

---

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*

---

Defined in *FastLeaderElection.tla*:

\ \* The set of server identifiers

CONSTANT *Server*

\ \* *Server* states CONSTANTS *LOOKING*, *FOLLOWING*, *LEADING*

\ \* Message types

CONSTANTS *NOTIFICATION*

\ \* *Timeout* signal

CONSTANT *NONE*

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*

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

---

Defined in *FastLeaderElection.tla*: *Quorums*, *NullPoint*

Return the maximum value from the set *S*

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

Return the minimum value from the set *S*

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

$MAXEPOCH \triangleq 10$

---

Defined in *FastLeaderElection.tla*: *serverVars*:  $\langle state, currentEpoch, lastZxid \rangle$ , *electionVars*:  $\langle currentVote, logicalClock, receiveVotes, outOfElection, recvQueue, waitNotmsg \rangle$ , *leaderVars*:  $\langle leadingVoteSet \rangle$ , *electionMsgs*, *idTable*

The current phase of *server*(*ELECTION*, *DISCOVERY*, *SYNCHRONIZATION*, *BROADCAST*)  
VARIABLE *zabState*

The epoch number of the last *NEWEPOCH(LEADERINFO)* packet accepted  
namely *f.p* in paper, and *currentEpoch* in *Zab.tla*.  
VARIABLE *acceptedEpoch*

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

*commitIndex[i]*: The maximum index of transactions that have been saved in a quorum of servers  
in the perspective of server *i*.(*increases monotonically before restarting*)  
VARIABLE *commitIndex*

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.

*learners[i]*: The set of servers which leader *i* think are connected with *i*.  
VARIABLE *learners*

The messages representing requests and responses sent from one server to another.  
*msgs[i][j]* means the input buffer of server *j* from server *i*.  
VARIABLE *msgs*

The set of followers who has successfully sent *CEPOCH(FOLLOWERINFO)* to leader.(*equals to connectingFollowers in code*)  
VARIABLE *epochRecv*

The set of followers who has successfully sent *ACK-E* to leader.(*equals to electingFollowers in code*)  
VARIABLE *ackRecv*

The set of followers who has successfully sent *ACK-LD* to leader in leader.(*equals to newLeaderProposal in code*)  
VARIABLE *ackldRecv*

The set of servers which leader *i* broadcasts *PROPOSAL* and *COMMIT* to.(*equals to forwardingFollowers in code*)  
VARIABLE *forwarding*

*ackIndex[i][j]*: The latest index that leader *i* has received from follower *j* via *ACK*.  
VARIABLE *ackIndex*

*currentCounter[i]*: The count of transactions that clients request leader *i*.  
VARIABLE *currentCounter*

*sendCounter[i]*: The count of transactions that leader *i* has broadcast in *PROPOSAL*.  
VARIABLE *sendCounter*

*committedIndex[i]*: The maximum index of trasactions that leader *i* has broadcast in *COMMIT*.

VARIABLE *committedIndex*

*committedCounter*[*i*][*j*]: The latest counter of transaction that leader *i* has confirmed that follower *j* has committed.

VARIABLE *committedCounter*

*initialHistory*[*i*]: The initial history if leader *i* in epoch *acceptedEpoch*[*i*].

VARIABLE *initialHistory*

the maximum epoch in *CEPOCH* the prospective leader received from followers.

VARIABLE *tempMaxEpoch*

*ceepochSent*[*i*] = TRUE means follower *i* has sent *CEPOCH(FOLLOWERINFO)* to leader.

VARIABLE *ceepochSent*

*leaderAddr*[*i*]: The leader *id* of follower *i*. We use *leaderAddr* to express whether follower *i* has connected or lost connection.

VARIABLE *leaderAddr*

*synced*[*i*] = TRUE: follower *i* has completed sync with leader.

VARIABLE *synced*

The set of leaders in every epoch, only used in verifying properties.

VARIABLE *epochLeader*

The set of all broadcast messages, only used in verifying properties.

VARIABLE *proposalMsgsLog*

A variable used to check whether there are conditions contrary to the facts.

VARIABLE *inherentViolated*

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

7 variables

$electionVarsZ \triangleq electionVars$  6 variables

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

11 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

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

Add a message to *msgs* – add a message *m* to *msgs*[*i*][*j*].

$Send(i, j, m) \triangleq msgs' = [msgs \text{ EXCEPT } ![i][j] = Append(msgs[i][j], m)]$

Remove a message from *msgs* – discard head of *msgs*[*i*][*j*].

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

ELSE  $msgs$

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

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

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

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

Combination of *Send* and *Discard* – discard head of  $msgs[j][i]$  and add  $m$  into  $msgs[i][j]$ .

$Reply(i, j, m) \triangleq msgs' = [msgs \text{ EXCEPT } ![j][i] = Tail(msgs[j][i]), \\ ![i][j] = Append(msgs[i][j], m)]$

shuffle the input buffer from server  $j(i)$  in server  $i(j)$ .

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

---

$PZxidEqual(p, z) \triangleq p.epoch = z[1] \wedge p.counter = z[2]$

$TransactionEqual(t1, t2) \triangleq \wedge t1.epoch = t2.epoch \\ \wedge t1.counter = t2.counter$

$TransactionPrecede(t1, t2) \triangleq \vee t1.epoch < t2.epoch \\ \vee \wedge t1.epoch = t2.epoch \\ \wedge t1.counter < t2.counter$

---

Define initial values for all variables

$InitServerVarsZ \triangleq \wedge InitServerVars \\ \wedge zabState = [s \in Server \mapsto ELECTION] \\ \wedge acceptedEpoch = [s \in Server \mapsto 0] \\ \wedge history = [s \in Server \mapsto \langle \rangle] \\ \wedge commitIndex = [s \in Server \mapsto 0]$

$InitLeaderVarsZ \triangleq \wedge InitLeaderVars$

$$\begin{aligned}
\wedge \text{learners} &= [s \in \text{Server} \mapsto \{\}] \\
\wedge \text{cepocheRecv} &= [s \in \text{Server} \mapsto \{\}] \\
\wedge \text{ackeRecv} &= [s \in \text{Server} \mapsto \{\}] \\
\wedge \text{ackldRecv} &= [s \in \text{Server} \mapsto \{\}] \\
\wedge \text{ackIndex} &= [s \in \text{Server} \mapsto [v \in \text{Server} \mapsto 0]] \\
\wedge \text{currentCounter} &= [s \in \text{Server} \mapsto 0] \\
\wedge \text{sendCounter} &= [s \in \text{Server} \mapsto 0] \\
\wedge \text{committedIndex} &= [s \in \text{Server} \mapsto 0] \\
\wedge \text{committedCounter} &= [s \in \text{Server} \mapsto [v \in \text{Server} \mapsto 0]] \\
\wedge \text{forwarding} &= [s \in \text{Server} \mapsto \{\}]
\end{aligned}$$

$$\text{InitElectionVarsZ} \triangleq \text{InitElectionVars}$$

$$\begin{aligned}
\text{InitTempVarsZ} &\triangleq \wedge \text{initialHistory} = [s \in \text{Server} \mapsto \langle \rangle] \\
&\wedge \text{tempMaxEpoch} = [s \in \text{Server} \mapsto 0]
\end{aligned}$$

$$\begin{aligned}
\text{InitFollowerVarsZ} &\triangleq \wedge \text{cepocheSent} = [s \in \text{Server} \mapsto \text{FALSE}] \\
&\wedge \text{leaderAddr} = [s \in \text{Server} \mapsto \text{NullPoint}] \\
&\wedge \text{synced} = [s \in \text{Server} \mapsto \text{FALSE}]
\end{aligned}$$

$$\begin{aligned}
\text{InitVerifyVarsZ} &\triangleq \wedge \text{proposalMsgsLog} = \{\} \\
&\wedge \text{epochLeader} = [i \in 1 \dots \text{MAXEPOCH} \mapsto \{\}] \\
&\wedge \text{inherentViolated} = \text{FALSE}
\end{aligned}$$

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

$$\begin{aligned}
\text{InitZ} &\triangleq \wedge \text{InitServerVarsZ} \\
&\wedge \text{InitLeaderVarsZ} \\
&\wedge \text{InitElectionVarsZ} \\
&\wedge \text{InitTempVarsZ} \\
&\wedge \text{InitFollowerVarsZ} \\
&\wedge \text{InitVerifyVarsZ} \\
&\wedge \text{InitMsgVarsZ} \\
&\wedge \text{idTable} = \text{InitializeIdTable}(\text{Server})
\end{aligned}$$

