionToSynchronization}(i)$   
 $\vee \exists i, j \in \text{Server} : \text{RECOVERYSYNC}(i, j)$   
 $\vee \exists i, j \in \text{Server} : \text{FollowerHandleNEWLEADER}(i, j)$   
 $\vee \exists i, j \in \text{Server} : \text{LeaderHandleACKLD}(i, j)$   
 $\vee \exists i \in \text{Server} : \text{LeaderTransitionToBroadcast}(i)$   
 $\vee \exists i, j \in \text{Server} : \text{FollowerHandleUPTODATE}(i, j)$

Zab module – Broadcast part

$\vee \exists i \in \text{Server}, v \in \text{Value} : \text{ClientRequestAndLeaderBroadcastProposal}(i, v)$   
 $\vee \exists i \in \text{Server}, v \in \text{Value} : \text{ClientRequest}(i, v)$   
 $\vee \exists i \in \text{Server} : \text{LeaderBroadcastProposal}(i)$   
 $\vee \exists i, j \in \text{Server} : \text{FollowerHandlePROPOSAL}(i, j)$   
 $\vee \exists i, j \in \text{Server} : \text{LeaderProcessACK}(i, j)$   
 $\vee \exists i \in \text{Server} : \text{LeaderAdvanceCommit}(i)$   
 $\vee \exists i \in \text{Server} : \text{LeaderBroadcastCommit}(i)$   
 $\vee \exists i, j \in \text{Server} : \text{FollowerHandleCOMMIT}(i, j)$

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

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

$\text{SpecZ} \triangleq \text{InitZ} \wedge \Box[\text{NextZ}]_{\text{vars}}$

---

Define safety properties of Zab 1.0 protocol.

$ShouldNotBeTriggered \triangleq inherentViolated = FALSE$

There is most one established leader for a certain epoch.

$Leadership1 \triangleq \forall i, j \in Server :$

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

$Leadership2 \triangleq \forall epoch \in 1 \dots MAXEPOCH : Cardinality(epochLeader[epoch]) < 2$

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

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

LET  $smaller \triangleq Minimum(\{commitIndex[i], commitIndex[j]\})$

IN  $\vee smaller = 0$

$\vee \wedge smaller > 0$

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

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

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

$\wedge IsFollower(i)$

$\wedge commitIndex[i] > 0$

$\Rightarrow \forall index \in 1 \dots commitIndex[i] : \exists msg \in proposalMsgsLog :$

$\vee \wedge msg.mtype = PROPOSAL$

$\wedge TransactionEqual(msg.mproposal, history[i][index])$

$\vee \wedge msg.mtype = "RECOVERYSYNC"$

$\wedge \exists tindex \in 1 \dots Len(msg.mproposals) : TransactionEqual(msg.mproposals[tindex], history[i][index])$

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

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

$\wedge IsFollower(i) \wedge commitIndex[i] > 0$

$\wedge IsFollower(j) \wedge commitIndex[j] > 0$

$\Rightarrow$

$\forall index1 \in 1 \dots commitIndex[i], index2 \in 1 \dots commitIndex[j] :$

$\vee \exists indexj \in 1 \dots commitIndex[j] :$

$TransactionEqual(history[j][indexj], history[i][index1])$

$\vee \exists indexi \in 1 \dots commitIndex[i] :$

$TransactionEqual(history[i][indexi], history[j][index2])$

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

$TotalOrder \triangleq \forall i, j \in Server : commitIndex[i] \geq 2 \wedge commitIndex[j] \geq 2$

$\Rightarrow \forall indexi1 \in 1 \dots (commitIndex[i] - 1) : \forall indexi2 \in (indexi1 + 1) \dots commitIndex[i] :$

LET  $logOk \triangleq \exists index \in 1 \dots commitIndex[j] : TransactionEqual(history[i][indexi2], history[j][index])$

