

---

MODULE *Zab*

---

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, 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 *ackRecv*

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*

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*

VARIABLE *recoverySent*

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*  $\triangleq$   $\langle recoveryRespRecv, recoveryMaxEpoch, recoveryMEOracle, recoverySent \rangle$

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

---

*LastZxid(his)*  $\triangleq$  IF *Len(his)* > 0 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* EXCEPT ![*i*][*j*] = *Append(msgs[i][j], m)*]

*Send2(i, j, m1, m2)*  $\triangleq$  *msgs'* = [*msgs* 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'* = IF *msgs[i][j]*  $\neq \langle \rangle$  THEN [*msgs* EXCEPT ![*i*][*j*] = *Tail(msgs[i][j])*]  
 ELSE *msgs*

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

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

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

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

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

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

---

Define initial values for all variables

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

---

A server becomes pleader and a quorum servers knows that.

$$\begin{aligned}
Election(i, Q) &\triangleq \\
&\wedge i \in Q \\
&\wedge state' \\
&= [s \in Server \mapsto \text{IF } s = i \text{ THEN } ProspectiveLeader \\
&\quad \text{ELSE IF } s \in Q \text{ THEN } Follower \\
&\quad \text{ELSE } state[s]] \\
&\wedge cluster' \\
&= [cluster \text{ EXCEPT } ![i] = Q] \quad \text{cluster is first initialized in election, not } phase1. \\
&\wedge cepochRecv' \\
&= [cepochRecv \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge ackeRecv' \\
&= [ackeRecv \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge ackldRecv' \\
&= [ackldRecv \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge ackIndex' \\
&= [ii \in Server \mapsto [ij \in Server \mapsto \\
&\quad \text{IF } ii = i \text{ THEN } 0 \\
&\quad \text{ELSE } ackIndex[ii][ij]]] \\
&\wedge committedIndex' \\
&= [committedIndex \text{ EXCEPT } ![i] = 0] \\
&\wedge initialHistory' \\
&= [initialHistory \text{ EXCEPT } ![i] = \langle \rangle] \\
&\wedge tempMaxEpoch' \\
&= [tempMaxEpoch \text{ EXCEPT } ![i] = currentEpoch[i]] \\
&\wedge tempMaxLastEpoch' \\
&= [tempMaxLastEpoch \text{ EXCEPT } ![i] = currentEpoch[i]] \\
&\wedge tempInitialHistory' \\
&= [tempInitialHistory \text{ EXCEPT } ![i] = history[i]] \\
&\wedge leaderOracle' \\
&= [s \in Server \mapsto \text{IF } s \in Q \text{ THEN } i \\
&\quad \text{ELSE } leaderOracle[s]] \\
&\wedge leaderEpoch' \\
&= [s \in Server \mapsto \text{IF } s \in Q \text{ THEN } currentEpoch[s] \\
&\quad \text{ELSE } leaderEpoch[s]] \\
&\wedge cepochSent' \\
&= [s \in Server \mapsto \text{IF } s \in Q \text{ THEN } FALSE \\
&\quad \text{ELSE } cepochSent[s]] \\
&\wedge msgs' \\
&= [ii \in Server \mapsto [ij \in Server \mapsto \\
&\quad \text{IF } ii \in Q \vee ij \in Q \text{ THEN } \langle \rangle \\
&\quad \text{ELSE } msgs[ii][ij]]]
\end{aligned}$$

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.

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

The leader finds timeout with another follower.

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

$$\begin{aligned}
& \vee \wedge \text{newCluster} \notin \text{Quorums} \\
& \wedge \text{LET } Q \triangleq \text{CHOOSE } q \in \text{Quorums}: i \in q \\
& \quad v \triangleq \text{CHOOSE } s \in Q: \text{TRUE} \\
& \text{IN } \text{Election}(v, Q) \\
& \exists Q \in \text{Quorums} : \wedge i \in Q \\
& \quad \wedge \exists v \in Q : \text{Election}(v, Q) \\
& \wedge \text{UNCHANGED } \langle \text{currentEpoch}, \text{history}, \text{commitIndex}, \text{currentCounter}, \text{sendCounter}, \text{recoveryVars}, \text{proposals} \rangle
\end{aligned}$$

A follower finds timeout with the leader.

