

The leader's epoch or the last new epoch proposal the follower acknowledged
(namely epoch of the last *NEWEPOCH* accepted, $f.p$ in paper).

VARIABLE *currentEpoch*

The last new leader proposal the follower acknowledged
(namely epoch of the last *NEWLEADER* accepted, $f.a$ in paper).

VARIABLE *leaderEpoch*

The identifier of the leader for followers.

VARIABLE *leaderOracle*

The history of servers as the sequence of transactions.

VARIABLE *history*

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 servers which the leader think follow itself (Q in paper).

VARIABLE *cluster*

The set of followers who has successfully sent *CEPOCH* to pleader in pleader.

VARIABLE *ceepochRecv*

The set of followers who has successfully sent *ACK-E* to pleader in pleader.

VARIABLE *ackeRecv*

The set of followers who has successfully sent *ACK-LD* to pleader in pleader.

VARIABLE *ackldRecv*

$ackIndex[i][j]$ means leader i has received how many *ACK* messages from follower j .
So $ackIndex[i][i]$ is not used.

VARIABLE *ackIndex*

$currentCounter[i]$ means the count of transactions client requests leader.

VARIABLE *currentCounter*

$sendCounter[i]$ means the count of transactions leader has broadcast.

VARIABLE *sendCounter*

$initialHistory[i]$ means the initial history of leader i in epoch $currentEpoch[i]$.

VARIABLE *initialHistory*

$commitIndex[i]$ means leader/follower i should commit how many proposals and sent *COMMIT* messages.

It should be more formal to add variable *applyIndex/deliverIndex* to represent the prefix entries of the history that has applied to state machine, but we can tolerate that $applyIndex(deliverIndex\ here) = commitIndex$.

This does not violate correctness. (*commitIndex* increases monotonically before restarting)

VARIABLE *commitIndex*

$commitIndex[i]$ means leader i has committed how many proposals and sent *COMMIT* messages.
 VARIABLE *committedIndex*

Helper matrix for follower to stop sending *CEPOCH* to pleader in followers.
 Because *CEPOCH* is the sole message which follower actively sends to pleader.
 VARIABLE *ceepochSent*

the maximum epoch in *CEPOCH* pleader received from followers.
 VARIABLE *tempMaxEpoch*

the maximum *leaderEpoch* and most up-to-date history in *ACKE* pleader received from followers.
 VARIABLE *tempMaxLastEpoch*

Because pleader updates state and broadcasts *NEWLEADER* when it receives *ACKE* from a quorum of followers,
 and *initialHistory* is determined. But *tempInitialHistory* may change when receiving other *ACKEs* after entering into *phase2*.
 So it is necessary to split *initialHistory* with *tempInitialHistory*.
 VARIABLE *tempInitialHistory*

the set of all broadcast messages whose type is proposal that any leader has sent, only used in verifying properties.
 So the variable will only be changed in transition *LeaderBroadcast1*.
 VARIABLE *proposalMsgsLog*

Helper set for server who restarts to collect which servers has responded to it.
 VARIABLE *recoveryRespRecv*

the maximum epoch and corresponding *leaderOracle* in *RECOVERYRESPONSE* from followers.
 VARIABLE *recoveryMaxEpoch*

VARIABLE *recoveryMEOracle*

Helper variable for server after restart to stop sending *RECOVERYREQUEST* continuously.
 VARIABLE *recoverySent*

variables that uniquely used for constraining state space in model checking
 VARIABLES *electionNum*, the round of leader election, not equal to $Maximum\{currentEpoch[i] : i \in Server\}$,
 because *currentEpoch* will increase only when follower receives *NEWEPOCH*,
 and it is common that some round of election ends without leader broadcasting *NEWEPOCH*
 or follower receiving *NEWEPOCH*.
totalRestartNum the number of restart from all servers, also as a global variable.

Persistent state of a server: history, *currentEpoch*, *leaderEpoch*
 $serverVars \triangleq \langle state, currentEpoch, leaderEpoch, leaderOracle, history, commitIndex \rangle$
 $leaderVars \triangleq \langle cluster, cepochRecv, ackRecv, ackldRecv, ackIndex, currentCounter, sendCounter, initialHistory \rangle$
 $tempVars \triangleq \langle tempMaxEpoch, tempMaxLastEpoch, tempInitialHistory \rangle$
 $recoveryVars \triangleq \langle recoveryRespRecv, recoveryMaxEpoch, recoveryMEOracle, recoverySent \rangle$
 $testVars \triangleq \langle electionNum, totalRestartNum \rangle$
 $vars \triangleq \langle serverVars, msgs, leaderVars, tempVars, recoveryVars, cepochSent, proposalMsgsLog, testVars \rangle$

$LastZxid(his) \triangleq \text{IF } Len(his) > 0 \text{ THEN } \langle his[Len(his)].epoch, his[Len(his)].counter \rangle$
 $\text{ELSE } \langle -1, -1 \rangle$

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

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

$Send2(i, j, m1, m2) \triangleq msgs' = [msgs \text{ EXCEPT } ![i][j] = Append(Append(msgs[i][j], m1), m2)]$

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

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

Leader/Pleader broadcasts a message to all other servers in Q

$Broadcast(i, m) \triangleq msgs' = [ii \in Server \mapsto [ij \in Server \mapsto \text{IF } \wedge ii = i$
 $\wedge ij \neq i$
 $\wedge ij \in cluster[i] \text{ THEN } Append(msgs[ii][ij], m)$
 $\text{ELSE } msgs[ii][ij]]]$

$BroadcastToAll(i, m) \triangleq msgs' = [ii \in Server \mapsto [ij \in Server \mapsto \text{IF } \wedge ii = i \wedge ij \neq i \text{ THEN } Append(msgs[ii][ij], m)$
 $\text{ELSE } msgs[ii][ij]]]$

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

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

$Reply2(i, j, m1, m2) \triangleq msgs' = [msgs \text{ EXCEPT } ![j][i] = Tail(msgs[j][i]),$
 $![i][j] = Append(Append(msgs[i][j], m1), m2)]$

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

Define initial values for all variables

$Init \triangleq \wedge state = [s \in Server \mapsto Follower]$
 $\wedge currentEpoch = [s \in Server \mapsto 0]$
 $\wedge leaderEpoch = [s \in Server \mapsto 0]$
 $\wedge leaderOracle = [s \in Server \mapsto NullPoint]$
 $\wedge history = [s \in Server \mapsto \langle \rangle]$
 $\wedge msgs = [i \in Server \mapsto [j \in Server \mapsto \langle \rangle]]$
 $\wedge cluster = [i \in Server \mapsto \{\}]$
 $\wedge epochRecv = [s \in Server \mapsto \{\}]$
 $\wedge ackRecv = [s \in Server \mapsto \{\}]$
 $\wedge ackldRecv = [s \in Server \mapsto \{\}]$
 $\wedge ackIndex = [i \in Server \mapsto [j \in Server \mapsto 0]]$
 $\wedge currentCounter = [s \in Server \mapsto 0]$
 $\wedge sendCounter = [s \in Server \mapsto 0]$
 $\wedge commitIndex = [s \in Server \mapsto 0]$
 $\wedge committedIndex = [s \in Server \mapsto 0]$

$$\begin{aligned}
\wedge \text{initialHistory} &= [s \in \text{Server} \mapsto \langle \rangle] \\
\wedge \text{epochSent} &= [s \in \text{Server} \mapsto \text{FALSE}] \\
\wedge \text{tempMaxEpoch} &= [s \in \text{Server} \mapsto 0] \\
\wedge \text{tempMaxLastEpoch} &= [s \in \text{Server} \mapsto 0] \\
\wedge \text{tempInitialHistory} &= [s \in \text{Server} \mapsto \langle \rangle] \\
\wedge \text{recoveryRespRecv} &= [s \in \text{Server} \mapsto \{\}] \\
\wedge \text{recoveryMaxEpoch} &= [s \in \text{Server} \mapsto 0] \\
\wedge \text{recoveryMEOracle} &= [s \in \text{Server} \mapsto \text{NullPoint}] \\
\wedge \text{recoverySent} &= [s \in \text{Server} \mapsto \text{FALSE}] \\
\wedge \text{proposalMsgsLog} &= \{\} \\
\wedge \text{electionNum} &= 0 \\
\wedge \text{totalRestartNum} &= 0
\end{aligned}$$

A server becomes pleader and a quorum servers knows that.

