
MODULE *ZabWithQ*

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

This work is driven by *Flavio P. Junqueira*, “*Zab*: High-performance broadcast for primary-backup systems”

EXTENDS *Integers, FiniteSets, Sequences, Naturals, TLC*

The set of server identifiers
CONSTANT *Server*

The set of requests that can go into history
CONSTANT *Value*

Server states
It is unnecessary to add state ELECTION, we can own it by setting *leaderOracle* to Null.
CONSTANTS *Follower, Leader, ProspectiveLeader*

Message types
CONSTANTS *CEPOCH, NEWEPOCH, ACKE, NEWLEADER, ACKLD, COMMITLD, PROPOSE, ACK, C*

Additional Message types used for recovery when restarting
CONSTANTS *RECOVERYREQUEST, RECOVERYRESPONSE*

the maximum round of epoch (initially {0, 1, 2}), currently not used
CONSTANT *Epoches*

Return the maximum value from the set *S*
 $Maximum(S) \triangleq \text{IF } S = \{\} \text{ THEN } -1$
 $\text{ELSE CHOOSE } n \in S : \forall m \in S : n \geq m$

Return the minimum value from the set *S*
 $Minimum(S) \triangleq \text{IF } S = \{\} \text{ THEN } -1$
 $\text{ELSE CHOOSE } n \in S : \forall m \in S : n \leq m$

$Quorums \triangleq \{Q \in \text{SUBSET } Server : Cardinality(Q) * 2 > Cardinality(Server)\}$
 ASSUME $QuorumsAssumption \triangleq \wedge \forall Q \in Quorums : Q \subseteq Server$
 $\wedge \forall Q1, Q2 \in Quorums : Q1 \cap Q2 \neq \{\}$

$None \triangleq \text{CHOOSE } v : v \notin Value$

$NullPoint \triangleq \text{CHOOSE } p : p \notin Server$

The server's *state(Follower, Leader, ProspectiveLeader)*.
VARIABLE *state*

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 $applyIndex(deliverIndex \text{ here}) = commitIndex$.
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*

Hepler 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*

Persistent state of a server: history, *currentEpoch*, *leaderEpoch*

$serverVars \triangleq \langle state, currentEpoch, leaderEpoch, leaderOracle, history, commitIndex \rangle$

$leaderVars \triangleq \langle cluster, cepochRecv, ackRecv, ackldRecv, ackIndex, currentCounter, sendCounter, initialHistory \rangle$

$tempVars \triangleq \langle tempMaxEpoch, tempMaxLastEpoch, tempInitialHistory \rangle$

$recoveryVars \triangleq \langle recoveryRespRecv, recoveryMaxEpoch, recoveryMEOracle \rangle$

$vars \triangleq \langle serverVars, msgs, leaderVars, tempVars, recoveryVars, cepochSent, proposalMsgsLog \rangle$

$LastZxid(his) \triangleq \text{IF } Len(his) > 0 \text{ THEN } \langle his[Len(his)].epoch, his[Len(his)].counter \rangle$
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*]
 $Discard(i, j) \triangleq msgs' = \text{IF } msgs[i][j] \neq \langle \rangle \text{ THEN } [msgs \text{ EXCEPT } ![i][j] = Tail(msgs[i][j])]$
ELSE *msgs*

Leader/Pleader broadcasts a message to all other servers in *Q*
 $Broadcast(i, m) \triangleq msgs' = [ii \in Server \mapsto [ij \in Server \mapsto \text{IF } ii = i \wedge ij \neq i$

$$\wedge ij \in cluster[i] \text{ THEN } Append(msgs[i][ij], m) \\ \text{ELSE } msgs[i][ij]]$$

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

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

$$\begin{aligned} Init \triangleq & \wedge state = [s \in Server \mapsto Follower] \\ & \wedge currentEpoch = [s \in Server \mapsto 0] \\ & \wedge leaderEpoch = [s \in Server \mapsto 0] \\ & \wedge leaderOracle = [s \in Server \mapsto NullPoint] \\ & \wedge history = [s \in Server \mapsto \langle \rangle] \\ & \wedge msgs = [i \in Server \mapsto [j \in Server \mapsto \langle \rangle]] \\ & \wedge cluster = [i \in Server \mapsto \{\}] \\ & \wedge cepochRecv = [s \in Server \mapsto \{\}] \\ & \wedge ackRecv = [s \in Server \mapsto \{\}] \\ & \wedge ackldRecv = [s \in Server \mapsto \{\}] \\ & \wedge ackIndex = [i \in Server \mapsto [j \in Server \mapsto 0]] \\ & \wedge currentCounter = [s \in Server \mapsto 0] \\ & \wedge sendCounter = [s \in Server \mapsto 0] \\ & \wedge commitIndex = [s \in Server \mapsto 0] \\ & \wedge committedIndex = [s \in Server \mapsto 0] \\ & \wedge initialHistory = [s \in Server \mapsto \langle \rangle] \\ & \wedge cepochSent = [s \in Server \mapsto FALSE] \\ & \wedge tempMaxEpoch = [s \in Server \mapsto 0] \\ & \wedge tempMaxLastEpoch = [s \in Server \mapsto 0] \\ & \wedge tempInitialHistory = [s \in Server \mapsto \langle \rangle] \\ & \wedge recoveryRespRecv = [s \in Server \mapsto \{\}] \\ & \wedge recoveryMaxEpoch = [s \in Server \mapsto 0] \\ & \wedge recoveryMEOracle = [s \in Server \mapsto NullPoint] \\ & \wedge proposalMsgsLog = \{\} \end{aligned}$$