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


---

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.

$$\begin{aligned}
\text{Restart}(i) & \triangleq \\
& \wedge \text{state}' = [\text{state} \text{ EXCEPT } ![i] = \text{Follower}] \\
& \wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = \text{NullPoint}] \\
& \wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = 0] \\
& \wedge \text{ceepochSent}' = [\text{ceepochSent} \text{ EXCEPT } ![i] = \text{FALSE}] \\
& \wedge \text{msgs}' = [ii \in \text{Server} \mapsto [ij \in \text{Server} \mapsto \text{IF } ij = i \text{ THEN } \langle \rangle \\
& \quad \text{ELSE } \text{msgs}[ii][ij]]] \\
& \wedge \text{recoverySent}' = [\text{recoverySent} \text{ EXCEPT } ![i] = \text{FALSE}] \\
& \wedge \text{UNCHANGED } \langle \text{currentEpoch}, \text{leaderEpoch}, \text{history}, \text{leaderVars}, \text{tempVars}, \\
& \quad \text{recoveryRespRecv}, \text{recoveryMaxEpoch}, \text{recoveryMEOracle}, \text{proposalMsgsLog} \rangle
\end{aligned}$$

$$\begin{aligned}
\text{RecoveryAfterRestart}(i) & \triangleq \\
& \wedge \text{state}[i] = \text{Follower} \\
& \wedge \text{leaderOracle}[i] = \text{NullPoint} \\
& \wedge \neg \text{recoverySent}[i] \\
& \wedge \text{recoveryRespRecv}' = [\text{recoveryRespRecv} \text{ EXCEPT } ![i] = \{\}] \\
& \wedge \text{recoveryMaxEpoch}' = [\text{recoveryMaxEpoch} \text{ EXCEPT } ![i] = \text{currentEpoch}[i]] \\
& \wedge \text{recoveryMEOracle}' = [\text{recoveryMEOracle} \text{ EXCEPT } ![i] = \text{NullPoint}] \\
& \wedge \text{recoverySent}' = [\text{recoverySent} \text{ EXCEPT } ![i] = \text{TRUE}] \\
& \wedge \text{BroadcastToAll}(i, [\text{mtype} \mapsto \text{RECOVERYREQUEST}]) \\
& \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{proposalMsgsLog} \rangle
\end{aligned}$$

$$\begin{aligned}
\text{HandleRecoveryRequest}(i, j) & \triangleq \\
& \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \wedge \text{msgs}[j][i][1].\text{mtype} = \text{RECOVERYREQUEST} \\
& \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{RECOVERYRESPONSE}], \text{msgs}[j][i][2..])
\end{aligned}$$

$$\begin{aligned}
& \text{moracle} \mapsto \text{leaderOracle}[i], \\
& \text{mepoch} \mapsto \text{currentEpoch}[i]) \\
& \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{leaderVars}, \text{tempVars}, \text{cePOCHSent}, \text{recoveryVars}, \text{proposalMsgsLog} \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] \\
& \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 \wedge \text{recoveryMaxEpoch}' = [\text{recoveryMaxEpoch} \text{ EXCEPT } ![i] = \text{msg.mepoch}] \\
& \quad \wedge \text{recoveryMEOracle}' = [\text{recoveryMEOracle} \text{ EXCEPT } ![i] = \text{msg.moracle}] \\
& \quad \vee \wedge \neg \text{infoOk} \\
& \quad \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{recoveryRe} \\
& \quad \quad \quad \text{ELSE } \text{recoveryRe} \\
& \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{leaderVars}, \text{tempVars}, \text{cePOCHSent}, \text{recoverySent}, \text{proposalMsgsLog} \rangle \\
\text{FindCluster}(i) & \triangleq \\
& \wedge \text{state}[i] = \text{Follower} \\
& \wedge \text{leaderOracle}[i] = \text{NullPoint} \\
& \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  $\text{CEPOCH}$ .

$$\begin{aligned}
\text{FollowerDiscovery1}(i) & \triangleq \\
& \wedge \text{state}[i] = \text{Follower} \\
& \wedge \text{leaderOracle}[i] \neq \text{NullPoint} \\
& \wedge \neg \text{cePOCHSent}[i] \\
& \wedge \text{LET } \text{leader} \triangleq \text{leaderOracle}[i]
\end{aligned}$$

IN  $Send(i, leader, [mtype \mapsto CEPOCH,$   
 $mepoch \mapsto currentEpoch[i]])$   
 $\wedge cepochSent' = [ceepochSent \text{ EXCEPT } ![i] = \text{TRUE}]$   
 $\wedge \text{UNCHANGED } \langle serverVars, leaderVars, tempVars, recoveryVars, proposalMsgsLog \rangle$

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

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

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*.

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



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.