$\text{Election}(i, Q) \triangleq$

test restrictions

$\wedge \text{electionNum} < \text{MaxElectionNum}$

$\wedge i \in Q$

$\wedge \text{electionNum}' = \text{electionNum} + 1$

$\wedge \text{state}' = [s \in \text{Server} \mapsto \text{IF } s = i \text{ THEN } \text{ProspectiveLeader} \\ \text{ELSE IF } s \in Q \text{ THEN } \text{Follower} \\ \text{ELSE } \text{state}[s]]$

$\wedge \text{cluster}' = [\text{cluster} \text{ EXCEPT } ![i] = Q] \text{ cluster is first initialized in election, not phase1.}$

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

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

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

$\wedge \text{ackIndex}' = [ii \in \text{Server} \mapsto [ij \in \text{Server} \mapsto \\ \text{IF } ii = i \text{ THEN } 0$

$\text{ELSE } \text{ackIndex}[ii][ij]]]$

$\wedge \text{committedIndex}' = [\text{committedIndex} \text{ EXCEPT } ![i] = 0]$

$\wedge \text{initialHistory}' = [\text{initialHistory} \text{ EXCEPT } ![i] = \langle \rangle]$

$\wedge \text{tempMaxEpoch}' = [\text{tempMaxEpoch} \text{ EXCEPT } ![i] = \text{currentEpoch}[i]]$

$\wedge \text{tempMaxLastEpoch}' = [\text{tempMaxLastEpoch} \text{ EXCEPT } ![i] = \text{currentEpoch}[i]]$

$\wedge \text{tempInitialHistory}' = [\text{tempInitialHistory} \text{ EXCEPT } ![i] = \text{history}[i]]$

$\wedge \text{leaderOracle}' = [s \in \text{Server} \mapsto \text{IF } s \in Q \text{ THEN } i \\ \text{ELSE } \text{leaderOracle}[s]]$

$\wedge \text{leaderEpoch}' = [s \in \text{Server} \mapsto \text{IF } s \in Q \text{ THEN } \text{currentEpoch}[s] \\ \text{ELSE } \text{leaderEpoch}[s]]$

$\wedge \text{epochSent}' = [s \in \text{Server} \mapsto \text{IF } s \in Q \text{ THEN } \text{FALSE} \\ \text{ELSE } \text{epochSent}[s]]$

$\wedge \text{msgs}' = [ii \in \text{Server} \mapsto [ij \in \text{Server} \mapsto \\ \text{IF } ii \in Q \vee ij \in Q \text{ THEN } \langle \rangle \\ \text{ELSE } \text{msgs}[ii][ij]]]$

$\wedge \text{UNCHANGED } \langle \text{currentEpoch}, \text{history}, \text{commitIndex}, \text{currentCounter}, \text{sendCounter}, \text{proposalMsgsLog} \rangle$

The action should be triggered once at the beginning.

Because we abstract the part of leader election, we can use global variables in this action.

$$\begin{aligned} \text{InitialElection}(i, Q) &\triangleq \\ &\wedge \forall s \in \text{Server} : \text{state}[s] = \text{Follower} \wedge \text{leaderOracle}[s] = \text{NullPoint} \\ &\wedge \text{Election}(i, Q) \\ &\wedge \text{UNCHANGED} \langle \text{currentEpoch}, \text{history}, \text{commitIndex}, \text{currentCounter}, \text{sendCounter}, \text{recoveryVars}, \text{pro} \rangle \end{aligned}$$

The leader finds timeout with another follower.

$$\begin{aligned} \text{LeaderTimeout}(i, j) &\triangleq \\ &\wedge \text{state}[i] \neq \text{Follower} \\ &\wedge j \neq i \\ &\wedge j \in \text{cluster}[i] \\ &\wedge \text{LET } \text{newCluster} \triangleq \text{cluster}[i] \setminus \{j\} \\ &\quad \text{IN } \wedge \vee \wedge \text{newCluster} \in \text{Quorums} \\ &\quad \wedge \text{cluster}' = [\text{cluster} \text{ EXCEPT } ![i] = \text{newCluster}] \\ &\quad \wedge \text{clean}(i, j) \\ &\quad \wedge \text{UNCHANGED} \langle \text{state}, \text{cepcvRecv}, \text{ackRecv}, \text{ackldRecv}, \text{ackIndex}, \text{committedIndex}, \text{initialHistory}, \\ &\quad \quad \text{tempMaxEpoch}, \text{tempMaxLastEpoch}, \text{tempInitialHistory}, \text{leaderOracle}, \text{leaderOracle}, \text{pro} \rangle \\ &\quad \vee \wedge \text{newCluster} \notin \text{Quorums} \\ &\quad \wedge \text{LET } Q \triangleq \text{CHOOSE } q \in \text{Quorums} : i \in q \\ &\quad \quad v \triangleq \text{CHOOSE } s \in Q : \text{TRUE} \\ &\quad \quad \text{IN } \text{Election}(v, Q) \\ &\quad \exists Q \in \text{Quorums} : \wedge i \in Q \\ &\quad \quad \wedge \exists v \in Q : \text{Election}(v, Q) \\ &\wedge \text{UNCHANGED} \langle \text{currentEpoch}, \text{history}, \text{commitIndex}, \text{currentCounter}, \text{sendCounter}, \text{recoveryVars}, \text{pro} \rangle \end{aligned}$$

A follower finds timeout with the leader.

$$\begin{aligned} \text{FollowerTimeout}(i) &\triangleq \\ &\wedge \text{state}[i] = \text{Follower} \\ &\wedge \text{leaderOracle}[i] \neq \text{NullPoint} \\ &\wedge \exists Q \in \text{Quorums} : \wedge i \in Q \\ &\quad \wedge \exists v \in Q : \text{Election}(v, Q) \\ &\wedge \text{UNCHANGED} \langle \text{currentEpoch}, \text{history}, \text{commitIndex}, \text{currentCounter}, \text{sendCounter}, \text{recoveryVars}, \text{pro} \rangle \end{aligned}$$

A server halts and restarts.

Like Recovery protocol in View-stamped Replication, we let a server join in cluster by broadcast recovery and wait until receiving responses from a quorum of servers.

$$\begin{aligned} \text{Restart}(i) &\triangleq \\ &\text{test restrictions} \\ &\wedge \text{totalRestartNum} < \text{MaxTotalRestartNum} \\ &\wedge \text{totalRestartNum}' = \text{totalRestartNum} + 1 \\ &\wedge \text{state}' = [\text{state} \text{ EXCEPT } ![i] = \text{Follower}] \\ &\wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = \text{NullPoint}] \\ &\wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = 0] \\ &\wedge \text{cepcvSent}' = [\text{cepcvSent} \text{ EXCEPT } ![i] = \text{FALSE}] \end{aligned}$$