A server becomes pleader and a quorum servers knows that.

$$\begin{aligned} Election(i, Q) \triangleq & \\ & \wedge i \in Q \\ & \wedge state' = [s \in Server \mapsto \text{IF } s = i \text{ THEN } ProspectiveLeader] \end{aligned}$$

ELSE IF $s \in Q$ THEN *Follower*
ELSE $state[s]$]]

$\wedge cluster'$ = $[cluster \text{ EXCEPT } ![i] = Q]$ cluster is first initialized in election, not *phase1*.
 $\wedge cepochRecv'$ = $[ceepochRecv \text{ EXCEPT } ![i] = \{i\}]$
 $\wedge ackeRecv'$ = $[ackeRecv \text{ EXCEPT } ![i] = \{i\}]$
 $\wedge ackldRecv'$ = $[ackldRecv \text{ EXCEPT } ![i] = \{i\}]$
 $\wedge ackIndex'$ = $[ii \in Server \mapsto [ij \in Server \mapsto$
IF $ii = i$ THEN 0
ELSE $ackIndex[ii][ij]]]$

$\wedge committedIndex'$ = $[committedIndex \text{ EXCEPT } ![i] = 0]$
 $\wedge initialHistory'$ = $[initialHistory \text{ EXCEPT } ![i] = \langle \rangle]$
 $\wedge tempMaxEpoch'$ = $[tempMaxEpoch \text{ EXCEPT } ![i] = currentEpoch[i]]$
 $\wedge tempMaxLastEpoch'$ = $[tempMaxLastEpoch \text{ EXCEPT } ![i] = currentEpoch[i]]$
 $\wedge tempInitialHistory'$ = $[tempInitialHistory \text{ EXCEPT } ![i] = history[i]]$
 $\wedge leaderOracle'$ = $[s \in Server \mapsto \text{IF } s \in Q \text{ THEN } i$
ELSE $leaderOracle[s]]$

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

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

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

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

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.

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

The leader finds timeout with another follower.

$LeaderTimeout(i, j) \triangleq$
 $\wedge state[i] \neq Follower$
 $\wedge j \neq i$
 $\wedge j \in cluster[i]$
 $\wedge \text{LET } newCluster \triangleq cluster[i] \setminus \{j\}$
IN $\wedge \vee \wedge newCluster \in Quorums$
 $\wedge cluster' = [cluster \text{ EXCEPT } ![i] = newCluster]$
 $\wedge clean(i, j)$
 $\wedge \text{UNCHANGED } \langle state, cepochRecv, ackeRecv, ackldRecv, ackIndex, committedIndex, initialHistory,$
 $tempMaxEpoch, tempMaxLastEpoch, tempInitialHistory, leaderOracle, leaderEpoch, cepochSent, msgs \rangle$
 $\vee \wedge newCluster \notin Quorums$
 $\wedge \text{LET } Q \triangleq \text{CHOOSE } q \in Quorums : i \in q$

$v \triangleq \text{CHOOSE } s \in Q : \text{TRUE}$
 IN $Election(v, Q)$
 $\exists Q \in Quorums : \wedge i \in Q$
 $\wedge \exists v \in Q : Election(v, Q)$
 $\wedge \text{UNCHANGED } \langle currentEpoch, history, commitIndex, currentCounter, sendCounter, recoveryVars, proposalMsgs \rangle$

A follower finds timeout with the leader.
 $FollowerTimeout(i) \triangleq$
 $\wedge state[i] = Follower$
 $\wedge leaderOracle[i] \neq NullPoint$
 $\wedge \exists Q \in Quorums : \wedge i \in Q$
 $\wedge \exists v \in Q : Election(v, Q)$
 $\wedge \text{UNCHANGED } \langle currentEpoch, history, commitIndex, currentCounter, sendCounter, recoveryVars, proposalMsgs \rangle$

A server halts and restarts.
 $Restart(i) \triangleq$
 $\wedge state' = [state \text{ EXCEPT } ![i] = Follower]$
 $\wedge leaderOracle' = [leaderOracle \text{ EXCEPT } ![i] = NullPoint]$
 $\wedge commitIndex' = [commitIndex \text{ EXCEPT } ![i] = 0]$
 $\wedge cepochSent' = [ceepochSent \text{ EXCEPT } ![i] = \text{FALSE}]$
 $\wedge msgs' = [ii \in Server \mapsto [ij \in Server \mapsto \text{IF } ij = i \text{ THEN } \langle \rangle$
 $\text{ELSE } msgs[ii][ij]]]$
 $\wedge \text{UNCHANGED } \langle currentEpoch, leaderEpoch, history, leaderVars, tempVars, recoveryVars, proposalMsgs \rangle$

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.
 $RecoveryAfterRestart(i) \triangleq$
 $\wedge state[i] = Follower$
 $\wedge leaderOracle[i] = NullPoint$
 $\wedge recoveryRespRecv' = [recoveryRespRecv \text{ EXCEPT } ![i] = \{\}]$
 $\wedge recoveryMaxEpoch' = [recoveryMaxEpoch \text{ EXCEPT } ![i] = currentEpoch[i]]$
 $\wedge recoveryMEOracle' = [recoveryMEOracle \text{ EXCEPT } ![i] = NullPoint]$
 $\wedge BroadcastToAll(i, [mtype \mapsto RECOVERYREQUEST])$
 $\wedge \text{UNCHANGED } \langle serverVars, leaderVars, tempVars, cepochSent, proposalMsgsLog \rangle$