---


$$\begin{aligned}
\text{ZabTurnToLeading}(i) &\triangleq \\
&\wedge \text{zabState}' = [\text{zabState} \text{ EXCEPT } ![i] = \text{DISCOVERY}] \\
&\wedge \text{learners}' = [\text{learners} \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge \text{cepocheRecv}' = [\text{cepocheRecv} \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge \text{ackeRecv}' = [\text{ackeRecv} \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge \text{ackldRecv}' = [\text{ackldRecv} \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge \text{forwarding}' = [\text{forwarding} \text{ EXCEPT } ![i] = \{\}] \\
&\wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i] = [v \in \text{Server} \mapsto \text{IF } v = i \text{ THEN } \text{Len}(\text{history}[i]) \\
&\hspace{15em} \text{ELSE } 0]] \\
&\wedge \text{currentCounter}' = [\text{currentCounter} \text{ EXCEPT } ![i] = 0]
\end{aligned}$$

$$\begin{aligned}
\wedge \text{sendCounter}' &= [\text{sendCounter} \quad \text{EXCEPT } ![i] = 0] \\
\wedge \text{commitIndex}' &= [\text{commitIndex} \quad \text{EXCEPT } ![i] = 0] \\
\wedge \text{committedIndex}' &= [\text{committedIndex} \quad \text{EXCEPT } ![i] = 0] \\
\wedge \text{committedCounter}' &= [\text{committedCounter} \quad \text{EXCEPT } ![i] = [v \in \text{Server} \mapsto \text{IF } v = i \text{ THEN } \text{Len}(\text{history}[v]) \text{ ELSE } 0]] \\
\wedge \text{initialHistory}' &= [\text{initialHistory} \quad \text{EXCEPT } ![i] = \text{history}[i]] \\
\wedge \text{tempMaxEpoch}' &= [\text{tempMaxEpoch} \quad \text{EXCEPT } ![i] = \text{acceptedEpoch}[i]]
\end{aligned}$$

$$\begin{aligned}
\text{ZabTurnToFollowing}(i) &\triangleq \\
\wedge \text{zabState}' &= [\text{zabState} \quad \text{EXCEPT } ![i] = \text{DISCOVERY}] \\
\wedge \text{cepochSent}' &= [\text{cepochSent} \quad \text{EXCEPT } ![i] = \text{FALSE}] \\
\wedge \text{synced}' &= [\text{synced} \quad \text{EXCEPT } ![i] = \text{FALSE}] \\
\wedge \text{commitIndex}' &= [\text{commitIndex} \quad \text{EXCEPT } ![i] = 0]
\end{aligned}$$

#### Fast Leader Election

$$\begin{aligned}
\text{FLEReceiveNotmsg}(i, j) &\triangleq \\
\wedge \text{ReceiveNotmsg}(i, j) & \\
\wedge \text{UNCHANGED} \langle \text{zabState}, \text{acceptedEpoch}, \text{history}, \text{commitIndex}, \text{learners}, \text{cepochRecv}, \text{ackRecv}, \text{ackldRecv}, & \\
&\text{sendCounter}, \text{committedIndex}, \text{committedCounter}, \text{tempVarsZ}, \text{followerVarsZ}, \text{verifyVarsZ} \rangle
\end{aligned}$$

$$\begin{aligned}
\text{FLENotmsgTimeout}(i) &\triangleq \\
\wedge \text{NotmsgTimeout}(i) & \\
\wedge \text{UNCHANGED} \langle \text{zabState}, \text{acceptedEpoch}, \text{history}, \text{commitIndex}, \text{learners}, \text{cepochRecv}, \text{ackRecv}, \text{ackldRecv}, & \\
&\text{sendCounter}, \text{committedIndex}, \text{committedCounter}, \text{tempVarsZ}, \text{followerVarsZ}, \text{verifyVarsZ} \rangle
\end{aligned}$$

$$\begin{aligned}
\text{FLEHandleNotmsg}(i) &\triangleq \\
\wedge \text{HandleNotmsg}(i) & \\
\wedge \text{LET } \text{newState} &\triangleq \text{state}'[i] \\
\text{IN} & \\
\vee \wedge \text{newState} = \text{LEADING} & \\
\wedge \text{ZabTurnToLeading}(i) & \\
\wedge \text{UNCHANGED} \langle \text{cepochSent}, \text{synced} \rangle & \\
\vee \wedge \text{newState} = \text{FOLLOWING} & \\
\wedge \text{ZabTurnToFollowing}(i) & \\
\wedge \text{UNCHANGED} \langle \text{learners}, \text{cepochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{forwarding}, \text{ackIndex}, \text{currentCounter} \rangle & \\
\vee \wedge \text{newState} = \text{LOOKING} & \\
\wedge \text{UNCHANGED} \langle \text{zabState}, \text{learners}, \text{cepochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{forwarding}, \text{ackIndex}, \text{currentCounter}, & \\
&\text{committedIndex}, \text{committedCounter}, \text{tempVarsZ}, \text{cepochSent}, \text{synced} \rangle \\
\wedge \text{UNCHANGED} \langle \text{acceptedEpoch}, \text{history}, \text{leaderAddr}, \text{verifyVarsZ}, \text{msgs} \rangle
\end{aligned}$$

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

$$\begin{aligned}
\text{FLEWaitNewNotmsg}(i) &\triangleq \\
\wedge \text{WaitNewNotmsg}(i) & \\
\wedge \text{UNCHANGED} \langle \text{zabState}, \text{acceptedEpoch}, \text{history}, \text{commitIndex}, \text{learners}, \text{cepochRecv}, \text{ackRecv}, \text{ackldRecv}, & \\
&\text{sendCounter}, \text{committedIndex}, \text{committedCounter}, \text{tempVarsZ}, \text{followerVarsZ}, \text{verifyVarsZ} \rangle
\end{aligned}$$

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

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


---

A sub-action describing how a server transitions from *LEADING/FOLLOWING* to *LOOKING*. Initially I call it 'ZabTimeoutZ', but it will be called not only when timeout, but also when finding a low epoch from leader.

$$\begin{aligned}
& FollowerShutdown(i) \triangleq \\
& \quad \wedge ZabTimeout(i) \\
& \quad \wedge zabState' = [zabState \text{ EXCEPT } ![i] = ELECTION] \\
& \quad \wedge leaderAddr' = [leaderAddr \text{ EXCEPT } ![i] = NullPoint] \\
& LeaderShutdown(i) \triangleq \\
& \quad \wedge ZabTimeout(i) \\
& \quad \wedge zabState' = [zabState \text{ EXCEPT } ![i] = ELECTION] \\
& \quad \wedge leaderAddr' = [s \in Server \mapsto \text{IF } s \in learners[i] \text{ THEN } NullPoint \text{ ELSE } leaderAddr[s]] \\
& \quad \wedge learners' = [learners \text{ EXCEPT } ![i] = \{\}] \\
& \quad \wedge forwarding' = [forwarding \text{ EXCEPT } ![i] = \{\}] \\
& \quad \wedge msgs' = [s \in Server \mapsto [v \in Server \mapsto \text{IF } v \in learners[i] \vee s \in learners[i] \text{ THEN } \langle \rangle \text{ ELSE } msgs[s][v]]] \\
& FollowerTimeout(i) \triangleq \\
& \quad \wedge state[i] = FOLLOWING \\
& \quad \wedge leaderAddr[i] = NullPoint \\
& \quad \wedge FollowerShutdown(i) \\
& \quad \wedge msgs' = [s \in Server \mapsto [v \in Server \mapsto \text{IF } v = i \text{ THEN } \langle \rangle \text{ ELSE } msgs[s][v]]] \\
& \quad \wedge \text{UNCHANGED } \langle acceptedEpoch, history, commitIndex, learners, cepochRecv, ackeRecv, ackldRecv, forwarding, \\
& \quad \quad currentCounter, sendCounter, committedIndex, committedCounter, tempVarsZ, cepochSent, synced \rangle \\
& LeaderTimeout(i) \triangleq \\
& \quad \wedge state[i] = LEADING \\
& \quad \wedge learners[i] \notin Quorums \\
& \quad \wedge LeaderShutdown(i)
\end{aligned}$$

$\wedge$  UNCHANGED  $\langle \text{acceptedEpoch}, \text{history}, \text{commitIndex}, \text{ceepochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{ackIndex}, \text{currentCounter}, \text{sendCounter}, \text{committedIndex}, \text{committedCounter}, \text{tempVarsZ}, \text{ceepochSent}, \text{synced}, \text{verifyVarsZ} \rangle$

---

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

$\text{EstablishConnection}(i, j) \triangleq$   
 $\wedge \text{state}[i] = \text{LEADING} \quad \wedge \text{state}[j] = \text{FOLLOWING}$   
 $\wedge j \notin \text{learners}[i] \quad \wedge \text{leaderAddr}[j] = \text{NullPoint}$   
 $\wedge \text{currentVote}[j].\text{proposedLeader} = i$   
 $\wedge \text{learners}' = [\text{learners} \text{ EXCEPT } ![i] = \text{learners}[i] \cup \{j\}]$  Leader: '*addLearnerHandler(peer)'*  
 $\wedge \text{leaderAddr}' = [\text{leaderAddr} \text{ EXCEPT } ![j] = i]$  Follower: '*connectToLeader(addr, hostname)'*  
 $\wedge$  UNCHANGED  $\langle \text{serverVarsZ}, \text{electionVarsZ}, \text{leadingVoteSet}, \text{ceepochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{forwarding}, \text{currentCounter}, \text{sendCounter}, \text{committedIndex}, \text{committedCounter}, \text{tempVarsZ}, \text{ceepochSent} \rangle$