$$\begin{aligned}
& \wedge \text{msgs}' &= [ii \in \text{Server} \mapsto [ij \in \text{Server} \mapsto \text{IF } ij = i \text{ THEN } \langle \rangle \\
& & \hspace{15em} \text{ELSE } \text{msgs}[ii][ij]]] \\
& \wedge \text{recoverySent}' &= [\text{recoverySent} \text{ EXCEPT } ![i] = \text{FALSE}] \\
& \wedge \text{UNCHANGED} &\langle \text{currentEpoch}, \text{leaderEpoch}, \text{history}, \text{leaderVars}, \text{tempVars}, \\
& & \text{recoveryRespRecv}, \text{recoveryMaxEpoch}, \text{recoveryMEOracle}, \text{proposalMsgsLog}, \text{election} \rangle \\
\text{RecoveryAfterRestart}(i) &\triangleq \\
& \text{test restrictions} \\
& \wedge \text{totalRestartNum} < \text{MaxTotalRestartNum} \\
& \wedge \text{state}[i] = \text{Follower} \\
& \wedge \text{leaderOracle}[i] = \text{NullPoint} \\
& \wedge \neg \text{recoverySent}[i] \\
& \wedge \text{recoveryRespRecv}' = [\text{recoveryRespRecv} \text{ EXCEPT } ![i] = \{\}] \\
& \wedge \text{recoveryMaxEpoch}' = [\text{recoveryMaxEpoch} \text{ EXCEPT } ![i] = \text{currentEpoch}[i]] \\
& \wedge \text{recoveryMEOracle}' = [\text{recoveryMEOracle} \text{ EXCEPT } ![i] = \text{NullPoint}] \\
& \wedge \text{recoverySent}' = [\text{recoverySent} \text{ EXCEPT } ![i] = \text{TRUE}] \\
& \wedge \text{BroadcastToAll}(i, [\text{mtype} \mapsto \text{RECOVERYREQUEST}]) \\
& \wedge \text{UNCHANGED} \langle \text{serverVars}, \text{leaderVars}, \text{tempVars}, \text{epochSent}, \text{proposalMsgsLog}, \text{testVars} \rangle \\
\text{HandleRecoveryRequest}(i, j) &\triangleq \\
& \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \wedge \text{msgs}[j][i][1].\text{mtype} = \text{RECOVERYREQUEST} \\
& \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{RECOVERYRESPONSE}, \\
& \hspace{10em} \text{moracle} \mapsto \text{leaderOracle}[i], \\
& \hspace{10em} \text{mepoch} \mapsto \text{currentEpoch}[i]]) \\
& \wedge \text{UNCHANGED} \langle \text{serverVars}, \text{leaderVars}, \text{tempVars}, \text{epochSent}, \text{recoveryVars}, \text{proposalMsgsLog}, \text{testVars} \rangle \\
\text{HandleRecoveryResponse}(i, j) &\triangleq \\
& \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \wedge \text{msgs}[j][i][1].\text{mtype} = \text{RECOVERYRESPONSE} \\
& \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \hspace{2em} \text{infoOk} \triangleq \wedge \text{msg.mepoch} \geq \text{recoveryMaxEpoch}[i] \\
& \hspace{4em} \wedge \text{msg.moracle} \neq \text{NullPoint} \\
& \text{IN } \vee \wedge \text{infoOk} \\
& \hspace{2em} \wedge \text{recoveryMaxEpoch}' = [\text{recoveryMaxEpoch} \text{ EXCEPT } ![i] = \text{msg.mepoch}] \\
& \hspace{2em} \wedge \text{recoveryMEOracle}' = [\text{recoveryMEOracle} \text{ EXCEPT } ![i] = \text{msg.moracle}] \\
& \vee \wedge \neg \text{infoOk} \\
& \hspace{2em} \wedge \text{UNCHANGED} \langle \text{recoveryMaxEpoch}, \text{recoveryMEOracle} \rangle \\
& \wedge \text{Discard}(j, i) \\
& \wedge \text{recoveryRespRecv}' = [\text{recoveryRespRecv} \text{ EXCEPT } ![i] = \text{IF } j \in \text{recoveryRespRecv}[i] \text{ THEN } \text{recoveryRespRecv}[i] \\
& \hspace{15em} \text{ELSE } \text{recoveryRespRecv}[i]] \\
& \wedge \text{UNCHANGED} \langle \text{serverVars}, \text{leaderVars}, \text{tempVars}, \text{epochSent}, \text{recoverySent}, \text{proposalMsgsLog}, \text{testVars} \rangle \\
\text{FindCluster}(i) &\triangleq \\
& \wedge \text{state}[i] = \text{Follower} \\
& \wedge \text{leaderOracle}[i] = \text{NullPoint}
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{ recoveryRespRecv}[i] \in \text{Quorums} \\
& \wedge \text{ LET } \text{infoOk} \triangleq \wedge \text{ recoveryMEOracle}[i] \neq i \\
& \quad \wedge \text{ recoveryMEOracle}[i] \neq \text{NullPoint} \\
& \quad \wedge \text{ currentEpoch}[i] \leq \text{ recoveryMaxEpoch}[i] \\
& \text{ IN } \quad \vee \wedge \neg \text{infoOk} \\
& \quad \wedge \text{ recoverySent}' = [\text{ recoverySent } \text{ EXCEPT } ![i] = \text{FALSE}] \\
& \quad \wedge \text{ UNCHANGED } \langle \text{ currentEpoch}, \text{ leaderOracle}, \text{ msgs} \rangle \\
& \quad \vee \wedge \text{infoOk} \\
& \quad \wedge \text{ currentEpoch}' = [\text{ currentEpoch } \text{ EXCEPT } ![i] = \text{ recoveryMaxEpoch}[i]] \\
& \quad \wedge \text{ leaderOracle}' = [\text{ leaderOracle } \text{ EXCEPT } ![i] = \text{ recoveryMEOracle}[i]] \\
& \quad \wedge \text{ Send}(i, \text{ recoveryMEOracle}[i], [\text{ mtype } \mapsto \text{CEPOCH}, \\
& \quad \quad \quad \text{ mepoch} \mapsto \text{ recoveryMaxEpoch}[i]]) \\
& \quad \wedge \text{ UNCHANGED } \text{ recoverySent} \\
& \wedge \text{ UNCHANGED } \langle \text{ state}, \text{ leaderEpoch}, \text{ history}, \text{ commitIndex}, \text{ leaderVars}, \text{ tempVars}, \\
& \quad \text{ recoveryRespRecv}, \text{ recoveryMaxEpoch}, \text{ recoveryMEOracle}, \text{ cepochSent}, \text{ proposalMsgs} \rangle
\end{aligned}$$

In phase $f11$, follower sends $f.p$ to pleader via CEPOCH .

$$\begin{aligned}
& \text{FollowerDiscovery1}(i) \triangleq \\
& \quad \wedge \text{ state}[i] = \text{Follower} \\
& \quad \wedge \text{ leaderOracle}[i] \neq \text{NullPoint} \\
& \quad \wedge \neg \text{ cepochSent}[i] \\
& \quad \wedge \text{ LET } \text{leader} \triangleq \text{ leaderOracle}[i] \\
& \quad \text{ IN } \quad \text{ Send}(i, \text{leader}, [\text{ mtype } \mapsto \text{CEPOCH}, \\
& \quad \quad \quad \text{ mepoch} \mapsto \text{ currentEpoch}[i]]) \\
& \quad \wedge \text{ cepochSent}' = [\text{ cepochSent } \text{ EXCEPT } ![i] = \text{TRUE}] \\
& \quad \wedge \text{ UNCHANGED } \langle \text{ serverVars}, \text{ leaderVars}, \text{ tempVars}, \text{ recoveryVars}, \text{ proposalMsgsLog}, \text{ testVars} \rangle
\end{aligned}$$

In phase $l11$, pleader receives CEPOCH from a quorum, and choose a new epoch e' as its own $l.p$ and sends NEWPOCH to followers.

$$\begin{aligned}
& \text{LeaderHandleCEPOCH}(i, j) \triangleq \\
& \quad \wedge \text{ state}[i] = \text{ProspectiveLeader} \\
& \quad \wedge \text{ msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{ msgs}[j][i][1].\text{mtype} = \text{CEPOCH} \\
& \quad \wedge \vee \text{ new message - modify } \text{ tempMaxEpoch} \text{ and } \text{ cepochRecv} \\
& \quad \quad \wedge \text{ NullPoint} \notin \text{ cepochRecv}[i] \\
& \quad \quad \wedge \text{ LET } \text{ newEpoch} \triangleq \text{ Maximum}(\{\text{ tempMaxEpoch}[i], \text{ msgs}[j][i][1].\text{mepoch}\}) \\
& \quad \quad \text{ IN } \quad \text{ tempMaxEpoch}' = [\text{ tempMaxEpoch } \text{ EXCEPT } ![i] = \text{ newEpoch}] \\
& \quad \quad \wedge \text{ cepochRecv}' = [\text{ cepochRecv } \text{ EXCEPT } ![i] = \text{ IF } j \in \text{ cepochRecv}[i] \text{ THEN } \text{ cepochRecv}[i] \\
& \quad \quad \quad \text{ ELSE } \text{ cepochRecv}[i] \cup \{j\}] \\
& \quad \wedge \text{ Discard}(j, i) \\
& \quad \vee \text{ new follower who joins in cluster / follower whose history and } \text{ commitIndex} \text{ do not match} \\
& \quad \quad \wedge \text{ NullPoint} \in \text{ cepochRecv}[i] \\
& \quad \quad \wedge \vee \wedge \text{ NullPoint} \notin \text{ ackeRecv}[i] \\
& \quad \quad \quad \wedge \text{ Reply}(i, j, [\text{ mtype } \mapsto \text{NEWPOCH}],
\end{aligned}$$

$$\begin{aligned}
& mepoch \mapsto leaderEpoch[i]) \\
\vee \wedge \text{NullPoint} \in \text{ackRecv}[i] \\
& \wedge \text{Reply2}(i, j, [mtype \mapsto \text{NEWEPOCH}, \\
& \quad mepoch \mapsto leaderEpoch[i], \\
& \quad [mtype \mapsto \text{NEWLEADER}, \\
& \quad mepoch \mapsto currentEpoch[i], \\
& \quad minitialHistory \mapsto initialHistory[i]]) \\
& \wedge \text{UNCHANGED } \langle \text{ceepochRecv}, \text{tempMaxEpoch} \rangle \\
& \wedge \text{cluster}' = [\text{cluster} \text{ EXCEPT } ![i] = \text{IF } j \in \text{cluster}[i] \text{ THEN } \text{cluster}[i] \text{ ELSE } \text{cluster}[i] \cup \{j\}] \\
& \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{ackRecv}, \text{ackldRecv}, \text{ackIndex}, \text{currentCounter}, \text{sendCounter}, \text{initialHist}, \\
& \quad \text{committedIndex}, \text{ceepochSent}, \text{tempMaxLastEpoch}, \text{tempInitialHistory}, \text{recoveryVars}, p \rangle
\end{aligned}$$