$HandleRecoveryRequest(i, j) \triangleq$
 $\wedge msgs[j][i] \neq \langle \rangle$
 $\wedge msgs[j][i][1].mtype = RECOVERYREQUEST$
 $\wedge Reply(i, j, [mtype \mapsto RECOVERYRESPONSE,$
 $\text{moracle} \mapsto leaderOracle[i],$
 $\text{mepoch} \mapsto currentEpoch[i]])$
 $\wedge \text{UNCHANGED } \langle serverVars, leaderVars, tempVars, cepochSent, recoveryVars, proposalMsgsLog \rangle$

$HandleRecoveryResponse(i, j) \triangleq$
 $\wedge msgs[j][i] \neq \langle \rangle$

$$\begin{aligned}
& \wedge \text{msgs}[j][i][1].\text{mtype} = \text{RECOVERYRESPONSE} \\
& \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \text{infoOk} \triangleq \wedge \text{msg.mepoch} \geq \text{recoveryMaxEpoch}[i] \\
& \quad \quad \wedge \text{msg.moracle} \neq \text{NullPoint} \\
& \text{IN} \quad \vee \wedge \text{infoOk} \\
& \quad \quad \wedge \text{recoveryMaxEpoch}' = [\text{recoveryMaxEpoch} \text{ EXCEPT } ![i] = \text{msg.mepoch}] \\
& \quad \quad \wedge \text{recoveryMEOracle}' = [\text{recoveryMEOracle} \text{ EXCEPT } ![i] = \text{msg.moracle}] \\
& \quad \vee \wedge \neg \text{infoOk} \\
& \quad \quad \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{recoveryRe} \\
& \quad \quad \quad \text{ELSE } \text{recoveryRe} \\
& \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{proposalMsgsLog} \rangle \\
& \text{FindCluster}(i) \triangleq \\
& \quad \wedge \text{state}[i] = \text{Follower} \\
& \quad \wedge \text{leaderOracle}[i] = \text{NullPoint} \\
& \quad \wedge \text{recoveryRespRecv}[i] \in \text{Quorums} \\
& \quad \wedge \text{LET } \text{infoOk} \triangleq \wedge \text{recoveryMEOracle}[i] \neq i \\
& \quad \quad \wedge \text{recoveryMEOracle}[i] \neq \text{NullPoint} \\
& \quad \text{IN} \quad \vee \wedge \neg \text{infoOk} \\
& \quad \quad \text{UNCHANGED } \langle \text{currentEpoch}, \text{leaderOracle}, \text{msgs} \rangle \\
& \quad \vee \wedge \text{infoOk} \\
& \quad \quad \wedge \text{currentEpoch}' = [\text{currentEpoch} \text{ EXCEPT } ![i] = \text{recoveryMaxEpoch}[i]] \\
& \quad \quad \wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = \text{recoveryMEOracle}[i]] \\
& \quad \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 } \langle \text{state}, \text{leaderEpoch}, \text{history}, \text{commitIndex}, \text{leaderVars}, \text{tempVars}, \text{recoveryVars}, \text{ceepoch} \rangle \\
\hline
& \text{In phase } f11, \text{ follower sends } f.p \text{ to pleader via } \text{CEPOCH}. \\
& \text{FollowerDiscovery1}(i) \triangleq \\
& \quad \wedge \text{state}[i] = \text{Follower} \\
& \quad \wedge \text{leaderOracle}[i] \neq \text{NullPoint} \\
& \quad \wedge \neg \text{ceepochSent}[i] \\
& \quad \wedge \text{LET } \text{leader} \triangleq \text{leaderOracle}[i] \\
& \quad \quad \text{IN } \text{Send}(i, \text{leader}, [\text{mtype} \mapsto \text{CEPOCH}, \\
& \quad \quad \quad \text{mepoch} \mapsto \text{currentEpoch}[i]]) \\
& \quad \wedge \text{ceepochSent}' = [\text{ceepochSent} \text{ EXCEPT } ![i] = \text{TRUE}] \\
& \quad \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{leaderVars}, \text{tempVars}, \text{recoveryVars}, \text{proposalMsgsLog} \rangle \\
& \text{In phase } l11, \text{ pleader receives } \text{CEPOCH} \text{ from a quorum, and choose a new epoch } e' \\
& \text{as its own } l.p \text{ and sends } \text{NEWPOCH} \text{ to followers.} \\
& \text{LeaderHandleCEPOCH}(i, j) \triangleq \\
& \quad \wedge \text{state}[i] = \text{ProspectiveLeader} \\
& \quad \wedge \text{msgs}[j][i] \neq \langle \rangle
\end{aligned}$$