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

$\text{Timeout}(i, j) \triangleq$   
 $\wedge \text{state}[i] = \text{LEADING} \wedge \text{state}[j] = \text{FOLLOWING}$   
 $\wedge j \in \text{learners}[i] \quad \wedge \text{leaderAddr}[j] = i$   
 The action of leader  $i$ . (corresponding to function '*removeLearnerHandler(peer)'*.)  
 $\wedge \text{learners}' = [\text{learners} \text{ EXCEPT } ![i] = \text{learners}[i] \setminus \{j\}]$   
 $\wedge \text{forwarding}' = [\text{forwarding} \text{ EXCEPT } ![i] = \text{IF } j \in \text{forwarding}[i] \text{ THEN } \text{forwarding}[i] \setminus \{j\} \text{ ELSE } \text{forwarding}[i]]$   
 $\wedge \text{ceepochRecv}' = [\text{ceepochRecv} \text{ EXCEPT } ![i] = \text{IF } j \in \text{ceepochRecv}[i] \text{ THEN } \text{ceepochRecv}[i] \setminus \{j\} \text{ ELSE } \text{ceepochRecv}[i]]$   
 $\wedge \text{ackRecv}' = [\text{ackRecv} \text{ EXCEPT } ![i] = \text{IF } j \in \text{ackRecv}[i] \text{ THEN } \text{ackRecv}[i] \setminus \{j\} \text{ ELSE } \text{ackRecv}[i]]$   
 $\wedge \text{ackldRecv}' = [\text{ackldRecv} \text{ EXCEPT } ![i] = \text{IF } j \in \text{ackldRecv}[i] \text{ THEN } \text{ackldRecv}[i] \setminus \{j\} \text{ ELSE } \text{ackldRecv}[i]]$   
 $\wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i][j] = 0]$   
 $\wedge \text{committedCounter}' = [\text{committedCounter} \text{ EXCEPT } ![i][j] = 0]$   
 The action of follower  $j$ .  
 $\wedge \text{FollowerShutdown}(j)$   
 Clean input buffer.  
 $\wedge \text{Clean}(i, j)$   
 $\wedge$  UNCHANGED  $\langle \text{acceptedEpoch}, \text{history}, \text{commitIndex}, \text{currentCounter}, \text{sendCounter}, \text{committedIndex}, \text{committedCounter}, \text{tempVarsZ}, \text{ceepochSent}, \text{synced}, \text{verifyVarsZ} \rangle$

---

In phase  $f11$ , follower sends  $f.p$  to leader via *FOLLOWERINFO(CEPOCH)*.

$\text{FollowerSendFOLLOWERINFO}(i) \triangleq$   
 $\wedge \text{state}[i] = \text{FOLLOWING}$   
 $\wedge \text{zabState}[i] = \text{DISCOVERY}$   
 $\wedge \text{leaderAddr}[i] \neq \text{NullPoint}$   
 $\wedge \neg \text{ceepochSent}[i]$   
 $\wedge \text{Send}(i, \text{leaderAddr}[i], [\text{mtype} \mapsto \text{FOLLOWERINFO}, \text{mepoch} \mapsto \text{acceptedEpoch}[i]])$   
 $\wedge \text{ceepochSent}' = [\text{ceepochSent} \text{ EXCEPT } ![i] = \text{TRUE}]$   
 $\wedge$  UNCHANGED  $\langle \text{serverVarsZ}, \text{leaderVarsZ}, \text{electionVarsZ}, \text{tempVarsZ}, \text{leaderAddr}, \text{synced}, \text{verifyVarsZ} \rangle$



In phase *l11*, leader waits for receiving *FOLLOWERINFO* from a quorum, and then chooses a new epoch  $e'$  as its own epoch and broadcasts *LEADERINFO*.

$LeaderHandleFOLLOWERINFO(i, j) \triangleq$

- $\wedge state[i] = LEADING$
- $\wedge msgs[j][i] \neq \langle \rangle$
- $\wedge msgs[j][i][1].mtype = FOLLOWERINFO$
- $\wedge LET\ msg \triangleq msgs[j][i][1]$ 
  - IN  $\vee \wedge NullPoint \notin cepochRecv[i]$  1. has not broadcast *LEADERINFO* – modify *tempMaxEpoch*
    - $\wedge LET\ newEpoch \triangleq Maximum(\{tempMaxEpoch[i], msg.mepoch\})$
    - IN  $tempMaxEpoch' = [tempMaxEpoch\ EXCEPT\ ![i] = newEpoch]$
    - $\wedge Discard(j, i)$
  - $\vee \wedge NullPoint \in cepochRecv[i]$  2. has broadcast *LEADERINFO* – no need to handle the *msg*, just
    - $\wedge Reply(i, j, [mtype \mapsto LEADERINFO,$
    - $mepoch \mapsto acceptedEpoch[i]])$
    - $\wedge UNCHANGED\ tempMaxEpoch$
- $\wedge cepochRecv' = [ceepochRecv\ EXCEPT\ ![i] = IF\ j \in cepochRecv[i]\ THEN\ cepochRecv[i]$   
 $ELSE\ cepochRecv[i] \cup \{j\}]$
- $\wedge UNCHANGED\ \langle serverVarsZ, followerVarsZ, electionVarsZ, initialHistory, leadingVoteSet, learners, forwarding, ackIndex, currentCounter, sendCounter, committedIndex, committedCounter \rangle$

$LeaderDiscovery1(i) \triangleq$

- $\wedge state[i] = LEADING$
- $\wedge zabState[i] = DISCOVERY$
- $\wedge cepochRecv[i] \in Quorums$
- $\wedge acceptedEpoch' = [acceptedEpoch\ EXCEPT\ ![i] = tempMaxEpoch[i] + 1]$
- $\wedge cepochRecv' = [ceepochRecv\ EXCEPT\ ![i] = cepochRecv[i] \cup \{NullPoint\}]$
- $\wedge BroadcastLEADERINFO(i, [mtype \mapsto LEADERINFO,$   
 $mepoch \mapsto acceptedEpoch'[i]])$
- $\wedge UNCHANGED\ \langle state, currentEpoch, lastZxid, zabState, history, commitIndex, electionVarsZ, leadingVoteSet, learners, forwarding, ackIndex, currentCounter, sendCounter, committedIndex, committedCounter, tempVarsZ, followerVarsZ, verifyVarsZ, electionMsgs, idTable \rangle$

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*.

$FollowerHandleLEADERINFO(i, j) \triangleq$

- $\wedge state[i] = FOLLOWING$
- $\wedge msgs[j][i] \neq \langle \rangle$
- $\wedge msgs[j][i][1].mtype = LEADERINFO$
- $\wedge LET\ msg \triangleq msgs[j][i][1]$ 
  - $infoOk \triangleq j = leaderAddr[i]$
  - $epochOk \triangleq \wedge infoOk$
  - $\wedge msg.mepoch \geq acceptedEpoch[i]$
  - $correct \triangleq \wedge epochOk$
  - $\wedge zabState[i] = DISCOVERY$
- IN  $\wedge infoOk$

$\wedge \vee \wedge \text{epochOk}$  1. Normal case  
 $\wedge \vee \wedge \text{correct}$   
 $\wedge \text{acceptedEpoch}' = [\text{acceptedEpoch} \text{ EXCEPT } ![i] = \text{msg.mepoch}]$   
 $\wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{ACKEPOCH},$   
 $\text{mepoch} \mapsto \text{msg.mepoch},$   
 $\text{mlastEpoch} \mapsto \text{currentEpoch}[i],$   
 $\text{mlastZxid} \mapsto \text{lastZxid}[i]])$   
 $\wedge \text{cephochSent}' = [\text{cephochSent} \text{ EXCEPT } ![i] = \text{TRUE}]$   
 $\wedge \text{UNCHANGED } \text{inherentViolated}$   
 $\vee \wedge \neg \text{correct}$   
 $\wedge \text{PrintT}(\text{"Exception: Condition correct is false in FollowerHandleLEADERINFO("} \circ \text{To})$   
 $\wedge \text{inherentViolated}' = \text{TRUE}$   
 $\wedge \text{Discard}(j, i)$   
 $\wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{cephochSent} \rangle$   
 $\wedge \text{zabState}' = [\text{zabState} \text{ EXCEPT } ![i] = \text{IF } \text{zabState}[i] = \text{DISCOVERY} \text{ THEN } \text{SYNCHRONIZATION}$   
ELSE } \text{zabState}[i]]  
 $\wedge \text{UNCHANGED } \langle \text{varsL}, \text{leaderAddr} \rangle$   
 $\vee \wedge \neg \text{epochOk}$  2. Abnormal case - go back to election  
 $\wedge \text{FollowerShutdown}(i)$   
 $\wedge \text{Clean}(i, j)$   
 $\wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{cephochSent}, \text{inherentViolated} \rangle$   
 $\wedge \text{UNCHANGED } \langle \text{history}, \text{commitIndex}, \text{learners}, \text{cephochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{forwarding}, \text{ackIndex},$   
 $\text{committedCounter}, \text{tempVarsZ}, \text{syncd}, \text{proposalMsgsLog}, \text{epochLeader} \rangle$