Here I decide to change leader's epoch in $l12 \& l21$, otherwise there may exist an old leader and a new leader who share the same epoch. So here I just change $leaderEpoch$, and use it in handling $ACK-E$.

$$\begin{aligned}
\text{LeaderDiscovery1}(i) & \triangleq \\
& \wedge \text{state}[i] = \text{ProspectiveLeader} \\
& \wedge \text{ceepochRecv}[i] \in \text{Quorums} \\
& \wedge \text{leaderEpoch}' = [\text{leaderEpoch} \text{ EXCEPT } ![i] = \text{tempMaxEpoch}[i] + 1] \\
& \wedge \text{ceepochRecv}' = [\text{ceepochRecv} \text{ EXCEPT } ![i] = \{\text{NullPoint}\}] \\
& \wedge \text{Broadcast}(i, [mtype \mapsto \text{NEWEPOCH}, \\
& \quad mepoch \mapsto leaderEpoch'[i]]) \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{leaderOracle}, \text{history}, \text{cluster}, \text{ackRecv}, \text{ackldRecv}, \text{ackIndex}, \text{currentCounter}, \\
& \quad \text{initialHistory}, \text{commitIndex}, \text{committedIndex}, \text{ceepochSent}, \text{tempVars}, \text{recoveryVars}, p \rangle
\end{aligned}$$

In phase $f12$, follower receives $NEWEPOCH$. If $e' > f.p$ then sends back $ACKE$, and $ACKE$ contains $f.a$ and hf to help pleader choose a newer history.

$$\begin{aligned}
\text{FollowerDiscovery2}(i, j) & \triangleq \\
& \wedge \text{state}[i] = \text{Follower} \\
& \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \wedge \text{msgs}[j][i][1].mtype = \text{NEWEPOCH} \\
& \wedge \text{LET } msg \triangleq \text{msgs}[j][i][1] \\
& \text{IN } \vee \text{ new } \text{NEWEPOCH} - \text{accept and reply} \\
& \wedge \text{currentEpoch}[i] < msg.mepoch \\
& \wedge \vee \wedge \text{leaderOracle}[i] = j \\
& \quad \wedge \text{currentEpoch}' = [\text{currentEpoch} \text{ EXCEPT } ![i] = msg.mepoch] \\
& \quad \wedge \text{Reply}(i, j, [mtype \mapsto \text{ACKE}, \\
& \quad \quad mepoch \mapsto msg.mepoch, \\
& \quad \quad mlastEpoch \mapsto leaderEpoch[i], \\
& \quad \quad mhf \mapsto \text{history}[i]]) \\
& \vee \wedge \text{leaderOracle}[i] \neq j \\
& \quad \wedge \text{Discard}(j, i) \\
& \quad \wedge \text{UNCHANGED } \text{currentEpoch} \\
& \vee \wedge \text{currentEpoch}[i] = msg.mepoch \\
& \quad \wedge \vee \wedge \text{leaderOracle}[i] = j \\
& \quad \quad \wedge \text{Reply}(i, j, [mtype \mapsto \text{ACKE},
\end{aligned}$$

$$\begin{aligned}
& mepoch \mapsto msg.mepoch, \\
& mlastEpoch \mapsto leaderEpoch[i], \\
& mh f \mapsto history[i]) \\
& \wedge \text{UNCHANGED } currentEpoch \\
\vee & \text{ It may happen when a leader do not update new epoch to all followers in } Q, \text{ and a new election begins} \\
& \wedge leaderOracle[i] \neq j \\
& \wedge Discard(j, i) \\
& \wedge \text{UNCHANGED } currentEpoch \\
\vee & \text{ stale } NEWPOCH - discard \\
& \wedge currentEpoch[i] > msg.mepoch \\
& \wedge Discard(j, i) \\
& \wedge \text{UNCHANGED } currentEpoch \\
& \wedge \text{UNCHANGED } \langle state, leaderEpoch, leaderOracle, history, leaderVars, \\
& \quad commitIndex, cepochSent, tempVars, recoveryVars, proposalMsgsLog, testVars \rangle
\end{aligned}$$

In phase *l12*, pleader receives *ACKE* from a quorum,
and select the history of one most up-to-date follower to be the initial history.

$$\begin{aligned}
LeaderHandleACKE(i, j) & \triangleq \\
& \wedge state[i] = ProspectiveLeader \\
& \wedge msgs[j][i] \neq \langle \rangle \\
& \wedge msgs[j][i][1].mtype = ACKE \\
& \wedge \text{LET } msg \triangleq msgs[j][i][1] \\
& \quad infoOk \triangleq \vee msg.mlastEpoch > tempMaxLastEpoch[i] \\
& \quad \vee \wedge msg.mlastEpoch = tempMaxLastEpoch[i] \\
& \quad \quad \wedge \vee LastZxid(msg.mhf)[1] > LastZxid(tempInitialHistory[i])[1] \\
& \quad \quad \vee \wedge LastZxid(msg.mhf)[1] = LastZxid(tempInitialHistory[i])[1] \\
& \quad \quad \wedge LastZxid(msg.mhf)[2] \geq LastZxid(tempInitialHistory[i])[2] \\
\text{IN } & \vee \wedge leaderEpoch[i] = msg.mepoch \\
& \wedge \vee \wedge infoOk \\
& \quad \wedge tempMaxLastEpoch' = [tempMaxLastEpoch \text{ EXCEPT } ![i] = msg.mlastEpoch] \\
& \quad \wedge tempInitialHistory' = [tempInitialHistory \text{ EXCEPT } ![i] = msg.mhf] \\
& \vee \wedge \neg infoOk \\
& \quad \wedge \text{UNCHANGED } \langle tempMaxLastEpoch, tempInitialHistory \rangle \\
& \quad \text{Followers not in } Q \text{ will not receive } NEWPOCH, \text{ so leader will receive } ACKE \text{ only when the source is in } Q \\
& \quad \wedge ackeRecv' = [ackeRecv \text{ EXCEPT } ![i] = \text{IF } j \notin ackeRecv[i] \text{ THEN } ackeRecv[i] \cup \{j\} \\
& \quad \quad \quad \text{ELSE } ackeRecv[i]}] \\
& \vee \wedge leaderEpoch[i] \neq msg.mepoch \\
& \quad \wedge \text{UNCHANGED } \langle tempMaxLastEpoch, tempInitialHistory, ackeRecv \rangle \\
& \wedge Discard(j, i) \\
& \wedge \text{UNCHANGED } \langle serverVars, cluster, cepochRecv, ackldRecv, ackIndex, currentCounter, \\
& \quad sendCounter, initialHistory, committedIndex, cepochSent, tempMaxEpoch, recoveryVars \rangle
\end{aligned}$$

$$\begin{aligned}
LeaderDiscovery2Sync1(i) & \triangleq \\
& \wedge state[i] = ProspectiveLeader \\
& \wedge ackeRecv[i] \in Quorums
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{currentEpoch}' = [\text{currentEpoch} \text{ EXCEPT } ![i] = \text{leaderEpoch}[i]] \\
& \wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{tempInitialHistory}[i]] \\
& \wedge \text{initialHistory}' = [\text{initialHistory} \text{ EXCEPT } ![i] = \text{tempInitialHistory}[i]] \\
& \wedge \text{ackRecv}' = [\text{ackRecv} \text{ EXCEPT } ![i] = \{\text{NullPoint}\}] \\
& \wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i][i] = \text{Len}(\text{tempInitialHistory}[i])] \\
& \text{until now, phase1(Discovery) ends} \\
& \wedge \text{Broadcast}(i, [\text{mtype} \mapsto \text{NEWLEADER}, \\
& \quad \text{mepoch} \mapsto \text{currentEpoch}'[i], \\
& \quad \text{minitialHistory} \mapsto \text{history}'[i]]) \\
& \wedge \text{LET } m \triangleq [\text{msource} \mapsto i, \text{mtype} \mapsto \text{NEWLEADER}, \text{mepoch} \mapsto \text{currentEpoch}'[i], \text{mproposals} \mapsto \text{histo} \\
& \quad \text{IN } \text{proposalMsgsLog}' = \text{IF } m \in \text{proposalMsgsLog} \text{ THEN } \text{proposalMsgsLog} \\
& \quad \quad \text{ELSE } \text{proposalMsgsLog} \cup \{m\} \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{leaderEpoch}, \text{leaderOracle}, \text{commitIndex}, \text{cluster}, \text{cepocheRecv}, \text{ackldRecv}, \\
& \quad \text{currentCounter}, \text{sendCounter}, \text{committedIndex}, \text{cepocheSent}, \text{tempVars}, \text{recoveryVars}, p \rangle
\end{aligned}$$

Note1: Delete the change of *commitIndex* in *LeaderDiscovery2Sync1* and *FollowerSync1*, then we can promise that *commitIndex* of every server increases monotonically, except that some server halts and restarts.