$\wedge \text{msgs}[j][i][1].\text{mtype} = \text{CEPOCH}$
 $\wedge \vee$ new message - modify tempMaxEpoch and cepochRecv
 $\wedge \text{NullPoint} \notin \text{cepochRecv}[i]$
 $\wedge \text{LET } \text{newEpoch} \triangleq \text{Maximum}(\{\text{tempMaxEpoch}[i], \text{msgs}[j][i][1].\text{mepoch}\})$
 $\text{IN } \text{tempMaxEpoch}' = [\text{tempMaxEpoch} \text{ EXCEPT } ![i] = \text{newEpoch}]$
 $\wedge \text{cepochRecv}' = [\text{cepochRecv} \text{ EXCEPT } ![i] = \text{IF } j \in \text{cepochRecv}[i] \text{ THEN } \text{cepochRecv}[i] \text{ ELSE } \text{cepochRecv}[i] \cup \{j\}]$
 $\wedge \text{Discard}(j, i)$
 \vee new follower who joins in cluster / follower whose history and commitIndex do not match
 $\wedge \text{NullPoint} \in \text{cepochRecv}[i]$
 $\wedge \vee \wedge \text{NullPoint} \notin \text{ackRecv}[i]$
 $\quad \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{NEWEPOCH},$
 $\quad \quad \text{mepoch} \mapsto \text{leaderEpoch}[i]])$
 $\vee \wedge \text{NullPoint} \in \text{ackRecv}[i]$
 $\quad \wedge \text{Reply2}(i, j, [\text{mtype} \mapsto \text{NEWEPOCH},$
 $\quad \quad \text{mepoch} \mapsto \text{leaderEpoch}[i],$
 $\quad \quad \text{[mtype} \mapsto \text{NEWLEADER},$
 $\quad \quad \text{mepoch} \mapsto \text{currentEpoch}[i],$
 $\quad \quad \text{minitialHistory} \mapsto \text{initialHistory}[i]])$
 $\wedge \text{UNCHANGED } \langle \text{cepochRecv}, \text{tempMaxEpoch} \rangle$
 $\wedge \text{cluster}' = [\text{cluster} \text{ EXCEPT } ![i] = \text{IF } j \in \text{cluster}[i] \text{ THEN } \text{cluster}[i] \text{ ELSE } \text{cluster}[i] \cup \{j\}]$
 $\wedge \text{UNCHANGED } \langle \text{serverVars}, \text{ackRecv}, \text{ackldRecv}, \text{ackIndex}, \text{currentCounter}, \text{sendCounter}, \text{initialHistory},$
 $\quad \text{committedIndex}, \text{cepochSent}, \text{tempMaxLastEpoch}, \text{tempInitialHistory}, \text{recoveryVars}, p \rangle$

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 .

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

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.

$\text{FollowerDiscovery2}(i, j) \triangleq$
 $\wedge \text{state}[i] = \text{Follower}$
 $\wedge \text{msgs}[j][i] \neq \langle \rangle$
 $\wedge \text{msgs}[j][i][1].\text{mtype} = \text{NEWEPOCH}$
 $\wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1]$
 $\text{IN } \vee$ new NEWEPOCH - accept and reply
 $\quad \wedge \text{currentEpoch}[i] < \text{msg.mepoch}$

$$\begin{aligned}
& \wedge \text{currentEpoch}' = [\text{currentEpoch} \text{ EXCEPT } ![i] = \text{msg.mepoch}] \\
& \wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = j] \\
& \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{ACKE}, \\
& \quad \text{mepoch} \mapsto \text{msg.mepoch}, \\
& \quad \text{mlastEpoch} \mapsto \text{leaderEpoch}[i], \\
& \quad \text{mhf} \mapsto \text{history}[i]]) \\
\vee & \wedge \text{currentEpoch}[i] = \text{msg.mepoch} \\
& \wedge \vee \wedge \text{leaderOracle}[i] = j \\
& \quad \wedge \text{Discard}(j, i) \\
& \quad \wedge \text{UNCHANGED } \langle \text{currentEpoch}, \text{leaderOracle} \rangle \\
\vee & \text{It may happen when a leader do not update new epoch to all followers in } Q, \text{ and a new election begins} \\
& \wedge \text{leaderOracle}[i] \neq j \\
& \wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = j] \\
& \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{ACKE}, \\
& \quad \text{mepoch} \mapsto \text{msg.mepoch}, \\
& \quad \text{mlastEpoch} \mapsto \text{leaderEpoch}[i], \\
& \quad \text{mhf} \mapsto \text{history}[i]]) \\
& \quad \wedge \text{UNCHANGED } \text{currentEpoch} \\
\vee & \text{stale NEWEPOCH} - \text{diacard} \\
& \wedge \text{currentEpoch}[i] > \text{msg.mepoch} \\
& \wedge \text{Discard}(j, i) \\
& \quad \wedge \text{UNCHANGED } \langle \text{currentEpoch}, \text{leaderOracle} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{leaderEpoch}, \text{history}, \text{leaderVars}, \text{commitIndex}, \text{ceepochSent}, \text{tempVars}, \text{recover} \rangle
\end{aligned}$$

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.

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

$$\begin{aligned}
& \wedge \text{state}[i] = \text{ProspectiveLeader} \\
& \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \wedge \text{msgs}[j][i][1].\text{mtype} = \text{ACKE} \\
& \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \text{infoOk} \triangleq \vee \text{msg.mlastEpoch} > \text{tempMaxLastEpoch}[i] \\
& \quad \quad \vee \wedge \text{msg.mlastEpoch} = \text{tempMaxLastEpoch}[i] \\
& \quad \quad \quad \wedge \vee \text{LastZxid}(\text{msg.mhf})[1] > \text{LastZxid}(\text{tempInitialHistory}[i])[1] \\
& \quad \quad \quad \vee \wedge \text{LastZxid}(\text{msg.mhf})[1] = \text{LastZxid}(\text{tempInitialHistory}[i])[1] \\
& \quad \quad \quad \quad \wedge \text{LastZxid}(\text{msg.mhf})[2] \geq \text{LastZxid}(\text{tempInitialHistory}[i])[2] \\
\text{IN } & \vee \wedge \text{leaderEpoch}[i] = \text{msg.mepoch} \\
& \quad \wedge \vee \wedge \text{infoOk} \\
& \quad \quad \wedge \text{tempMaxLastEpoch}' = [\text{tempMaxLastEpoch} \text{ EXCEPT } ![i] = \text{msg.mlastEpoch}] \\
& \quad \quad \wedge \text{tempInitialHistory}' = [\text{tempInitialHistory} \text{ EXCEPT } ![i] = \text{msg.mhf}] \\
& \quad \vee \wedge \neg \text{infoOk} \\
& \quad \quad \wedge \text{UNCHANGED } \langle \text{tempMaxLastEpoch}, \text{tempInitialHistory} \rangle
\end{aligned}$$