$FollowerDiscovery2(i, j) \triangleq$

- $\wedge state[i] = Follower$
- $\wedge msgs[j][i] \neq \langle \rangle$
- $\wedge msgs[j][i][1].mtype = NEWEPOCH$
- $\wedge LET\ msg \triangleq msgs[j][i][1]$
- IN  $\vee$   $\text{new } NEWEPOCH - \text{accept and reply}$ 
  - $\wedge currentEpoch[i] < msg.mepoch$
  - $\wedge currentEpoch' = [currentEpoch\ EXCEPT\ ![i] = msg.mepoch]$
  - $\wedge leaderOracle' = [leaderOracle\ EXCEPT\ ![i] = j]$
  - $\wedge Reply(i, j, [mtype \mapsto ACKE,$ 
    - $mepoch \mapsto msg.mepoch,$
    - $mlastEpoch \mapsto leaderEpoch[i],$
    - $mhf \mapsto history[i]])$
- $\vee \wedge currentEpoch[i] = msg.mepoch$
- $\wedge \vee \wedge leaderOracle[i] = j$ 
  - $\wedge Reply(i, j, [mtype \mapsto ACKE,$ 
    - $mepoch \mapsto msg.mepoch,$
    - $mlastEpoch \mapsto leaderEpoch[i],$
    - $mhf \mapsto history[i]])$
  - $\wedge UNCHANGED \langle currentEpoch, leaderOracle \rangle$
- $\vee$   $\text{It may happen when a leader do not update new epoch to all followers in } Q, \text{ and a new election begins}$ 
  - $\wedge leaderOracle[i] \neq j$
  - $\wedge leaderOracle' = [leaderOracle\ EXCEPT\ ![i] = j]$
  - $\wedge Reply(i, j, [mtype \mapsto ACKE,$ 
    - $mepoch \mapsto msg.mepoch,$
    - $mlastEpoch \mapsto leaderEpoch[i],$
    - $mhf \mapsto history[i]])$
  - $\wedge UNCHANGED\ currentEpoch$
- $\vee$   $\text{stale } NEWEPOCH - \text{discard}$ 
  - $\wedge currentEpoch[i] > msg.mepoch$
  - $\wedge Discard(j, i)$
  - $\wedge UNCHANGED \langle currentEpoch, leaderOracle \rangle$

$\wedge UNCHANGED \langle state, leaderEpoch, history, leaderVars, commitIndex, cepochSent, tempVars, recovery \rangle$

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.

$LeaderHandleACKE(i, j) \triangleq$

- $\wedge state[i] = ProspectiveLeader$
- $\wedge msgs[j][i] \neq \langle \rangle$
- $\wedge msgs[j][i][1].mtype = ACKE$
- $\wedge LET\ msg \triangleq msgs[j][i][1]$
- $infoOk \triangleq \vee msg.mlastEpoch > tempMaxLastEpoch[i]$
- $\vee \wedge msg.mlastEpoch = tempMaxLastEpoch[i]$

$$\begin{aligned}
& \wedge \vee LastZxid(msg.mhf)[1] > LastZxid(tempInitialHistory[i])[1] \\
& \vee \wedge LastZxid(msg.mhf)[1] = LastZxid(tempInitialHistory[i])[1] \\
& \wedge LastZxid(msg.mhf)[2] \geq LastZxid(tempInitialHistory[i])[2] \\
IN \quad & \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 UNCHANGED \langle tempMaxLastEpoch, tempInitialHistory \rangle \\
& \quad \text{Followers not in } Q \text{ will not receive } NEWEPOCH, \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 UNCHANGED \langle tempMaxLastEpoch, tempInitialHistory, ackeRecv \rangle \\
& \wedge Discard(j, i) \\
& \wedge UNCHANGED \langle serverVars, cluster, cepochRecv, ackldRecv, ackIndex, currentCounter, \\
& \quad \quad \quad sendCounter, initialHistory, committedIndex, cepochSent, tempMaxEpoch, recoveryVars \rangle \\
LeaderDiscovery2Sync1(i) & \triangleq \\
& \wedge state[i] = ProspectiveLeader \\
& \wedge ackeRecv[i] \in Quorums \\
& \wedge currentEpoch' = [currentEpoch \text{ EXCEPT } ![i] = leaderEpoch[i]] \\
& \wedge history' = [history \text{ EXCEPT } ![i] = tempInitialHistory[i]] \\
& \wedge initialHistory' = [initialHistory \text{ EXCEPT } ![i] = tempInitialHistory[i]] \\
& \wedge ackeRecv' = [ackeRecv \text{ EXCEPT } ![i] = \{NullPoint\}] \\
& \wedge ackIndex' = [ackIndex \text{ EXCEPT } ![i][i] = Len(tempInitialHistory[i])] \\
& \text{until now, phase1(Discovery) ends} \\
& \wedge Broadcast(i, [mtype \mapsto NEWLEADER, \\
& \quad \quad \quad mepoch \mapsto currentEpoch'[i], \\
& \quad \quad \quad minitialHistory \mapsto history'[i]]) \\
& \wedge UNCHANGED \langle state, leaderEpoch, leaderOracle, commitIndex, cluster, cepochRecv, ackldRecv, \\
& \quad \quad \quad currentCounter, sendCounter, committedIndex, cepochSent, tempVars, recoveryVars,
\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 *cepochRecv*, *ackeRecv*, *ackldRecv* to  $\{NullPoint\}$  in corresponding three actions to make sure that the prospective leader will not broadcast *NEWEPOCH*/*NEWLEADER*/*COMMITLD* twice.