Note2: Set *cepocheRecv*, *ackRecv*, *ackldRecv* to $\{\text{NullPoint}\}$ in corresponding three actions to make sure that the prospective leader will not broadcast *NEWLEADER/COMMITLD* twice.

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

$$\begin{aligned}
& \text{FollowerSync1}(i, j) \triangleq \\
& \quad \wedge \text{state}[i] = \text{Follower} \\
& \quad \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{NEWLEADER} \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \quad \text{replyOk} \triangleq \wedge \text{currentEpoch}[i] \leq \text{msg.mepoch} \\
& \quad \quad \quad \wedge \text{leaderOracle}[i] = j \\
& \quad \text{IN } \vee \text{new NEWLEADER} - \text{accept and reply} \\
& \quad \quad \wedge \text{replyOk} \\
& \quad \quad \wedge \text{currentEpoch}' = [\text{currentEpoch} \text{ EXCEPT } ![i] = \text{msg.mepoch}] \\
& \quad \quad \wedge \text{leaderEpoch}' = [\text{leaderEpoch} \text{ EXCEPT } ![i] = \text{msg.mepoch}] \\
& \quad \quad \wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{msg.minitialHistory}] \\
& \quad \quad \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{ACKLD}, \\
& \quad \quad \quad \text{mepoch} \mapsto \text{msg.mepoch}, \\
& \quad \quad \quad \text{mhistory} \mapsto \text{msg.minitialHistory}]) \\
& \quad \vee \text{stale NEWLEADER} - \text{discard} \\
& \quad \quad \wedge \neg \text{replyOk} \\
& \quad \quad \wedge \text{Discard}(j, i) \\
& \quad \quad \wedge \text{UNCHANGED } \langle \text{currentEpoch}, \text{leaderEpoch}, \text{history} \rangle \\
& \quad \wedge \text{UNCHANGED } \langle \text{state}, \text{commitIndex}, \text{leaderOracle}, \text{leaderVars}, \text{tempVars}, \text{cepocheSent}, \text{recoveryVars}, p \rangle
\end{aligned}$$

In phase *l22*, pleader receives *ACK-LD* from a quorum of followers, and sends *COMMIT-LD* to followers.

$$\text{LeaderHandleACKLD}(i, j) \triangleq$$

$\wedge state[i] = ProspectiveLeader$
 $\wedge msgs[j][i] \neq \langle \rangle$
 $\wedge msgs[j][i][1].mtype = ACKLD$
 $\wedge LET\ msg \triangleq msgs[j][i][1]$
 IN \vee new ACK-LD - accept
 $\wedge currentEpoch[i] = msg.mepoch$
 $\wedge ackIndex' = [ackIndex\ EXCEPT\ ![i][j] = Len(initialHistory[i])]$
 $\wedge ackldRecv' = [ackldRecv\ EXCEPT\ ![i] = IF\ j \notin ackldRecv[i]\ THEN\ ackldRecv[i] \cup \{j\}$
 $ELSE\ ackldRecv[i]]$
 \vee stale ACK-LD - discard
 $\wedge currentEpoch[i] \neq msg.mepoch$
 $\wedge UNCHANGED\ \langle ackldRecv, ackIndex \rangle$
 $\wedge Discard(j, i)$
 $\wedge UNCHANGED\ \langle serverVars, cluster, cepochRecv, ackeRecv, currentCounter,$
 $sendCounter, initialHistory, committedIndex, tempVars, cepochSent, recoveryVars, p$

$LeaderSync2(i) \triangleq$
 $\wedge state[i] = ProspectiveLeader$
 $\wedge ackldRecv[i] \in Quorums$
 $\wedge commitIndex' = [commitIndex\ EXCEPT\ ![i] = Len(history[i])]$
 $\wedge committedIndex' = [committedIndex\ EXCEPT\ ![i] = Len(history[i])]$
 $\wedge state' = [state\ EXCEPT\ ![i] = Leader]$
 $\wedge currentCounter' = [currentCounter\ EXCEPT\ ![i] = 0]$
 $\wedge sendCounter' = [sendCounter\ EXCEPT\ ![i] = 0]$
 $\wedge ackldRecv' = [ackldRecv\ EXCEPT\ ![i] = \{NullPoint\}]$
 $\wedge Broadcast(i, [mtype \mapsto COMMITLD,$
 $mepoch \mapsto currentEpoch[i],$
 $mlength \mapsto Len(history[i])])$
 $\wedge UNCHANGED\ \langle currentEpoch, leaderEpoch, leaderOracle, history, cluster, cepochRecv,$
 $ackeRecv, ackIndex, initialHistory, tempVars, cepochSent, recoveryVars, proposalM$

In phase $f22$, follower receives $COMMIT$ -LD and delivers all unprocessed transaction.
 $FollowerSync2(i, j) \triangleq$
 $\wedge state[i] = Follower$
 $\wedge msgs[j][i] \neq \langle \rangle$
 $\wedge msgs[j][i][1].mtype = COMMITLD$
 $\wedge LET\ msg \triangleq msgs[j][i][1]$
 $replyOk \triangleq \wedge currentEpoch[i] = msg.mepoch$
 $\wedge leaderOracle[i] = j$
 IN \vee new $COMMIT$ -LD - commit all transactions in initial history
Regardless of $Restart$, it must be true because one will receive $NEWLEADER$ before receiving $COMMIT$ -LD
 $\wedge replyOk$
 $\wedge \vee \wedge Len(history[i]) = msg.mlength$
 $\wedge commitIndex' = [commitIndex\ EXCEPT\ ![i] = Len(history[i])]$
 $\wedge Discard(j, i)$

$$\begin{aligned}
& \vee \wedge \text{Len}(\text{history}[i]) \neq \text{msg.mlength} \\
& \quad \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{CEPOCH}, \\
& \quad \quad \quad \text{mepoch} \mapsto \text{currentEpoch}[i]]) \\
& \quad \wedge \text{UNCHANGED } \text{commitIndex} \\
\vee & \text{ } > : \text{stale } \text{COMMIT-LD} - \text{discard} \\
& < : \text{In our implementation, ' < ' does not exist due to the guarantee of } \text{Restart} \\
& \wedge \neg \text{replyOk} \\
& \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED } \text{commitIndex} \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{leaderEpoch}, \text{leaderOracle}, \text{history}, \\
& \quad \text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, \text{proposalMsgsLog}, \text{testVars} \rangle
\end{aligned}$$

In phase *l31*, leader receives client request and broadcasts *PROPOSE*.

$$\begin{aligned}
\text{ClientRequest}(i, v) & \triangleq \\
& \text{test restrictions} \\
& \wedge \text{Len}(\text{history}[i]) < \text{MaxTransactionNum} \\
& \wedge \text{state}[i] = \text{Leader} \\
& \wedge \text{currentCounter}' = [\text{currentCounter} \text{ EXCEPT } ![i] = \text{currentCounter}[i] + 1] \\
& \wedge \text{LET } \text{newTransaction} \triangleq [\text{epoch} \mapsto \text{currentEpoch}[i], \\
& \quad \quad \quad \text{counter} \mapsto \text{currentCounter}'[i], \\
& \quad \quad \quad \text{value} \mapsto v] \\
& \text{IN } \wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{Append}(\text{history}[i], \text{newTransaction})] \\
& \quad \wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i][i] = \text{Len}(\text{history}'[i])] \text{ necessary, to push } \text{commitIndex} \\
& \wedge \text{UNCHANGED } \langle \text{msgs}, \text{state}, \text{currentEpoch}, \text{leaderEpoch}, \text{leaderOracle}, \text{commitIndex}, \text{cluster}, \text{ceepochR}, \\
& \quad \text{ackeRecv}, \text{ackldRecv}, \text{sendCounter}, \text{initialHistory}, \text{committedIndex}, \text{tempVars}, \text{ceepoch} \rangle \\
\text{LeaderBroadcast1}(i) & \triangleq \\
& \wedge \text{state}[i] = \text{Leader} \\
& \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{PROPOSE}, \\
& \quad \quad \quad \text{mepoch} \mapsto \text{currentEpoch}[i], \\
& \quad \quad \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{mtype} \mapsto \text{PROPOSE}, \text{mepoch} \mapsto \text{currentEpoch}[i], \text{mproposal} \mapsto \text{toBeSentEntry}] \\
& \quad \text{IN } \text{proposalMsgsLog}' = \text{proposalMsgsLog} \cup \{m\} \\
& \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{ceepochRecv}, \text{cluster}, \text{ackeRecv}, \text{ackldRecv}, \text{ackIndex}, \\
& \quad \text{currentCounter}, \text{initialHistory}, \text{committedIndex}, \text{tempVars}, \text{recoveryVars}, \text{ceepochSent} \rangle
\end{aligned}$$