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

$\text{subRECOVERYSYNC}(i, j) \triangleq$   
 $\text{LET } \text{canSync} \triangleq \wedge \text{state}[i] = \text{LEADING} \quad \wedge \text{zabState}[i] \neq \text{DISCOVERY} \quad \wedge j \in \text{learners}[i]$   
 $\wedge \text{state}[j] = \text{FOLLOWING} \wedge \text{zabState}[j] = \text{SYNCHRONIZATION} \wedge \text{leaderAddr}[j] =$   
 $\text{IN}$   
 $\vee \wedge \text{canSync}$   
 $\wedge \text{history}' = [\text{history} \text{ EXCEPT } ![j] = \text{history}[i]]$   
 $\wedge \text{lastZxid}' = [\text{lastZxid} \text{ EXCEPT } ![j] = \text{lastZxid}[i]]$   
 $\wedge \text{UpdateProposal}(j, \text{leaderAddr}[j], \text{lastZxid}'[j], \text{currentEpoch}[j])$   
 $\wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![j] = \text{commitIndex}[i]]$   
 $\wedge \text{syncd}' = [\text{syncd} \text{ EXCEPT } ![j] = \text{TRUE}]$   
 $\wedge \text{forwarding}' = [\text{forwarding} \text{ EXCEPT } ![i] = \text{forwarding}[i] \cup \{j\}]$  j will receive PROPOSAL and  
 $\wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i][j] = \text{Len}(\text{history}[i])]$   
 $\wedge \text{committedCounter}' = [\text{committedCounter} \text{ EXCEPT } ![i][j] = \text{Maximum}(\{\text{commitIndex}[i] - \text{Len}(\text{history}[i])\})]$   
 $\wedge \text{LET } \text{ms} \triangleq [\text{msource} \mapsto i, \text{mtype} \mapsto \text{"RECOVERYSYNC"}, \text{mepoch} \mapsto \text{acceptedEpoch}[i], \text{mproposals}$   
 $\text{IN } \text{proposalMsgsLog}' = \text{IF } \text{ms} \in \text{proposalMsgsLog} \text{ THEN } \text{proposalMsgsLog}$   
ELSE } \text{proposalMsgsLog} \cup \{\text{ms}\}  
 $\wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{NEWLEADER},$   
 $\text{mepoch} \mapsto \text{acceptedEpoch}[i],$   
 $\text{mlastZxid} \mapsto \text{lastZxid}[i]])$   
 $\vee \wedge \neg \text{canSync}$

$\wedge \text{Discard}(j, i)$   
 $\wedge \text{UNCHANGED } \langle \text{history}, \text{lastZxid}, \text{currentVote}, \text{commitIndex}, \text{synced}, \text{forwarding}, \text{ackIndex}, \text{commi} \rangle$

In phase *l12*, 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*.

$\text{LeaderHandleACKEPOCH}(i, j) \triangleq$   
 $\wedge \text{state}[i] = \text{LEADING}$   
 $\wedge \text{msgs}[j][i] \neq \langle \rangle$   
 $\wedge \text{msgs}[j][i][1].\text{mtype} = \text{ACKEPOCH}$   
 $\wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1]$   
 $\quad \text{infoOk} \triangleq \wedge j \in \text{learners}[i]$   
 $\quad \quad \wedge \text{acceptedEpoch}[i] = \text{msg.mepoch}$   
 $\quad \text{logOk} \triangleq \wedge \text{infoOk}$   $\text{logOk} = \text{TRUE}$  means leader is more up-to-date than follow  
 $\quad \quad \wedge \vee \text{currentEpoch}[i] > \text{msg.mlastEpoch}$   
 $\quad \quad \quad \vee \wedge \text{currentEpoch}[i] = \text{msg.mlastEpoch}$   
 $\quad \quad \quad \wedge \vee \text{lastZxid}[i][1] > \text{msg.mlastZxid}[1]$   
 $\quad \quad \quad \quad \vee \wedge \text{lastZxid}[i][1] = \text{msg.mlastZxid}[1]$   
 $\quad \quad \quad \quad \quad \wedge \text{lastZxid}[i][2] \geq \text{msg.mlastZxid}[2]$   
 $\quad \text{replyOk} \triangleq \wedge \text{infoOk}$   
 $\quad \quad \wedge \text{NullPoint} \in \text{ackeRecv}[i]$   
 $\text{IN } \wedge \text{infoOk}$   
 $\quad \wedge \vee \wedge \text{replyOk}$   
 $\quad \quad \wedge \text{subRECOVERYSYNC}(i, j)$   
 $\quad \quad \wedge \text{ackeRecv}' = [\text{ackeRecv} \text{ EXCEPT } ![i] = \text{IF } j \notin \text{ackeRecv}[i] \text{ THEN } \text{ackeRecv}[i] \cup \{j\}$   
 $\quad \quad \quad \text{ELSE } \text{ackeRecv}[i]]$   
 $\quad \quad \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{logicalClock}, \text{receiveVotes}, \text{outOfElection}, \text{recvQueue}$   
 $\quad \quad \quad \text{zabState}, \text{leaderAddr}, \text{learners} \rangle$   
 $\quad \vee \wedge \neg \text{replyOk}$   
 $\quad \wedge \vee \wedge \text{logOk}$   
 $\quad \quad \wedge \text{ackeRecv}' = [\text{ackeRecv} \text{ EXCEPT } ![i] = \text{IF } j \notin \text{ackeRecv}[i] \text{ THEN } \text{ackeRecv}[i] \cup \{j\}$   
 $\quad \quad \quad \text{ELSE } \text{ackeRecv}[i]]$   
 $\quad \quad \wedge \text{Discard}(j, i)$   
 $\quad \quad \wedge \text{UNCHANGED } \langle \text{varsL}, \text{zabState}, \text{leaderAddr}, \text{learners}, \text{forwarding} \rangle$   
 $\quad \vee \wedge \neg \text{logOk}$   $\text{go back to election}$   
 $\quad \quad \wedge \text{LeaderShutdown}(i)$   
 $\quad \quad \wedge \text{UNCHANGED } \text{ackeRecv}$   
 $\quad \quad \wedge \text{UNCHANGED } \langle \text{history}, \text{commitIndex}, \text{synced}, \text{forwarding}, \text{ackIndex}, \text{committedCounter},$   
 $\quad \quad \wedge \text{UNCHANGED } \langle \text{acceptedEpoch}, \text{ceepochRecv}, \text{ackldRecv}, \text{currentCounter}, \text{sendCounter}, \text{committedIndex} \rangle$

$\text{LeaderDiscovery2}(i) \triangleq$   
 $\wedge \text{state}[i] = \text{LEADING}$   
 $\wedge \text{zabState}[i] = \text{DISCOVERY}$   
 $\wedge \text{ackeRecv}[i] \in \text{Quorums}$   
 $\wedge \text{zabState}' = [\text{zabState} \text{ EXCEPT } ![i] = \text{SYNCHRONIZATION}]$   
 $\wedge \text{currentEpoch}' = [\text{currentEpoch} \text{ EXCEPT } ![i] = \text{acceptedEpoch}[i]]$

$\wedge \text{initialHistory}' = [\text{initialHistory} \text{ EXCEPT } ![i] = \text{history}[i]]$   
 $\wedge \text{ackRecv}' = [\text{ackRecv} \text{ EXCEPT } ![i] = \text{ackRecv}[i] \cup \{\text{NullPoint}\}]$   
 $\wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i][i] = \text{Len}(\text{history}[i])]$   
 $\wedge \text{UpdateProposal}(i, i, \text{lastZxid}[i], \text{currentEpoch}'[i])$   
 $\wedge \text{LET } \text{epoch} \triangleq \text{acceptedEpoch}[i]$   
 $\text{IN } \text{epochLeader}' = [\text{epochLeader} \text{ EXCEPT } ![epoch] = \text{epochLeader}[\text{epoch}] \cup \{i\}]$   
 $\wedge \text{UNCHANGED } \langle \text{state}, \text{lastZxid}, \text{acceptedEpoch}, \text{history}, \text{commitIndex}, \text{logicalClock}, \text{receiveVotes}, \text{outC}$   
 $\text{leadingVoteSet}, \text{learners}, \text{ceepochRecv}, \text{ackldRecv}, \text{forwarding}, \text{currentCounter}, \text{sendCo}$   
 $\text{tempMaxEpoch}, \text{followerVarsZ}, \text{proposalMsgsLog}, \text{inherentViolated}, \text{msgVarsZ}, \text{idTab}$