Followers not in Q will not receive *NEWEPOCH*, so leader will receive *ACKE* only when the source is in Q

$$\begin{aligned}
& \wedge \text{ackRecv}' = [\text{ackRecv} \text{ EXCEPT } ![i] = \text{IF } j \notin \text{ackRecv}[i] \text{ THEN } \text{ackRecv}[i] \cup \{j\} \\
& \quad \quad \quad \text{ELSE } \text{ackRecv}[i]]
\end{aligned}$$

$$\begin{aligned}
& \vee \wedge \text{leaderEpoch}[i] \neq \text{msg.mepoch} \\
& \wedge \text{UNCHANGED } \langle \text{tempMaxLastEpoch}, \text{tempInitialHistory}, \text{ackRecv} \rangle \\
& \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{cluster}, \text{cepochnRecv}, \text{ackldRecv}, \text{ackIndex}, \text{currentCounter}, \\
& \quad \text{sendCounter}, \text{initialHistory}, \text{committedIndex}, \text{cepochnSent}, \text{tempMaxEpoch}, \text{recoveryVars} \rangle \\
\text{LeaderDiscovery2Sync1}(i) & \triangleq \\
& \wedge \text{state}[i] = \text{ProspectiveLeader} \\
& \wedge \text{ackRecv}[i] \in \text{Quorums} \\
& \wedge \text{currentEpoch}' = [\text{currentEpoch} \text{ EXCEPT } ![i] = \text{leaderEpoch}[i]] \\
& \wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{tempInitialHistory}[i]] \\
& \wedge \text{initialHistory}' = [\text{initialHistory} \text{ EXCEPT } ![i] = \text{tempInitialHistory}[i]] \\
& \wedge \text{ackRecv}' = [\text{ackRecv} \text{ EXCEPT } ![i] = \{\text{NullPoint}\}] \\
& \wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i][i] = \text{Len}(\text{tempInitialHistory}[i])] \\
& \text{until now, phase1(Discovery) ends} \\
& \wedge \text{Broadcast}(i, [\text{mtype} \mapsto \text{NEWLEADER}, \\
& \quad \text{mepoch} \mapsto \text{currentEpoch}[i], \\
& \quad \text{minitialHistory} \mapsto \text{history}'[i]]) \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{leaderEpoch}, \text{leaderOracle}, \text{commitIndex}, \text{cluster}, \text{cepochnRecv}, \text{ackldRecv}, \\
& \quad \text{currentCounter}, \text{sendCounter}, \text{committedIndex}, \text{cepochnSent}, \text{tempVars}, \text{recoveryVars} \rangle
\end{aligned}$$

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 *cepochnRecv*, *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.

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

\wedge UNCHANGED $\langle \text{currentEpoch}, \text{leaderEpoch}, \text{leaderOracle}, \text{history} \rangle$
 \wedge UNCHANGED $\langle \text{state}, \text{commitIndex}, \text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, \text{proposalMsgsLd} \rangle$

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

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

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

In phase *f22*, follower receives *COMMIT-LD* and delivers all unprocessed transaction.

$\text{FollowerSync2}(i, j) \triangleq$
 $\wedge \text{state}[i] = \text{Follower}$
 $\wedge \text{msgs}[j][i] \neq \langle \rangle$
 $\wedge \text{msgs}[j][i][1].\text{mtype} = \text{COMMITLD}$
 $\wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1]$
 IN \vee $\text{new COMMIT-LD - commit all transactions in initial history}$
 Regardless of *Restart*, it must be true because one will receive *NEWLEADER* before receiving *COMMIT-LD*
 $\wedge \text{currentEpoch}[i] = \text{msg.mepoch}$
 $\wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = j] \text{ unnecessary}$

$$\begin{aligned}
& \wedge \vee \wedge \text{Len}(\text{history}[i]) = \text{msg.mlength} \\
& \quad \wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{Len}(\text{history}[i])] \\
& \quad \wedge \text{Discard}(j, i) \\
& \vee \wedge \text{Len}(\text{history}[i]) \neq \text{msg.mlength} \\
& \quad \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{CEPOCH}, \\
& \quad \quad \quad \text{mepoch} \mapsto \text{currentEpoch}[i]]) \\
& \quad \wedge \text{UNCHANGED } \text{commitIndex} \\
\vee & \text{ } > : \text{stale } \text{COMMIT-LD} - \text{discard} \\
& < : \text{In our implementation, ' < ' does not exist due to the guarantee of } \text{Restart} \\
& < : \text{If ' < ' exists, we can discard it and handle it in } \text{phase3} \\
& \wedge \text{currentEpoch}[i] \neq \text{msg.mepoch} \\
& \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED } \langle \text{commitIndex}, \text{leaderOracle} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{leaderOracle}, \text{history}, \text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{recover} \rangle
\end{aligned}$$