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

$$\begin{aligned}
\text{FollowerBroadcast1}(i, j) & \triangleq \\
& \wedge \text{state}[i] = \text{Follower} \\
& \wedge \text{msgs}[j][i] \neq \langle \rangle
\end{aligned}$$

$\wedge \text{msgs}[j][i][1].\text{mtype} = \text{PROPOSE}$
 $\wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1]$
 $\text{replyOk} \triangleq \wedge \text{currentEpoch}[i] = \text{msg.mepoch}$
 $\wedge \text{leaderOracle}[i] = j$
IN \vee $\text{It should be that } \vee \text{msg.mproposal.counter} = 1$
 $\vee \text{msg.mrpposal.counter} = \text{history}[\text{Len}(\text{history})].\text{counter} + 1$
 $\wedge \text{replyOk}$
 $\wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{Append}(\text{history}[i], \text{msg.mproposal})]$
 $\wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{ACK},$
 $\text{mepoch} \mapsto \text{currentEpoch}[i],$
 $\text{mindex} \mapsto \text{Len}(\text{history}'[i])])$
 \vee $\text{If happens, } \neq \text{ must be } >, \text{ namely a stale leader sends it.}$
 $\wedge \neg \text{replyOk}$
 $\wedge \text{Discard}(j, i)$
 $\wedge \text{UNCHANGED } \text{history}$
 $\wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{leaderEpoch}, \text{leaderOracle}, \text{commitIndex},$
 $\text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, \text{proposalMsgsLog}, \text{testVars} \rangle$

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

$\text{LeaderHandleACK}(i, j) \triangleq$
 $\wedge \text{state}[i] = \text{Leader}$
 $\wedge \text{msgs}[j][i] \neq \langle \rangle$
 $\wedge \text{msgs}[j][i][1].\text{mtype} = \text{ACK}$
 $\wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1]$
IN \vee $\text{It should be that } \text{ackIndex}[i][j] + 1 \triangleq \text{msg.mindex}$
 $\wedge \text{currentEpoch}[i] = \text{msg.mepoch}$
 $\wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i][j] = \text{Maximum}(\{\text{ackIndex}[i][j], \text{msg.mindex}\})]$
 \vee $\text{If happens, } \neq \text{ must be } >, \text{ namely a stale follower sends it.}$
 $\wedge \text{currentEpoch}[i] \neq \text{msg.mepoch}$
 $\wedge \text{UNCHANGED } \text{ackIndex}$
 $\wedge \text{Discard}(j, i)$
 $\wedge \text{UNCHANGED } \langle \text{serverVars}, \text{cluster}, \text{ceepochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{currentCounter},$
 $\text{sendCounter}, \text{initialHistory}, \text{committedIndex}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, p \rangle$
 $\text{LeaderAdvanceCommit}(i) \triangleq$
 $\wedge \text{state}[i] = \text{Leader}$
 $\wedge \text{commitIndex}[i] < \text{Len}(\text{history}[i])$
 $\wedge \text{LET } \text{Agree}(\text{index}) \triangleq \{i\} \cup \{k \in (\text{Server} \setminus \{i\}) : \text{ackIndex}[i][k] \geq \text{index}\}$
 $\text{agreeIndexes} \triangleq \{\text{index} \in (\text{commitIndex}[i] + 1) \dots \text{Len}(\text{history}[i]) : \text{Agree}(\text{index}) \in \text{Quorum}\}$
 $\text{newCommitIndex} \triangleq \text{IF } \text{agreeIndexes} \neq \{\} \text{ THEN } \text{Maximum}(\text{agreeIndexes})$
 $\text{ELSE } \text{commitIndex}[i]$
IN $\text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{newCommitIndex}]$
 $\wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{leaderEpoch}, \text{leaderOracle}, \text{history},$
 $\text{msgs}, \text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, \text{proposalMsgsLog}, \text{testVars} \rangle$

$$\begin{aligned}
& \text{LeaderBroadcast2}(i) \triangleq \\
& \quad \wedge \text{state}[i] = \text{Leader} \\
& \quad \wedge \text{committedIndex}[i] < \text{commitIndex}[i] \\
& \quad \wedge \text{LET } \text{newCommittedIndex} \triangleq \text{committedIndex}[i] + 1 \\
& \quad \text{IN } \wedge \text{Broadcast}(i, [\text{mtype} \mapsto \text{COMMIT}, \\
& \quad \quad \quad \text{mepoch} \mapsto \text{currentEpoch}[i], \\
& \quad \quad \quad \text{mindex} \mapsto \text{newCommittedIndex}, \\
& \quad \quad \quad \text{mcounter} \mapsto \text{history}[i][\text{newCommittedIndex}].\text{counter}]) \\
& \quad \wedge \text{committedIndex}' = [\text{committedIndex} \text{ EXCEPT } ![i] = \text{committedIndex}[i] + 1] \\
& \quad \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{cluster}, \text{cePOCHRecv}, \text{ackRecv}, \text{ackldRecv}, \text{ackIndex}, \text{currentCounter}, \\
& \quad \quad \text{sendCounter}, \text{initialHistory}, \text{tempVars}, \text{cePOCHSent}, \text{recoveryVars}, \text{proposalMsgsLog}, t \rangle
\end{aligned}$$

In phase $f32$, follower receives *COMMIT* and commits transaction.

$$\begin{aligned}
& \text{FollowerBroadcast2}(i, j) \triangleq \\
& \quad \wedge \text{state}[i] = \text{Follower} \\
& \quad \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{COMMIT} \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{COMMIT} \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \quad \text{replyOk} \triangleq \wedge \text{currentEpoch}[i] = \text{msg.mepoch} \\
& \quad \quad \quad \wedge \text{leaderOracle}[i] = j \\
& \quad \text{IN } \vee \wedge \text{replyOk} \\
& \quad \quad \wedge \text{LET } \text{infoOk} \triangleq \wedge \text{Len}(\text{history}[i]) \geq \text{msg.mindex} \\
& \quad \quad \quad \wedge \vee \wedge \text{msg.mindex} > 0 \\
& \quad \quad \quad \quad \wedge \text{history}[i][\text{msg.mindex}].\text{epoch} = \text{msg.mepoch} \\
& \quad \quad \quad \quad \wedge \text{history}[i][\text{msg.mindex}].\text{counter} = \text{msg.mcounter} \\
& \quad \quad \quad \vee \text{msg.mindex} = 0 \\
& \quad \quad \text{IN } \vee \text{new COMMIT} - \text{commit transaction in history} \\
& \quad \quad \quad \wedge \text{infoOk} \\
& \quad \quad \quad \wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{Maximum}(\{\text{commitIndex}[i], \text{msg.mindex}\})] \\
& \quad \quad \quad \wedge \text{Discard}(j, i) \\
& \quad \quad \vee \text{It may happen when the server is a new follower who joined in the cluster,} \\
& \quad \quad \quad \text{and it misses the corresponding PROPOSE.} \\
& \quad \quad \quad \wedge \neg \text{infoOk} \\
& \quad \quad \quad \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{CEPOCH}, \\
& \quad \quad \quad \quad \text{mepoch} \mapsto \text{currentEpoch}[i]]) \\
& \quad \quad \quad \wedge \text{UNCHANGED } \text{commitIndex} \\
& \quad \vee \text{stale COMMIT} - \text{discard} \\
& \quad \quad \wedge \neg \text{replyOk} \\
& \quad \quad \wedge \text{Discard}(j, i) \\
& \quad \quad \wedge \text{UNCHANGED } \text{commitIndex} \\
& \quad \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{leaderEpoch}, \text{history}, \text{leaderOracle}, \\
& \quad \quad \text{leaderVars}, \text{tempVars}, \text{cePOCHSent}, \text{recoveryVars}, \text{proposalMsgsLog}, \text{testVars} \rangle
\end{aligned}$$

There may be two ways to make sure all followers as up-to-date as the leader.
way1: choose *Send* not *Broadcast* when leader is going to send *PROPOSE* and *COMMIT*.
way2: When one follower receives *PROPOSE* or *COMMIT* which misses some entries between its history and the newest entry, the follower send *CEPOCH* to catch pace.
Here I choose *way2*, which I need not to rewrite *PROPOSE* and *COMMIT*, but need to modify the code when follower receives *COMMIT-LD* and *COMMIT*.