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

$\text{RECOVERYSYNC}(i, j) \triangleq$   
 $\wedge \text{state}[i] = \text{LEADING} \quad \wedge \text{zabState}[i] \neq \text{DISCOVERY} \quad \wedge j \in \text{learners}[i] \quad \wedge j \in \text{ackRecv}[i]$   
 $\wedge \text{state}[j] = \text{FOLLOWING} \wedge \text{zabState}[j] = \text{SYNCHRONIZATION} \wedge \text{leaderAddr}[j] = i \wedge \text{synced}[j] = \text{F}$   
 $\wedge \text{acceptedEpoch}[i] = \text{acceptedEpoch}[j] \setminus *$  This condition is unnecessary.  
 $\wedge \text{history}' = [\text{history} \text{ EXCEPT } ![j] = \text{history}[j]]$   
 $\wedge \text{lastZxid}' = [\text{lastZxid} \text{ EXCEPT } ![j] = \text{lastZxid}[j]]$   
 $\wedge \text{UpdateProposal}(j, \text{leaderAddr}[j], \text{lastZxid}'[j], \text{currentEpoch}[j])$   
 $\wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![j] = \text{commitIndex}[j]]$   
 $\wedge \text{synced}' = [\text{synced} \text{ EXCEPT } ![j] = \text{TRUE}]$   
 $\wedge \text{forwarding}' = [\text{forwarding} \text{ EXCEPT } ![i] = \text{forwarding}[i] \cup \{j\}]$   
 $\wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i][j] = \text{Len}(\text{history}[i])]$   
 $\wedge \text{committedCounter}' = [\text{committedCounter} \text{ EXCEPT } ![i][j] = \text{Maximum}(\{\text{commitIndex}[i] - \text{Len}(\text{initialHistory})\})]$   
 $\wedge \text{LET } ms \triangleq [\text{msource} \mapsto i, \text{mtype} \mapsto \text{"RECOVERYSYNC"}, \text{mepoch} \mapsto \text{acceptedEpoch}[i], \text{mproposals} \mapsto \text{proposalMsgsLog}]$   
 $\text{IN } \text{proposalMsgsLog}' = \text{IF } ms \in \text{proposalMsgsLog} \text{ THEN } \text{proposalMsgsLog}$   
 $\text{ELSE } \text{proposalMsgsLog} \cup \{ms\}$   
 $\wedge \text{Send}(i, j, [\text{mtype} \mapsto \text{NEWLEADER},$   
 $\text{mepoch} \mapsto \text{acceptedEpoch}[i],$   
 $\text{mlastZxid} \mapsto \text{lastZxid}[i]])$   
 $\wedge \text{UNCHANGED } \langle \text{state}, \text{zabState}, \text{acceptedEpoch}, \text{currentEpoch}, \text{logicalClock}, \text{receiveVotes}, \text{outOfElection},$   
 $\text{leadingVoteSet}, \text{learners}, \text{ceepochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{currentCounter}, \text{sendCounter},$   
 $\text{tempVarsZ}, \text{ceepochSent}, \text{leaderAddr}, \text{epochLeader}, \text{inherentViolated}, \text{electionMsgs}, \text{idTab}$

In phase *f21*, follower receives *NEWLEADER*. The follower updates its epoch and history, and sends back *ACK-LD* to leader.

$\text{FollowerHandleNEWLEADER}(i, j) \triangleq$   
 $\wedge \text{state}[i] = \text{FOLLOWING}$   
 $\wedge \text{msgs}[j][i] \neq \langle \rangle$   
 $\wedge \text{msgs}[j][i][1].\text{mtype} = \text{NEWLEADER}$   
 $\wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1]$   
 $\text{infoOk} \triangleq \wedge \text{leaderAddr}[i] = j$   
 $\wedge \text{acceptedEpoch}[i] = \text{msg.mepoch}$   
 $\text{correct} \triangleq \wedge \text{infoOk}$

$\wedge zabState[i] = SYNCHRONIZATION$   
 $\wedge synced[i]$   
 $\wedge ZxidEqual(lastZxid[i], msg.mlastZxid)$   
IN  $\wedge infoOk$   
 $\wedge currentEpoch' = [currentEpoch \text{ EXCEPT } ![i] = msg.mepoch]$   
 $\wedge UpdateProposal(i, j, lastZxid[i], currentEpoch'[i])$   
 $\wedge \vee \wedge correct$   
 $\wedge Reply(i, j, [mtype \mapsto ACKLD,$   
 $\quad mepoch \mapsto msg.mepoch])$   
 $\wedge UNCHANGED \text{ inherentViolated}$   
 $\vee \wedge \neg correct$   
 $\wedge PrintT(\text{"Exception: Condition correct is false in FollowerHandleNEWLEADER("} \circ ToString$   
 $\wedge inherentViolated' = TRUE$   
 $\wedge Discard(j, i)$   
 $\wedge UNCHANGED \langle state, lastZxid, zabState, acceptedEpoch, history, commitIndex, logicalClock, receiveV$   
 $\quad leaderVarsZ, tempVarsZ, followerVarsZ, proposalMsgsLog, epochLeader, electionMsgs$

In phase l22, leader receives *ACK-LD* from a quorum of followers, and sends *COMMIT-LD(UPTODATE)* to followers.

$LeaderHandleACKLD(i, j) \triangleq$   
 $\wedge state[i] = LEADING$   
 $\wedge msgs[j][i] \neq \langle \rangle$   
 $\wedge msgs[j][i][1].mtype = ACKLD$   
 $\wedge LET \text{ msg} \triangleq msgs[j][i][1]$   
 $\quad infoOk \triangleq \wedge acceptedEpoch[i] = msg.mepoch$   
 $\quad \wedge j \in learners[i]$   
 $\quad replyOk \triangleq \wedge infoOk$   
 $\quad \wedge NullPoint \in ackldRecv[i]$   
IN  $\wedge infoOk$   
 $\wedge \vee \wedge replyOk$   
 $\wedge Reply(i, j, [mtype \mapsto UPTODATE,$   
 $\quad mepoch \mapsto acceptedEpoch[i],$   
 $\quad mcommit \mapsto commitIndex[i]])$   
 $\wedge committedCounter' = [committedCounter \text{ EXCEPT } ![i][j] = Maximum(\{commitIndex[i]$   
 $\vee \wedge \neg replyOk$   
 $\wedge Discard(j, i)$   
 $\wedge UNCHANGED committedCounter$   
 $\wedge ackldRecv' = [ackldRecv \text{ EXCEPT } ![i] = \text{IF } j \notin ackldRecv[i] \text{ THEN } ackldRecv[i] \cup \{j\}$   
 $\quad \text{ELSE } ackldRecv[i}]$   
 $\wedge UNCHANGED \langle serverVarsZ, electionVarsZ, leadingVoteSet, learners, cepochRecv, ackeRecv, forward$   
 $\quad ackIndex, currentCounter, sendCounter, committedIndex, tempVarsZ, followerVarsZ$

$LeaderSync2(i) \triangleq$   
 $\wedge state[i] = LEADING$   
 $\wedge zabState[i] = SYNCHRONIZATION$   
 $\wedge ackldRecv[i] \in Quorums$

$$\begin{aligned}
& \wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{Len}(\text{history}[i])] \\
& \wedge \text{committedIndex}' = [\text{committedIndex} \text{ EXCEPT } ![i] = \text{Len}(\text{history}[i])] \\
& \wedge \text{zabState}' = [\text{zabState} \text{ EXCEPT } ![i] = \text{BROADCAST}] \\
& \wedge \text{currentCounter}' = [\text{currentCounter} \text{ EXCEPT } ![i] = 0] \\
& \wedge \text{sendCounter}' = [\text{sendCounter} \text{ EXCEPT } ![i] = 0] \\
& \wedge \text{ackldRecv}' = [\text{ackldRecv} \text{ EXCEPT } ![i] = \text{ackldRecv}[i] \cup \{\text{NullPoint}\}] \\
& \wedge \text{BroadcastUPTODATE}(i, [\text{mtype} \mapsto \text{UPTODATE}, \\
& \quad \text{mepoch} \mapsto \text{acceptedEpoch}[i], \\
& \quad \text{mcommit} \mapsto \text{Len}(\text{history}[i])]) \quad \text{In actual UPTODATE doesn't carry this info} \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{lastZxid}, \text{acceptedEpoch}, \text{history}, \text{electionVarsZ}, \text{leadingVoteS}, \\
& \quad \text{committedCounter}, \text{tempVarsZ}, \text{followerVarsZ}, \text{verifyVarsZ}, \text{electionMsgs}, \text{idTable} \rangle
\end{aligned}$$