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

$$\begin{aligned}
& \text{ClientRequest}(i, v) \triangleq \\
& \quad \wedge \text{state}[i] = \text{Leader} \\
& \quad \wedge \text{currentCounter}' = [\text{currentCounter} \text{ EXCEPT } ![i] = \text{currentCounter}[i] + 1] \\
& \quad \wedge \text{LET } \text{newTransaction} \triangleq [\text{epoch} \mapsto \text{currentEpoch}[i], \\
& \quad \quad \quad \text{counter} \mapsto \text{currentCounter}'[i], \\
& \quad \quad \quad \text{value} \mapsto v] \\
& \quad \text{IN } \wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{Append}(\text{history}[i], \text{newTransaction})] \\
& \quad \quad \wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i] = \text{Len}(\text{history}'[i])] \text{ necessary, to push } \text{commitIndex} \\
& \quad \wedge \text{UNCHANGED } \langle \text{msgs}, \text{state}, \text{currentEpoch}, \text{leaderEpoch}, \text{leaderOracle}, \text{commitIndex}, \text{cluster}, \text{ceepochRecv}, \\
& \quad \quad \quad \text{ackldRecv}, \text{sendCounter}, \text{initialHistory}, \text{committedIndex}, \text{tempVars}, \text{ceepochRecv} \rangle
\end{aligned}$$

$$\begin{aligned}
& \text{LeaderBroadcast1}(i) \triangleq \\
& \quad \wedge \text{state}[i] = \text{Leader} \\
& \quad \wedge \text{sendCounter}[i] < \text{currentCounter}[i] \\
& \quad \wedge \text{LET } \text{toBeSentCounter} \triangleq \text{sendCounter}[i] + 1 \\
& \quad \quad \text{toBeSentIndex} \triangleq \text{Len}(\text{initialHistory}[i]) + \text{toBeSentCounter} \\
& \quad \quad \text{toBeSentEntry} \triangleq \text{history}[i][\text{toBeSentIndex}] \\
& \quad \text{IN } \wedge \text{Broadcast}(i, [\text{mtype} \mapsto \text{PROPOSE}, \\
& \quad \quad \quad \text{mepoch} \mapsto \text{currentEpoch}[i], \\
& \quad \quad \quad \text{mproposal} \mapsto \text{toBeSentEntry}]) \\
& \quad \wedge \text{sendCounter}' = [\text{sendCounter} \text{ EXCEPT } ![i] = \text{toBeSentCounter}] \\
& \quad \wedge \text{LET } m \triangleq [\text{msource} \mapsto i, \text{mtype} \mapsto \text{PROPOSE}, \text{mepoch} \mapsto \text{currentEpoch}[i], \text{mproposal} \mapsto \text{toBeSentEntry}] \\
& \quad \quad \text{IN } \text{proposalMsgsLog}' = \text{proposalMsgsLog} \cup \{m\} \\
& \quad \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{ceepochRecv}, \text{cluster}, \text{ackRecv}, \text{ackldRecv}, \text{ackIndex}, \\
& \quad \quad \quad \text{currentCounter}, \text{initialHistory}, \text{committedIndex}, \text{tempVars}, \text{recoveryVars}, \text{ceepochSent} \rangle
\end{aligned}$$

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

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

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

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 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 If happens, \neq must be $>$, namely a stale follower sends it.
 $\wedge \text{currentEpoch}[i] \neq \text{msg.mepoch}$
 $\wedge \text{UNCHANGED } \text{ackIndex}$
 $\wedge \text{Discard}(j, i)$
 $\wedge \text{UNCHANGED } \langle \text{serverVars}, \text{cluster}, \text{ceepochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{currentCounter},$
 $\text{sendCounter}, \text{initialHistory}, \text{committedIndex}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, \text{proposalMsgsLog} \rangle$

$\text{LeaderAdvanceCommit}(i) \triangleq$
 $\wedge \text{state}[i] = \text{Leader}$
 $\wedge \text{commitIndex}[i] < \text{Len}(\text{history}[i])$
 $\wedge \text{LET } \text{Agree}(\text{index}) \triangleq \{i\} \cup \{k \in (\text{Server} \setminus \{i\}) : \text{ackIndex}[i][k] \geq \text{index}\}$
 $\text{agreeIndexes} \triangleq \{\text{index} \in (\text{commitIndex}[i] + 1) \dots \text{Len}(\text{history}[i]) : \text{Agree}(\text{index}) \in \text{Quorum}\}$
 $\text{newCommitIndex} \triangleq \text{IF } \text{agreeIndexes} \neq \{\} \text{ THEN } \text{Maximum}(\text{agreeIndexes})$
 $\text{ELSE } \text{commitIndex}[i]$
 IN $\text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{newCommitIndex}]$
 $\wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{leaderEpoch}, \text{leaderOracle}, \text{history},$
 $\text{msgs}, \text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, \text{proposalMsgsLog} \rangle$

$\text{LeaderBroadcast2}(i) \triangleq$

$$\begin{aligned}
& \wedge \text{state}[i] = \text{Leader} \\
& \wedge \text{committedIndex}[i] < \text{commitIndex}[i] \\
& \wedge \text{LET } \text{newCommittedIndex} \triangleq \text{committedIndex}[i] + 1 \\
& \quad \text{IN } \wedge \text{Broadcast}(i, [\text{mtype} \mapsto \text{COMMIT}, \\
& \quad \quad \text{mepoch} \mapsto \text{currentEpoch}[i], \\
& \quad \quad \text{mindex} \mapsto \text{newCommittedIndex}, \\
& \quad \quad \text{mcounter} \mapsto \text{history}[\text{newCommittedIndex}].\text{counter}]) \\
& \quad \wedge \text{committedIndex}' = [\text{committedIndex} \text{ EXCEPT } ![i] = \text{committedIndex}[i] + 1] \\
& \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{cluster}, \text{ceepochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{ackIndex}, \text{currentCounter}, \\
& \quad \text{sendCounter}, \text{initialHistory}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, \text{proposalMsgsLog} \rangle
\end{aligned}$$