In phase *l33*, upon receiving *CEPOCH*, leader *l* proposes back *NEWEPOCH* and *NEWLEADER*.
 $LeaderHandleCEPOCHinPhase3(i, j) \triangleq$

$$\begin{aligned} & \wedge state[i] = Leader \\ & \wedge msgs[j][i] \neq \langle \rangle \\ & \wedge msgs[j][i][1].mtype = CEPOCH \\ & \wedge LET \ msg \triangleq \ msgs[j][i][1] \\ & \quad IN \quad \vee \wedge currentEpoch[i] \geq msg.mepoch \\ & \quad \quad \wedge Reply2(i, j, [mtype \mapsto NEWEPOCH, \\ & \quad \quad \quad mepoch \mapsto currentEpoch[i], \\ & \quad \quad \quad [mtype \mapsto NEWLEADER, \\ & \quad \quad \quad mepoch \mapsto currentEpoch[i], \\ & \quad \quad \quad minitialHistory \mapsto history[i]]) \\ & \quad \quad \wedge LET \ m \triangleq [msource \mapsto i, mtype \mapsto NEWLEADER, mepoch \mapsto currentEpoch[i], mproposal \\ & \quad \quad \quad IN \quad proposalMsgsLog' = IF \ m \in proposalMsgsLog \ THEN \ proposalMsgsLog \\ & \quad \quad \quad \quad \quad \quad \quad ELSE \ proposalMsgsLog \cup \{m\} \\ & \quad \vee \wedge currentEpoch[i] < msg.mepoch \\ & \quad \quad \wedge UNCHANGED \langle msgs, proposalMsgsLog \rangle \\ & \wedge UNCHANGED \langle serverVars, leaderVars, tempVars, cepochSent, recoveryVars, testVars \rangle \end{aligned}$$

In phase *l34*, upon receiving ack from *f* of the *NEWLEADER*, it sends a commit message to *f*.
Leader *l* also makes $Q := Q \cup \{f\}$.

$LeaderHandleACKLDinPhase3(i, j) \triangleq$

$$\begin{aligned} & \wedge state[i] = Leader \\ & \wedge msgs[j][i] \neq \langle \rangle \\ & \wedge msgs[j][i][1].mtype = ACKLD \\ & \wedge LET \ msg \triangleq \ msgs[j][i][1] \\ & \quad aimCommitIndex \triangleq Minimum(\{commitIndex[i], Len(msg.mhistory)\}) \\ & \quad aimCommitCounter \triangleq IF \ aimCommitIndex = 0 \ THEN \ 0 \ ELSE \ history[i][aimCommitIndex].co \\ & \quad IN \quad \vee \wedge currentEpoch[i] = msg.mepoch \\ & \quad \quad \wedge ackIndex' = [ackIndex \ EXCEPT \ ![i][j] = Len(msg.mhistory)] \\ & \quad \quad \wedge Reply(i, j, [mtype \mapsto COMMIT, \\ & \quad \quad \quad mepoch \mapsto currentEpoch[i], \\ & \quad \quad \quad mindex \mapsto aimCommitIndex, \\ & \quad \quad \quad mcounter \mapsto aimCommitCounter]) \\ & \quad \vee \wedge currentEpoch[i] \neq msg.mepoch \\ & \quad \quad \wedge Discard(j, i) \\ & \quad \quad \wedge UNCHANGED \ ackIndex \\ & \wedge cluster' = [cluster \ EXCEPT \ ![i] = IF \ j \in cluster[i] \ THEN \ cluster[i] \end{aligned}$$

ELSE $cluster[i] \cup \{j\}$

\wedge UNCHANGED $\langle serverVars, cepochRecv, ackeRecv, ackldRecv, currentCounter, sendCounter, initialHistory, committedIndex, tempVars, cepochSent, recoveryVars, proposalMsgsLo,$

To ensure any follower can find the correct leader, the follower should modify *leaderOracle* anytime when it receive messages from leader, because a server may restart and join the cluster *Q* halfway and receive the first message which is not *NEWEPOCH*. But we can delete this restriction when we ensure *Broadcast* function acts on the followers in the cluster not any servers in the whole system, then one server must has correct *leaderOracle* before it receives messages.

Let me suppose two conditions when one follower sends *CEPOCH* to leader:

0. Usually, the server becomes follower in election and sends *CEPOCH* before receiving *NEWEPOCH*.
1. The follower wants to join the cluster halfway and get the newest history.
2. The follower has received *COMMIT*, but there exists the gap between its own history and *mindex*, which means there are some transactions before *mindex* miss. Here we choose to send *CEPOCH* again, to receive the newest history from leader.

$BecomeFollower(i) \triangleq$
 $\wedge state[i] \neq Follower$
 $\wedge \exists j \in Server \setminus \{i\} : \wedge msgs[j][i] \neq \langle \rangle$
 $\wedge msgs[j][i][1].mtype \neq RECOVERYREQUEST$
 $\wedge msgs[j][i][1].mtype \neq RECOVERYRESPONSE$
 $\wedge LET \ msg \triangleq \ msgs[j][i][1]$
 IN $\wedge NullPoint \in cepochRecv[i]$
 $\wedge Maximum(\{currentEpoch[i], leaderEpoch[i]\}) < msg.mepoch$
 $\wedge \vee msg.mtype = NEWEPOCH$
 $\vee msg.mtype = NEWLEADER$
 $\vee msg.mtype = COMMITLD$
 $\vee msg.mtype = PROPOSE$
 $\vee msg.mtype = COMMIT$
 $\wedge state' = [state \text{ EXCEPT } ![i] = Follower]$
 $\wedge currentEpoch' = [currentEpoch \text{ EXCEPT } ![i] = msg.mepoch]$
 $\wedge leaderOracle' = [leaderOracle \text{ EXCEPT } ![i] = j]$
 $\wedge Reply(i, j, [mtype \mapsto CEPOCH,$
 $\hspace{10em} mepoch \mapsto currentEpoch[i]])$

Here we should not use *Discard*.

\wedge UNCHANGED $\langle leaderEpoch, history, commitIndex, leaderVars, tempVars, cepochSent, recoveryVars,$

$DiscardStaleMessage(i) \triangleq$
 $\wedge \exists j \in Server \setminus \{i\} : \wedge msgs[j][i] \neq \langle \rangle$
 $\wedge msgs[j][i][1].mtype \neq RECOVERYREQUEST$
 $\wedge msgs[j][i][1].mtype \neq RECOVERYRESPONSE$
 $\wedge LET \ msg \triangleq \ msgs[j][i][1]$
 IN $\vee \wedge state[i] = Follower$
 $\wedge \vee msg.mepoch < currentEpoch[i] \setminus * \text{ Discussed before.}$
 $\vee msg.mtype = CEPOCH$

$$\begin{aligned}
& \vee \text{msg.mtype} = \text{ACKE} \\
& \vee \text{msg.mtype} = \text{ACKLD} \\
& \vee \text{msg.mtype} = \text{ACK} \\
& \vee \wedge \text{state}[i] \neq \text{Follower} \\
& \wedge \text{msg.mtype} \neq \text{CEPOCH} \\
& \wedge \vee \wedge \text{state}[i] = \text{ProspectiveLeader} \\
& \quad \wedge \vee \text{msg.mtype} = \text{ACK} \\
& \quad \quad \vee \wedge \text{msg.mepoch} \leq \text{Maximum}(\{\text{currentEpoch}[i], \text{leaderEpoch}[i]\}) \\
& \quad \quad \quad \wedge \vee \text{msg.mtype} = \text{NEWPOCH} \\
& \quad \quad \quad \vee \text{msg.mtype} = \text{NEWLEADER} \\
& \quad \quad \quad \vee \text{msg.mtype} = \text{COMMITLD} \\
& \quad \quad \quad \vee \text{msg.mtype} = \text{PROPOSE} \\
& \quad \quad \quad \vee \text{msg.mtype} = \text{COMMIT} \\
& \vee \wedge \text{state}[i] = \text{Leader} \\
& \quad \wedge \vee \text{msg.mtype} = \text{ACKE} \\
& \quad \quad \vee \wedge \text{msg.mepoch} \leq \text{currentEpoch}[i] \\
& \quad \quad \quad \wedge \vee \text{msg.mtype} = \text{NEWPOCH} \\
& \quad \quad \quad \vee \text{msg.mtype} = \text{NEWLEADER} \\
& \quad \quad \quad \vee \text{msg.mtype} = \text{COMMITLD} \\
& \quad \quad \quad \vee \text{msg.mtype} = \text{PROPOSE} \\
& \quad \quad \quad \vee \text{msg.mtype} = \text{COMMIT} \\
& \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, \text{proposalMsgsLog}, \text{testV} \rangle
\end{aligned}$$

Defines how the variables may transition.