$$\begin{aligned}
& \text{FollowerHandleUPTODATE}(i, j) \triangleq \\
& \quad \wedge \text{state}[i] = \text{FOLLOWING} \\
& \quad \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{UPTODATE} \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \quad \text{infoOk} \triangleq \wedge \text{leaderAddr}[i] = j \\
& \quad \quad \quad \wedge \text{acceptedEpoch}[i] = \text{msg.mepoch} \\
& \quad \quad \text{correct} \triangleq \wedge \text{infoOk} \\
& \quad \quad \quad \wedge \text{zabState}[i] = \text{SYNCHRONIZATION} \\
& \quad \quad \quad \wedge \text{currentEpoch}[i] = \text{msg.mepoch} \\
& \quad \text{IN } \wedge \text{infoOk} \\
& \quad \quad \wedge \vee \wedge \text{correct} \\
& \quad \quad \quad \wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{Maximum}(\{\text{commitIndex}[i], \text{msg.mcomm}\})] \\
& \quad \quad \quad \wedge \text{zabState}' = [\text{zabState} \text{ EXCEPT } ![i] = \text{BROADCAST}] \\
& \quad \quad \quad \wedge \text{UNCHANGED } \text{inherentViolated} \\
& \quad \quad \vee \wedge \neg \text{correct} \\
& \quad \quad \quad \wedge \text{PrintT}(\text{"Exception: Condition correct is false in FollowerHandleUPTODATE("} \circ \text{ToString}(\text{msg})) \\
& \quad \quad \quad \wedge \text{inherentViolated}' = \text{TRUE} \\
& \quad \quad \quad \wedge \text{UNCHANGED } \langle \text{commitIndex}, \text{zabState} \rangle \\
& \quad \wedge \text{Discard}(j, i) \\
& \quad \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{lastZxid}, \text{acceptedEpoch}, \text{history}, \text{electionVarsZ}, \text{leaderVarsZ}, \text{temp}, \\
& \quad \quad \text{proposalMsgsLog}, \text{epochLeader}, \text{electionMsgs}, \text{idTable} \rangle
\end{aligned}$$


---

In phase *l31*, leader receives client request and broadcasts *PROPOSAL*. Note: In production, any server in traffic can receive requests and forward it to leader if necessary. We choose to let leader be the sole one who can receive requests, to simplify spec and keep correctness at the same time.

$$\begin{aligned}
& \text{ClientRequest}(i, v) \triangleq \\
& \quad \wedge \text{state}[i] = \text{LEADING} \\
& \quad \wedge \text{zabState}[i] = \text{BROADCAST} \\
& \quad \wedge \text{currentCounter}' = [\text{currentCounter} \text{ EXCEPT } ![i] = \text{currentCounter}[i] + 1] \\
& \quad \wedge \text{LET } \text{newTransaction} \triangleq [\text{epoch} \mapsto \text{acceptedEpoch}[i], \\
& \quad \quad \text{counter} \mapsto \text{currentCounter}'[i],
\end{aligned}$$

$$\begin{aligned}
& \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{UNCHANGED } & \langle \text{state}, \text{currentEpoch}, \text{zabState}, \text{acceptedEpoch}, \text{commitIndex}, \text{logicalClock}, \text{receiveVotes}, \\
& \text{leadingVoteSet}, \text{learners}, \text{ceepochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{forwarding}, \text{sendCounter}, \\
& \text{tempVarsZ}, \text{followerVarsZ}, \text{verifyVarsZ}, \text{msgVarsZ}, \text{idTable} \rangle \\
\text{LeaderBroadcast1}(i) \triangleq & \\
& \wedge \text{state}[i] = \text{LEADING} \\
& \wedge \text{zabState}[i] = \text{BROADCAST} \\
& \wedge \text{sendCounter}[i] < \text{currentCounter}[i] \\
& \wedge \text{LET } \text{toBeSentCounter} \triangleq \text{sendCounter}[i] + 1 \\
& \quad \text{toBeSentIndex} \triangleq \text{Len}(\text{initialHistory}[i]) + \text{toBeSentCounter} \\
& \quad \text{toBeSentEntry} \triangleq \text{history}[i][\text{toBeSentIndex}] \\
\text{IN } & \wedge \text{Broadcast}(i, [\text{mtype} \mapsto \text{PROPOSAL}, \\
& \quad \text{mepoch} \mapsto \text{acceptedEpoch}[i], \\
& \quad \text{mproposal} \mapsto \text{toBeSentEntry}]) \\
& \wedge \text{sendCounter}' = [\text{sendCounter} \text{ EXCEPT } ![i] = \text{toBeSentCounter}] \\
& \wedge \text{LET } m \triangleq [\text{msource} \mapsto i, \text{mepoch} \mapsto \text{acceptedEpoch}[i], \text{mtype} \mapsto \text{PROPOSAL}, \text{mproposal} \mapsto \text{toBeSentEntry}] \\
& \quad \text{IN } \text{proposalMsgsLog}' = \text{proposalMsgsLog} \cup \{m\} \\
\wedge \text{UNCHANGED } & \langle \text{serverVarsZ}, \text{electionVarsZ}, \text{leadingVoteSet}, \text{learners}, \text{ceepochRecv}, \text{ackRecv}, \text{ackldRecv}, \\
& \text{committedIndex}, \text{committedCounter}, \text{tempVarsZ}, \text{followerVarsZ}, \text{epochLeader}, \text{inherent} \rangle
\end{aligned}$$

In phase *f31*, follower accepts proposal and append it to history.

$$\begin{aligned}
\text{FollowerHandlePROPOSAL}(i, j) \triangleq & \\
& \wedge \text{state}[i] = \text{FOLLOWING} \\
& \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \wedge \text{msgs}[j][i][1].\text{mtype} = \text{PROPOSAL} \\
& \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \text{infoOk} \triangleq \wedge \text{leaderAddr}[i] = j \\
& \quad \wedge \text{acceptedEpoch}[i] = \text{msg.mepoch} \\
& \quad \text{correct} \triangleq \wedge \text{infoOk} \\
& \quad \wedge \text{zabState}[i] \neq \text{DISCOVERY} \\
& \quad \wedge \text{syncd}[i] \\
& \quad \text{logOk} \triangleq \vee \wedge \text{msg.mproposal.counter} = 1 \quad \text{the first PROPOSAL in this epoch} \\
& \quad \wedge \vee \text{Len}(\text{history}[i]) = 0 \\
& \quad \vee \wedge \text{Len}(\text{history}[i]) > 0 \\
& \quad \wedge \text{history}[i][\text{Len}(\text{history}[i])].\text{epoch} < \text{msg.mepoch} \\
& \quad \vee \wedge \text{msg.mproposal.counter} > 1 \quad \text{not the first PROPOSAL in this epoch} \\
& \quad \wedge \text{Len}(\text{history}[i]) > 0 \\
& \quad \wedge \text{history}[i][\text{Len}(\text{history}[i])].\text{epoch} = \text{msg.mepoch} \\
& \quad \wedge \text{history}[i][\text{Len}(\text{history}[i])].\text{counter} = \text{msg.mproposal.counter} - 1 \\
\text{IN } & \wedge \text{infoOk}
\end{aligned}$$

$$\begin{aligned}
& \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: Condition logOk is false in FollowerHandlePROPOSAL("} \circ \text{ToString}(\text{state}, \text{currentEpoch}, \text{zabState}, \text{acceptedEpoch}, \text{commitIndex}, \text{logicalClock}, \text{receiveVotes}, \text{leaderVarsZ}, \text{tempVarsZ}, \text{followerVarsZ}, \text{proposalMsgsLog}, \text{epochLeader}, \text{electionMsgs}) \\
& \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: Condition correct is false in FollowerHandlePROPOSAL("} \circ \text{ToString}(\text{state}, \text{currentEpoch}, \text{zabState}, \text{acceptedEpoch}, \text{commitIndex}, \text{logicalClock}, \text{receiveVotes}, \text{leaderVarsZ}, \text{tempVarsZ}, \text{followerVarsZ}, \text{proposalMsgsLog}, \text{epochLeader}, \text{electionMsgs}) \\
& \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}, \text{leaderVarsZ}, \text{tempVarsZ}, \text{followerVarsZ}, \text{proposalMsgsLog}, \text{epochLeader}, \text{electionMsgs} \rangle
\end{aligned}$$

In phase *l32*, leader receives ack from a quorum of followers to a certain proposal, and commits the proposal.