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{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \quad \text{IN } \vee \wedge \text{currentEpoch}[i] = \text{msg.mepoch} \\
& \quad \quad \quad \wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = j] \\
& \quad \quad \quad \wedge \text{LET } \text{infoOk} \triangleq \wedge \text{Len}(\text{history}[i]) \geq \text{msg.mindex} \\
& \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 \text{IN } \vee \text{new COMMIT} - \text{commit transaction in history} \\
& \quad \quad \quad \quad \wedge \text{infoOk} \\
& \quad \quad \quad \quad \wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{Maximum}(\{\text{commitIndex}[i], \text{msg.mindex}\})] \\
& \quad \quad \quad \quad \wedge \text{Discard}(j, i) \\
& \quad \quad \quad \vee \text{It may happen when the server is a new follower who joined in the cluster,} \\
& \quad \quad \quad \quad \text{and it misses the corresponding PROPOSE.} \\
& \quad \quad \quad \quad \wedge \neg \text{infoOk} \\
& \quad \quad \quad \quad \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{CEPOCH}, \\
& \quad \quad \quad \quad \quad \text{mepoch} \mapsto \text{currentEpoch}[i]]) \\
& \quad \quad \quad \quad \wedge \text{UNCHANGED } \text{commitIndex} \\
& \quad \vee \text{stale COMMIT} - \text{discard} \\
& \quad \quad \wedge \text{currentEpoch}[i] \neq \text{msg.mepoch} \\
& \quad \quad \wedge \text{Discard}(j, i) \\
& \quad \quad \wedge \text{UNCHANGED } \langle \text{commitIndex}, \text{leaderOracle} \rangle \\
& \quad \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{leaderEpoch}, \text{history}, \\
& \quad \quad \text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, \text{proposalMsgsLog} \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*.
 $LeaderHandleCEPOCHinPhase3(i, j) \triangleq$

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

In phase *l34*, upon receiving ack from *f* of the *NEWLEADER*, it sends a commit message to *f*.

Leader *l* also makes $Q := Q \cup \{f\}$.

$LeaderHandleACKLDinPhase3(i, j) \triangleq$

$$\begin{aligned} & \wedge state[i] = Leader \\ & \wedge msgs[j][i] \neq \langle \rangle \\ & \wedge msgs[j][i][1].mtype = ACKLD \\ & \wedge LET \ msg \triangleq \ msgs[j][i][1] \\ & \quad aimCommitIndex \triangleq \ Minimum(\{commitIndex[i], Len(msg.mhistory)\}) \\ & \quad IN \quad \vee \wedge currentEpoch[i] = msg.mepoch \\ & \quad \quad \wedge ackIndex' = [ackIndex \text{ EXCEPT } ![i][j] = Len(msg.mhistory)] \\ & \quad \quad \wedge Reply(i, j, [mtype \mapsto COMMIT, \\ & \quad \quad \quad mepoch \mapsto currentEpoch[i], \\ & \quad \quad \quad mindex \mapsto aimCommitIndex, \\ & \quad \quad \quad mcounter \mapsto history[aimCommitIndex].counter]) \\ & \quad \vee \wedge currentEpoch[i] \neq msg.mepoch \\ & \quad \quad \wedge Discard(j, i) \\ & \quad \quad \wedge UNCHANGED \ ackIndex \\ & \wedge cluster' = [cluster \text{ EXCEPT } ![i] = \text{IF } j \in cluster[i] \text{ THEN } cluster[i] \\ & \quad \quad \quad \text{ELSE } cluster[i] \cup \{j\}] \\ & \wedge UNCHANGED \langle serverVars, cepochRecv, ackRecv, ackldRecv, currentCounter, sendCounter, \\ & \quad \quad \quad initialHistory, committedIndex, tempVars, cepochSent, recoveryVars, proposalMsgsLog \rangle \end{aligned}$$

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.

$BecomeFollower(i) \triangleq$

$$\begin{aligned}
& \wedge state[i] \neq \text{Follower} \\
& \wedge \exists j \in \text{Server} \setminus \{i\} : \wedge msgs[j][i] \neq \langle \rangle \\
& \quad \wedge msgs[j][i][1].mtype \neq \text{RECOVERYREQUEST} \\
& \quad \wedge msgs[j][i][1].mtype \neq \text{RECOVERYRESPONSE} \\
& \quad \wedge \text{LET } msg \triangleq msgs[j][i][1] \\
& \quad \text{IN} \quad \wedge currentEpoch[i] < msg.mepoch \\
& \quad \quad \wedge \vee msg.mtype = \text{NEWPOCH} \\
& \quad \quad \vee msg.mtype = \text{NEWLEADER} \\
& \quad \quad \vee msg.mtype = \text{COMMITLD} \\
& \quad \quad \vee msg.mtype = \text{PROPOSE} \\
& \quad \quad \vee msg.mtype = \text{COMMIT} \\
& \quad \wedge state' = [state \text{ EXCEPT } ![i] = \text{Follower}] \\
& \quad \wedge currentEpoch' = [currentEpoch \text{ EXCEPT } ![i] = msg.mepoch] \\
& \quad \wedge leaderOracle' = [leaderOracle \text{ EXCEPT } ![i] = j] \\
& \quad \text{Here we should not use } Discard.
\end{aligned}$$