$\text{Next} \triangleq$

$$\begin{aligned}
& \vee \exists i \in \text{Server}, Q \in \text{Quorums} : \text{InitialElection}(i, Q) \\
& \vee \exists i \in \text{Server} : \text{Restart}(i) \\
& \vee \exists i \in \text{Server} : \text{RecoveryAfterRestart}(i) \\
& \vee \exists i, j \in \text{Server} : \text{HandleRecoveryRequest}(i, j) \\
& \vee \exists i, j \in \text{Server} : \text{HandleRecoveryResponse}(i, j) \\
& \vee \exists i, j \in \text{Server} : \text{FindCluster}(i) \\
& \vee \exists i, j \in \text{Server} : \text{LeaderTimeout}(i, j) \\
& \vee \exists i \in \text{Server} : \text{FollowerTimeout}(i) \\
& \vee \exists i \in \text{Server} : \text{FollowerDiscovery1}(i) \\
& \vee \exists i, j \in \text{Server} : \text{LeaderHandleCEPOCH}(i, j) \\
& \vee \exists i \in \text{Server} : \text{LeaderDiscovery1}(i) \\
& \vee \exists i, j \in \text{Server} : \text{FollowerDiscovery2}(i, j) \\
& \vee \exists i, j \in \text{Server} : \text{LeaderHandleACKE}(i, j) \\
& \vee \exists i \in \text{Server} : \text{LeaderDiscovery2Sync1}(i) \\
& \vee \exists i, j \in \text{Server} : \text{FollowerSync1}(i, j) \\
& \vee \exists i, j \in \text{Server} : \text{LeaderHandleACKLD}(i, j) \\
& \vee \exists i \in \text{Server} : \text{LeaderSync2}(i)
\end{aligned}$$

$$\begin{aligned}
& \vee \exists i, j \in \text{Server} : \text{FollowerSync2}(i, j) \\
& \vee \exists i \in \text{Server}, v \in \text{Value} : \text{ClientRequest}(i, v) \\
& \vee \exists i \in \text{Server} : \text{LeaderBroadcast1}(i) \\
& \vee \exists i, j \in \text{Server} : \text{FollowerBroadcast1}(i, j) \\
& \vee \exists i, j \in \text{Server} : \text{LeaderHandleACK}(i, j) \\
& \vee \exists i \in \text{Server} : \text{LeaderAdvanceCommit}(i) \\
& \vee \exists i \in \text{Server} : \text{LeaderBroadcast2}(i) \\
& \vee \exists i, j \in \text{Server} : \text{FollowerBroadcast2}(i, j) \\
& \vee \exists i, j \in \text{Server} : \text{LeaderHandleCEPOCHinPhase3}(i, j) \\
& \vee \exists i, j \in \text{Server} : \text{LeaderHandleACKLDinPhase3}(i, j) \\
& \vee \exists i \in \text{Server} : \text{DiscardStaleMessage}(i) \\
& \vee \exists i \in \text{Server} : \text{BecomeFollower}(i)
\end{aligned}$$

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

Define some variants, safety propoties, and liveness propoties of *Zab* consensus algorithm.

Safety properties

There is most one leader/prospective leader in a certain epoch.

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

$$\begin{aligned}
& \wedge \vee \text{state}[i] = \text{Leader} \\
& \vee \wedge \text{state}[i] = \text{ProspectiveLeader} \\
& \wedge \text{NullPoint} \in \text{ackRecv}[i] \quad \text{prospective leader determines its epoch after broadcasting NEWLE} \\
& \wedge \vee \text{state}[j] = \text{Leader} \\
& \vee \wedge \text{state}[j] = \text{ProspectiveLeader} \\
& \wedge \text{NullPoint} \in \text{ackRecv}[j] \\
& \wedge \text{currentEpoch}[i] = \text{currentEpoch}[j] \\
& \Rightarrow i = j
\end{aligned}$$

Here, delivering means deliver some transaction from history to replica. We can assume $\text{deliverIndex} = \text{commitIndex}$.

So we can assume the set of delivered transactions is the prefix of history with index from 1 to commitIndex .

We can express a transaction by two-tuple $\langle \text{epoch}, \text{counter} \rangle$ according to its uniqueness.

$$\begin{aligned}
\text{equal}(\text{entry1}, \text{entry2}) & \triangleq \wedge \text{entry1.epoch} = \text{entry2.epoch} \\
& \wedge \text{entry1.counter} = \text{entry2.counter}
\end{aligned}$$

$$\begin{aligned}
\text{precede}(\text{entry1}, \text{entry2}) & \triangleq \vee \text{entry1.epoch} < \text{entry2.epoch} \\
& \vee \wedge \text{entry1.epoch} = \text{entry2.epoch} \\
& \wedge \text{entry1.counter} < \text{entry2.counter}
\end{aligned}$$

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

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

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

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

$$\begin{aligned} & state[i] = Follower \wedge commitIndex[i] > 0 \\ \Rightarrow & \forall index \in 1 \dots commitIndex[i] : \exists msg \in proposalMsgsLog : \\ & \quad \vee \wedge msg.mtype = PROPOSE \\ & \quad \wedge equal(msg.mproposal, history[i][index]) \\ & \quad \vee \wedge msg.mtype = NEWLEADER \\ & \quad \wedge \exists pindex \in 1 \dots Len(msg.mproposals) : equal(msg.mproposals[pindex], history[i][index]) \end{aligned}$$

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

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

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

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

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

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

$$\begin{aligned} LocalPrimaryOrder \triangleq & LET mset(i, e) \triangleq \{msg \in proposalMsgsLog : \wedge msg.mtype = PROPOSE \\ & \quad \wedge msg.msource = i \\ & \quad \wedge msg.mepoch = e\} \\ & mentries(i, e) \triangleq \{msg.mproposal : msg \in mset(i, e)\} \\ IN & \forall i \in Server : \forall e \in 1 \dots currentEpoch[i] : \\ & \quad \vee Cardinality(mentries(i, e)) < 2 \\ & \quad \vee \wedge Cardinality(mentries(i, e)) \geq 2 \\ & \quad \quad \wedge \exists tsc1, tsc2 \in mentries(i, e) : \\ & \quad \quad \quad \vee equal(tsc1, tsc2) \\ & \quad \quad \quad \vee \wedge \neg equal(tsc1, tsc2) \\ & \quad \quad \quad \wedge LET tscPre \triangleq IF precede(tsc1, tsc2) THEN tsc1 ELSE tsc2 \end{aligned}$$

$$\begin{aligned}
& tscNext \triangleq \text{IF } precede(tsc1, tsc2) \text{ THEN } tsc2 \text{ ELSE } tsc1 \\
\text{IN } & \forall j \in Server : \wedge commitIndex[j] \geq 2 \\
& \quad \wedge \exists index \in 1 \dots commitIndex[j] : equal(history[j][index], tsc1) \\
\Rightarrow & \exists index2 \in 1 \dots commitIndex[j] : \\
& \quad \wedge equal(history[j][index2], tscNext) \\
& \quad \wedge index2 > 1 \\
& \quad \wedge \exists index1 \in 1 \dots (index2 - 1) : equal(history[j][index1], tscPre)
\end{aligned}$$

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 .

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

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 .

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

Liveness property

Suppose that :

- A quorum Q of followers are up.
- The followers in Q elect the same process l and l is up.
- Messages between a follower in Q and l are received in a timely fashion.

If l proposes a transaction a , then a is eventually committed.

\ * Modification History

\ * Last modified Wed May 05 22:09:48 CST 2021 by Dell

\ * Created Sat Dec 05 13:32:08 CST 2020 by Dell