$$\begin{aligned}
& \text{LeaderHandleACK}(i, j) \triangleq \\
& \quad \wedge \text{state}[i] = \text{LEADING} \\
& \quad \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{ACK} \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \quad \text{infoOk} \triangleq \wedge j \in \text{forwarding}[i] \\
& \quad \quad \quad \wedge \text{acceptedEpoch}[i] = \text{msg.mepoch} \\
& \quad \quad \text{correct} \triangleq \wedge \text{infoOk} \\
& \quad \quad \quad \wedge \text{zabState}[i] = \text{BROADCAST} \\
& \quad \quad \quad \wedge \text{sendCounter}[i] \geq \text{msg.mzxid}[2] \\
& \quad \quad \text{logOk} \triangleq \wedge \text{infoOk} \\
& \quad \quad \quad \wedge \text{ackIndex}[i][j] + 1 = \text{Len}(\text{initialHistory}[i]) + \text{msg.mzxid}[2] \\
& \text{IN} \quad \wedge \text{infoOk} \\
& \quad \wedge \vee \wedge \text{correct} \\
& \quad \quad \wedge \vee \wedge \text{logOk} \\
& \quad \quad \quad \wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i][j] = \text{ackIndex}[i][j] + 1] \\
& \quad \quad \vee \wedge \neg \text{logOk} \\
& \quad \quad \quad \wedge \text{PrintT}(\text{"Note: redundant ACK."}) \\
& \quad \quad \quad \wedge \text{UNCHANGED } \text{ackIndex} \\
& \quad \quad \wedge \text{UNCHANGED } \text{inherentViolated} \\
& \quad \vee \wedge \neg \text{correct} \\
& \quad \quad \wedge \text{PrintT}(\text{"Exception: Condition correct is false in FollowerHandleACK("} \circ \text{ToString}(i) \circ \text{"}, \text{state}, \text{currentEpoch}, \text{zabState}, \text{acceptedEpoch}, \text{commitIndex}, \text{logicalClock}, \text{receiveVotes}, \text{leaderVarsZ}, \text{tempVarsZ}, \text{followerVarsZ}, \text{proposalMsgsLog}, \text{epochLeader}, \text{electionMsgs})
\end{aligned}$$



$$\begin{aligned}
& \wedge \text{inherentViolated}' = \text{TRUE} \\
& \wedge \text{UNCHANGED } \text{ackIndex} \\
& \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED } \langle \text{serverVarsZ}, \text{electionVarsZ}, \text{leadingVoteSet}, \text{learners}, \text{epochRecv}, \text{ackRecv}, \text{ackldRe} \\
& \quad \text{committedIndex}, \text{committedCounter}, \text{tempVarsZ}, \text{followerVarsZ}, \text{proposalMsgsLog}, \text{ep} \\
\text{LeaderAdvanceCommit}(i) \triangleq & \\
& \wedge \text{state}[i] = \text{LEADING} \\
& \wedge \text{zabState}[i] = \text{BROADCAST} \\
& \wedge \text{commitIndex}[i] < \text{Len}(\text{history}[i]) \\
& \wedge \text{LET } \text{Agree}(\text{index}) \triangleq \{i\} \cup \{k \in (\text{Server} \setminus \{i\}) : \text{ackIndex}[i][k] \geq \text{index}\} \\
& \quad \text{agreeIndexes} \triangleq \{\text{index} \in (\text{commitIndex}[i] + 1) \dots \text{Len}(\text{history}[i]) : \text{Agree}(\text{index}) \in \text{Quorum}\} \\
& \quad \text{newCommitIndex} \triangleq \text{IF } \text{agreeIndexes} \neq \{\} \text{ THEN } \text{Maximum}(\text{agreeIndexes}) \\
& \quad \quad \quad \text{ELSE } \text{commitIndex}[i] \\
& \text{IN } \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{newCommitIndex}] \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{lastZxid}, \text{zabState}, \text{acceptedEpoch}, \text{history}, \text{electionVarsZ}, \text{leader} \\
& \quad \text{verifyVarsZ}, \text{msgVarsZ}, \text{idTable} \rangle \\
\text{LeaderBroadcast2}(i) \triangleq & \\
& \wedge \text{state}[i] = \text{LEADING} \\
& \wedge \text{zabState}[i] = \text{BROADCAST} \\
& \wedge \text{committedIndex}[i] < \text{commitIndex}[i] \\
& \wedge \text{Len}(\text{initialHistory}[i]) + \text{sendCounter}[i] > \text{committedIndex}[i] \\
& \wedge \text{LET } \text{newCommittedIndex} \triangleq \text{committedIndex}[i] + 1 \\
& \text{IN } \wedge \text{Broadcast}(i, [\text{mtype} \mapsto \text{COMMIT}, \\
& \quad \text{mepoch} \mapsto \text{acceptedEpoch}[i], \\
& \quad \text{mzxid} \mapsto \langle \text{history}[i][\text{newCommittedIndex}].\text{epoch}, \text{history}[i][\text{newCommittedIndex}] \\
& \wedge \text{committedIndex}' = [\text{committedIndex} \text{ EXCEPT } ![i] = \text{committedIndex}[i] + 1] \\
& \wedge \text{committedCounter}' = [\text{committedCounter} \text{ EXCEPT } ![i] = [v \in \text{Server} \mapsto \text{IF } \wedge v \in \text{forwarding} \\
& \quad \wedge \text{committedCounter}[v] < \text{committedCounter}[i] \text{ THEN } \text{history}[i][\text{newCommittedIndex}] \\
& \quad \text{ELSE } \text{committedCounter}[v]] \\
& \wedge \text{UNCHANGED } \langle \text{serverVarsZ}, \text{electionVarsZ}, \text{leadingVoteSet}, \text{learners}, \text{epochRecv}, \text{ackRecv}, \text{ackldRe} \\
& \quad \text{sendCounter}, \text{tempVarsZ}, \text{followerVarsZ}, \text{verifyVarsZ}, \text{electionMsgs}, \text{idTable} \rangle \\
\text{In phase } f32, \text{ follower receives } \text{COMMIT} \text{ and commits transaction.} \\
\text{FollowerHandleCOMMIT}(i, j) \triangleq & \\
& \wedge \text{state}[i] = \text{FOLLOWING} \\
& \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \wedge \text{msgs}[j][i][1].\text{mtype} = \text{COMMIT} \\
& \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \text{infoOk} \triangleq \wedge \text{leaderAddr}[i] = j \\
& \quad \quad \wedge \text{acceptedEpoch}[i] = \text{msg.mepoch} \\
& \text{correct} \triangleq \wedge \text{infoOk} \\
& \quad \wedge \text{zabState}[i] \neq \text{DISCOVERY} \\
& \quad \wedge \text{syncd}[i]
\end{aligned}$$

$$\begin{aligned}
mindex &\triangleq \text{IF } Len(history[i]) = 0 \text{ THEN } -1 \\
&\quad \text{ELSE IF } \exists idx \in 1 \dots Len(history[i]) : PZxidEqual(history[i][idx], msg.mzxid) \\
&\quad \quad \text{THEN CHOOSE } idx \in 1 \dots Len(history[i]) : PZxidEqual(history[i][idx], msg.mzxid) \\
&\quad \quad \text{ELSE } -1 \\
logOk &\triangleq mindex > 0 \\
latest &\triangleq commitIndex[i] + 1 = mindex \\
\text{IN } &\wedge infoOk \\
&\wedge \vee \wedge correct \\
&\quad \wedge \vee \wedge logOk \\
&\quad \quad \wedge \vee \wedge latest \\
&\quad \quad \quad \wedge commitIndex' = [commitIndex \text{ EXCEPT } !i] = commitIndex[i] + 1 \\
&\quad \quad \quad \wedge \text{UNCHANGED } inherentViolated \\
&\quad \vee \wedge \neg latest \\
&\quad \quad \wedge PrintT(\text{"Note: Condition latest is false in FollowerHandleCOMMIT("} \circ ToString(i)) \\
&\quad \quad \wedge inherentViolated' = \text{TRUE} \\
&\quad \quad \wedge \text{UNCHANGED } commitIndex \\
&\vee \wedge \neg logOk \\
&\quad \wedge PrintT(\text{"Exception: Condition logOk is false in FollowerHandleCOMMIT("} \circ ToString(i)) \\
&\quad \quad \wedge inherentViolated' = \text{TRUE} \\
&\quad \quad \wedge \text{UNCHANGED } commitIndex \\
&\vee \wedge \neg correct \\
&\quad \wedge PrintT(\text{"Exception: Condition correct is false in FollowerHandleCOMMIT("} \circ ToString(i)) \\
&\quad \quad \wedge inherentViolated' = \text{TRUE} \\
&\quad \quad \wedge \text{UNCHANGED } commitIndex \\
&\wedge Discard(j, i) \\
&\wedge \text{UNCHANGED } \langle state, currentEpoch, lastZxid, zabState, acceptedEpoch, history, electionVarsZ, leader, \\
&\quad proposalMsgsLog, epochLeader, electionMsgs, idTable \rangle
\end{aligned}$$