$\wedge \text{UNCHANGED } \langle leaderEpoch, history, commitIndex, msgs, leaderVars, tempVars, cepochSent, recovery \rangle$

$DiscardStaleMessage(i) \triangleq$

$$\begin{aligned}
& \wedge \exists j \in \text{Server} \setminus \{i\} : \wedge msgs[j][i] \neq \langle \rangle \\
& \quad \wedge msgs[j][i][1].mtype \neq \text{RECOVERYREQUEST} \\
& \quad \wedge msgs[j][i][1].mtype \neq \text{RECOVERYRESPONSE} \\
& \quad \wedge \text{LET } msg \triangleq msgs[j][i][1] \\
& \quad \text{IN} \quad \vee \wedge state[i] = \text{Follower} \\
& \quad \quad \wedge \vee msg.mepoch < currentEpoch[i] \setminus * \text{ Discussed before.} \\
& \quad \quad \vee msg.mtype = \text{CEPOCH} \\
& \quad \quad \vee msg.mtype = \text{ACKE} \\
& \quad \quad \vee msg.mtype = \text{ACKLD} \\
& \quad \quad \vee msg.mtype = \text{ACK} \\
& \quad \vee \wedge state[i] = \text{Leader} \\
& \quad \quad \wedge msg.mtype \neq \text{CEPOCH} \\
& \quad \quad \wedge \vee msg.mepoch \leq currentEpoch[i] \\
& \quad \quad \quad \vee msg.mtype = \text{ACKE} \text{ response of NEWPOCH} \\
& \quad \vee \wedge state[i] = \text{ProspectiveLeader} \\
& \quad \quad \wedge msg.mtype \neq \text{CEPOCH} \\
& \quad \quad \wedge \vee msg.mepoch \leq currentEpoch[i] \\
& \quad \quad \quad \vee msg.mtype = \text{ACK} \\
& \wedge Discard(j, i)
\end{aligned}$$

$\wedge \text{UNCHANGED } \langle serverVars, leaderVars, tempVars, cepochSent, recoveryVars, proposalMsgsLog \rangle$

Defines how the variables may transition.

$Next \triangleq$

$\vee \exists i \in Server, Q \in Quorums : InitialElection(i, Q)$
 $\vee \exists i \in Server : Restart(i)$
 $\vee \exists i \in Server : RecoveryAfterRestart(i)$
 $\vee \exists i, j \in Server : HandleRecoveryRequest(i, j)$
 $\vee \exists i, j \in Server : HandleRecoveryResponse(i, j)$
 $\vee \exists i, j \in Server : FindCluster(i)$
 $\vee \exists i, j \in Server : LeaderTimeout(i, j)$
 $\vee \exists i \in Server : FollowerTimeout(i)$
 $\vee \exists i \in Server : FollowerDiscovery1(i)$
 $\vee \exists i, j \in Server : LeaderHandleCEPOCH(i, j)$
 $\vee \exists i \in Server : LeaderDiscovery1(i)$
 $\vee \exists i, j \in Server : FollowerDiscovery2(i, j)$
 $\vee \exists i, j \in Server : LeaderHandleACKE(i, j)$
 $\vee \exists i \in Server : LeaderDiscovery2Sync1(i)$
 $\vee \exists i, j \in Server : FollowerSync1(i, j)$
 $\vee \exists i, j \in Server : LeaderHandleACKLD(i, j)$
 $\vee \exists i \in Server : LeaderSync2(i)$
 $\vee \exists i, j \in Server : FollowerSync2(i, j)$
 $\vee \exists i \in Server, v \in Value : ClientRequest(i, v)$
 $\vee \exists i \in Server : LeaderBroadcast1(i)$
 $\vee \exists i, j \in Server : FollowerBroadcast1(i, j)$
 $\vee \exists i, j \in Server : LeaderHandleACK(i, j)$
 $\vee \exists i \in Server : LeaderAdvanceCommit(i)$
 $\vee \exists i \in Server : LeaderBroadcast2(i)$
 $\vee \exists i, j \in Server : FollowerBroadcast2(i, j)$
 $\vee \exists i, j \in Server : LeaderHandleCEPOCHinPhase3(i, j)$
 $\vee \exists i, j \in Server : LeaderHandleACKLDinPhase3(i, j)$
 $\vee \exists i \in Server : DiscardStaleMessage(i)$
 $\vee \exists i \in Server : BecomeFollower(i)$

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

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.

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

$\wedge state[i] = Leader \vee state[i] = ProspectiveLeader$
 $\wedge state[j] = Leader \vee state[j] = ProspectiveLeader$

$$\begin{aligned} & \wedge \text{currentEpoch}[i] = \text{currentEpoch}[j] \\ & \Rightarrow i = j \end{aligned}$$

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.

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

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.

$$\begin{aligned} \text{Integrity} & \triangleq \forall i \in \text{Server} : \\ & \text{state}[i] = \text{Follower} \wedge \text{commitIndex}[i] > 0 \\ & \Rightarrow \forall \text{index} \in 1 \dots \text{commitIndex}[i] : \exists \text{msg} \in \text{proposalMsgsLog} : \\ & \text{equal}(\text{msg.mproposal}, \text{history}[i][\text{index}]) \end{aligned}$$

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

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

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

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

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

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

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

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

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

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

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

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 *Mon Apr 26 21:48:14 CST 2021* by Dell
\ * Created Sat *Dec 05 13:32:08 CST 2020* by Dell