IN  $\vee \neg logOk$

$\vee \wedge logOk$

$$\begin{aligned} & \wedge \exists \text{index}j2 \in 1 \dots \text{commitIndex}[j] : \\ & \quad \wedge \text{TransactionEqual}(\text{history}[i][\text{index}i2], \text{history}[j][\text{index}j2]) \\ & \quad \wedge \exists \text{index}j1 \in 1 \dots (\text{index}j2 - 1) : \text{TransactionEqual}(\text{history}[i][\text{index}i1], \text{history}[j][\text{index}j1]) \end{aligned}$$

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

$$\begin{aligned} \text{LocalPrimaryOrder} &\triangleq \text{LET } \text{mset}(i, e) \triangleq \{ \text{msg} \in \text{proposalMsgsLog} : \wedge \text{msg.mtype} = \text{PROPOSAL} \\ & \quad \wedge \text{msg.msource} = i \\ & \quad \wedge \text{msg.mepoch} = e \} \\ & \quad \text{mentries}(i, e) \triangleq \{ \text{msg.mproposal} : \text{msg} \in \text{mset}(i, e) \} \\ \text{IN } & \forall i \in \text{Server} : \forall e \in 1 \dots \text{currentEpoch}[i] : \\ & \quad \vee \text{Cardinality}(\text{mentries}(i, e)) < 2 \\ & \quad \vee \wedge \text{Cardinality}(\text{mentries}(i, e)) \geq 2 \\ & \quad \wedge \exists \text{tsc1}, \text{tsc2} \in \text{mentries}(i, e) : \\ & \quad \quad \vee \text{TransactionEqual}(\text{tsc1}, \text{tsc2}) \\ & \quad \quad \vee \wedge \neg \text{TransactionEqual}(\text{tsc1}, \text{tsc2}) \\ & \quad \quad \wedge \text{LET } \text{tscPre} \triangleq \text{IF } \text{TransactionPrecede}(\text{tsc1}, \text{tsc2}) \text{ THEN } \text{tsc1} \text{ ELSE } \text{tsc2} \\ & \quad \quad \quad \text{tscNext} \triangleq \text{IF } \text{TransactionPrecede}(\text{tsc1}, \text{tsc2}) \text{ THEN } \text{tsc2} \text{ ELSE } \text{tsc1} \\ & \quad \quad \text{IN } \forall j \in \text{Server} : \wedge \text{commitIndex}[j] \geq 2 \\ & \quad \quad \quad \wedge \exists \text{index} \in 1 \dots \text{commitIndex}[j] : \text{TransactionEqual}(\text{history}[i][\text{index}], \text{history}[j][\text{index}]) \\ & \quad \quad \Rightarrow \exists \text{index}2 \in 1 \dots \text{commitIndex}[j] : \\ & \quad \quad \quad \wedge \text{TransactionEqual}(\text{history}[j][\text{index}2], \text{tscNext}) \\ & \quad \quad \quad \wedge \text{index}2 > 1 \\ & \quad \quad \quad \wedge \exists \text{index}1 \in 1 \dots (\text{index}2 - 1) : \text{TransactionEqual}(\text{history}[j][\text{index}1], \text{tscPre}) \end{aligned}$$

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

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

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

$$\begin{aligned} \text{PrimaryIntegrity} &\triangleq \forall i, j \in \text{Server} : \wedge \text{IsLeader}(i) \\ & \quad \wedge \text{IsFollower}(j) \\ & \quad \wedge \text{commitIndex}[j] \geq 1 \\ & \quad \Rightarrow \forall \text{index} \in 1 \dots \text{commitIndex}[j] : \vee \text{history}[j][\text{index}].\text{epoch} \geq \text{currentEpoch}[i] \\ & \quad \quad \vee \wedge \text{history}[j][\text{index}].\text{epoch} < \text{currentEpoch}[i] \\ & \quad \quad \wedge \exists \text{idx} \in 1 \dots \text{commitIndex}[i] : \text{TransactionEqual}(\text{history}[i][\text{idx}], \text{history}[j][\text{index}]) \end{aligned}$$

---

\ \* Modification History  
\ \* Last modified *Fri Oct 08 21:56:29 CST 2021* by Dell  
\ \* Created *Tue Jun 29 22:13:02 CST 2021* by Dell