---

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

$$\begin{aligned}
FilterNonexistentMessage(i) &\triangleq \\
&\wedge \exists j \in Server \setminus \{i\} : \wedge msgs[j][i] \neq \langle \rangle \\
&\quad \wedge \text{LET } msg \triangleq msgs[j][i][1] \\
&\quad \text{IN} \\
&\quad \vee \wedge state[i] = LEADING \\
&\quad \quad \wedge \text{LET } infoOk \triangleq \wedge j \in learners[i] \\
&\quad \quad \quad \wedge acceptedEpoch[i] = msg.mepoch \\
&\quad \text{IN} \\
&\quad \quad \vee msg.mtype = LEADERINFO \\
&\quad \quad \vee msg.mtype = NEWLEADER \\
&\quad \quad \vee msg.mtype = UPTODATE \\
&\quad \quad \vee msg.mtype = PROPOSAL \\
&\quad \quad \vee msg.mtype = COMMIT \\
&\quad \quad \vee \wedge j \notin learners[i]
\end{aligned}$$

$$\begin{aligned}
& \wedge msg.mtype = FOLLOWERINFO \\
& \vee \wedge \neg infoOk \\
& \wedge \vee msg.mtype = ACKEPOCH \\
& \quad \vee msg.mtype = ACKLD \\
& \quad \vee msg.mtype = ACK \\
& \vee \wedge state[i] = FOLLOWING \\
& \wedge LET infoOk \triangleq \wedge j = leaderAddr[i] \\
& \quad \wedge acceptedEpoch[i] = msg.mepoch \\
& IN \\
& \vee msg.mtype = FOLLOWERINFO \\
& \vee msg.mtype = ACKEPOCH \\
& \vee msg.mtype = ACKLD \\
& \vee msg.mtype = ACK \\
& \vee \wedge j \neq leaderAddr[i] \\
& \quad \wedge msg.mtype = LEADERINFO \\
& \vee \wedge \neg infoOk \\
& \quad \wedge \vee msg.mtype = NEWLEADER \\
& \quad \vee msg.mtype = UPTODATE \\
& \quad \vee msg.mtype = PROPOSAL \\
& \quad \vee msg.mtype = COMMIT \\
& \vee state[i] = LOOKING \\
& \wedge Discard(j, i) \\
& \wedge inherentViolated' = TRUE \\
& \wedge UNCHANGED \langle serverVarsZ, electionVarsZ, leaderVarsZ, tempVarsZ, followerVarsZ, proposalMsgsL
\end{aligned}$$


---

Defines how the variables may transition.

$NextZ \triangleq$

*FLE* modlue

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

Some conditions like failure, network delay

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

Zab module - Discovery and Synchronization part

$\vee \exists i, j \in Server : EstablishConnection(i, j)$   
 $\vee \exists i \in Server : FollowerSendFOLLOWERINFO(i)$   
 $\vee \exists i, j \in Server : LeaderHandleFOLLOWERINFO(i, j)$   
 $\vee \exists i \in Server : LeaderDiscovery1(i)$   
 $\vee \exists i, j \in Server : FollowerHandleLEADERINFO(i, j)$   
 $\vee \exists i, j \in Server : LeaderHandleACKEPOCH(i, j)$

$\forall \exists i \in \text{Server} : \text{LeaderDiscovery2}(i)$   
 $\forall \exists i, j \in \text{Server} : \text{RECOVERYSYNC}(i, j)$   
 $\forall \exists i, j \in \text{Server} : \text{FollowerHandleNEWLEADER}(i, j)$   
 $\forall \exists i, j \in \text{Server} : \text{LeaderHandleACKLD}(i, j)$   
 $\forall \exists i \in \text{Server} : \text{LeaderSync2}(i)$   
 $\forall \exists i, j \in \text{Server} : \text{FollowerHandleUPTODATE}(i, j)$   
 Zab module – Broadcast part  
 $\forall \exists i \in \text{Server}, v \in \text{Value} : \text{ClientRequest}(i, v)$   
 $\forall \exists i \in \text{Server} : \text{LeaderBroadcast1}(i)$   
 $\forall \exists i, j \in \text{Server} : \text{FollowerHandlePROPOSAL}(i, j)$   
 $\forall \exists i, j \in \text{Server} : \text{LeaderHandleACK}(i, j)$   
 $\forall \exists i \in \text{Server} : \text{LeaderAdvanceCommit}(i)$   
 $\forall \exists i \in \text{Server} : \text{LeaderBroadcast2}(i)$   
 $\forall \exists i, j \in \text{Server} : \text{FollowerHandleCOMMIT}(i, j)$   
 An action used to judge whether there are redundant messages in network  
 $\forall \exists i \in \text{Server} : \text{FilterNonexistentMessage}(i)$

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

---

Define safety properties of Zab 1.0 protocol.

$\text{ShouldNotBeTriggered} \triangleq \text{inherentViolated} = \text{FALSE}$

There is most one established leader for a certain epoch.

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

$\wedge \text{state}[i] = \text{LEADING} \wedge \text{zabState}[i] \in \{\text{SYNCHRONIZATION}, \text{BROADCAST}\}$   
 $\wedge \text{state}[j] = \text{LEADING} \wedge \text{zabState}[j] \in \{\text{SYNCHRONIZATION}, \text{BROADCAST}\}$   
 $\wedge \text{acceptedEpoch}[i] = \text{acceptedEpoch}[j]$   
 $\Rightarrow i = j$

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

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

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

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

IN  $\vee \text{smaller} = 0$

$\vee \wedge \text{smaller} > 0$

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

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

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

$\wedge \text{state}[i] = \text{FOLLOWING}$

$\wedge \text{commitIndex}[i] > 0$

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

$\vee \wedge \text{msg.mtype} = \text{PROPOSAL}$

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

$$\begin{aligned} & \vee \wedge msg.mtype = \text{"RECOVERYSYNC"} \\ & \wedge \exists tindex \in 1 \dots Len(msg.mproposals) : TransactionEqual(msg.mproposals[tindex], h \end{aligned}$$

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

$$\begin{aligned} Agreement &\triangleq \forall i, j \in Server : \\ & \wedge state[i] = FOLLOWING \wedge commitIndex[i] > 0 \\ & \wedge state[j] = FOLLOWING \wedge commitIndex[j] > 0 \\ & \Rightarrow \\ & \forall index1 \in 1 \dots commitIndex[i], index2 \in 1 \dots commitIndex[j] : \\ & \quad \vee \exists indexj \in 1 \dots commitIndex[j] : \\ & \quad \quad TransactionEqual(history[j][indexj], history[i][index1]) \\ & \quad \vee \exists indexi \in 1 \dots commitIndex[i] : \\ & \quad \quad TransactionEqual(history[i][indexi], history[j][index2]) \end{aligned}$$

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

$$\begin{aligned} TotalOrder &\triangleq \forall i, j \in Server : commitIndex[i] \geq 2 \wedge commitIndex[j] \geq 2 \\ & \Rightarrow \forall indexi1 \in 1 \dots (commitIndex[i] - 1) : \forall indexi2 \in (indexi1 + 1) \dots commitIndex[i] : \\ & \quad LET logOk \triangleq \exists index \in 1 \dots commitIndex[j] : TransactionEqual(history[i][indexi2], history[j][index]) \\ & \quad IN \quad \vee \neg logOk \\ & \quad \vee \wedge logOk \\ & \quad \quad \wedge \exists indexj2 \in 1 \dots commitIndex[j] : \\ & \quad \quad \quad \wedge TransactionEqual(history[i][indexi2], history[j][indexj2]) \\ & \quad \quad \quad \wedge \exists indexj1 \in 1 \dots (indexj2 - 1) : TransactionEqual(history[i][indexi2], history[j][indexj1]) \end{aligned}$$

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

$$\begin{aligned} LocalPrimaryOrder &\triangleq LET mset(i, e) \triangleq \{msg \in proposalMsgsLog : \wedge msg.mtype = PROPOSAL \\ & \quad \wedge msg.msource = i \\ & \quad \wedge msg.mepoch = e\} \\ & \quad mentries(i, e) \triangleq \{msg.mproposal : msg \in mset(i, e)\} \\ IN \quad & \forall i \in Server : \forall e \in 1 \dots currentEpoch[i] : \\ & \quad \vee Cardinality(mentries(i, e)) < 2 \\ & \quad \vee \wedge Cardinality(mentries(i, e)) \geq 2 \\ & \quad \quad \wedge \exists tsc1, tsc2 \in mentries(i, e) : \\ & \quad \quad \vee TransactionEqual(tsc1, tsc2) \\ & \quad \quad \vee \wedge \neg TransactionEqual(tsc1, tsc2) \\ & \quad \quad \quad \wedge LET tscPre \triangleq IF TransactionPrecede(tsc1, tsc2) THEN tsc1 ELSE tsc2 \\ & \quad \quad \quad \quad tscNext \triangleq IF TransactionPrecede(tsc1, tsc2) THEN tsc2 ELSE tsc1 \\ & \quad \quad IN \quad \forall j \in Server : \wedge commitIndex[j] \geq 2 \\ & \quad \quad \quad \wedge \exists index \in 1 \dots commitIndex[j] : TransactionEqual(history[j][index], tscPre) \\ & \quad \quad \quad \Rightarrow \exists index2 \in 1 \dots commitIndex[j] : \\ & \quad \quad \quad \quad \wedge TransactionEqual(history[j][index2], tscNext) \\ & \quad \quad \quad \quad \wedge index2 > 1 \\ & \quad \quad \quad \quad \wedge \exists index1 \in 1 \dots (index2 - 1) : TransactionEqual(history[j][index1], tscNext) \end{aligned}$$

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

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

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

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

---

\ \* Modification History  
\ \* Last modified *Thu Jul 15 16:10:01 CST 2021* by Dell  
\ \* Created *Tue Jun 29 22:13:02 CST 2021* by Dell