---

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

$$\begin{aligned}
FollowerSync1(i, j) & \triangleq \\
& \wedge state[i] = Follower \\
& \wedge msgs[j][i] \neq \langle \rangle \\
& \wedge msgs[j][i][1].mtype = NEWLEADER \\
& \wedge LET msg \triangleq msgs[j][i][1]
\end{aligned}$$

IN  $\vee$  **new NEWLEADER – accept and reply**  
 $\wedge \text{currentEpoch}[i] \leq \text{msg.mepoch}$   
 $\wedge \text{currentEpoch}' = [\text{currentEpoch} \text{ EXCEPT } ![i] = \text{msg.mepoch}]$   
 $\wedge \text{leaderEpoch}' = [\text{leaderEpoch} \text{ EXCEPT } ![i] = \text{msg.mepoch}]$   
 $\wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = j]$   
 $\wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{msg.minitialHistory}]$   
 $\wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{ACKLD},$   
 $\text{mepoch} \mapsto \text{msg.mepoch},$   
 $\text{mhistory} \mapsto \text{msg.minitialHistory}])$   
 $\vee$  **stale NEWLEADER – discard**  
 $\wedge \text{currentEpoch}[i] > \text{msg.mepoch}$   
 $\wedge \text{Discard}(j, i)$   
 $\wedge \text{UNCHANGED } \langle \text{currentEpoch}, \text{leaderEpoch}, \text{leaderOracle}, \text{history} \rangle$   
 $\wedge \text{UNCHANGED } \langle \text{state}, \text{commitIndex}, \text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, \text{proposalMsgsLd} \rangle$

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

$\text{LeaderSync2}(i) \triangleq$   
 $\wedge \text{state}[i] = \text{ProspectiveLeader}$   
 $\wedge \text{ackldRecv}[i] \in \text{Quorums}$   
 $\wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{Len}(\text{history}[i])]$   
 $\wedge \text{committedIndex}' = [\text{committedIndex} \text{ EXCEPT } ![i] = \text{Len}(\text{history}[i])]$   
 $\wedge \text{state}' = [\text{state} \text{ EXCEPT } ![i] = \text{Leader}]$   
 $\wedge \text{currentCounter}' = [\text{currentCounter} \text{ EXCEPT } ![i] = 0]$   
 $\wedge \text{sendCounter}' = [\text{sendCounter} \text{ EXCEPT } ![i] = 0]$   
 $\wedge \text{ackldRecv}' = [\text{ackldRecv} \text{ EXCEPT } ![i] = \{\text{NullPoint}\}]$   
 $\wedge \text{Broadcast}(i, [\text{mtype} \mapsto \text{COMMITLD},$   
 $\text{mepoch} \mapsto \text{currentEpoch}[i],$   
 $\text{mlength} \mapsto \text{Len}(\text{history}[i])])$

$\wedge$  UNCHANGED  $\langle \text{currentEpoch}, \text{leaderEpoch}, \text{leaderOracle}, \text{history}, \text{cluster}, \text{ceepochRecv}, \text{ackRecv}, \text{ackIndex}, \text{initialHistory}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, \text{proposalM}$

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  $\vee$  new COMMIT-LD - commit all transactions in initial history  
 $\wedge \text{currentEpoch}[i] = \text{msg.mepoch}$   
 $\wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = j] \text{ unnecessary}$   
 $\wedge \vee \wedge \text{Len}(\text{history}[i]) = \text{msg.mlength}$   
 $\wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{Len}(\text{history}[i])]$   
 $\wedge \text{Discard}(j, i)$   
 $\vee \wedge \text{Len}(\text{history}[i]) \neq \text{msg.mlength}$   
 $\wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{CEPOCH}, \text{mepoch} \mapsto \text{currentEpoch}[i]])$   
 $\wedge$  UNCHANGED  $\text{commitIndex}$   
 $\vee$   $> : \text{stale COMMIT-LD - discard}$   
 $< : \text{If } ' < ' \text{ exists, we can discard it and handle it in } \text{phase3}$   
 $\wedge \text{currentEpoch}[i] \neq \text{msg.mepoch}$   
 $\wedge \text{Discard}(j, i)$   
 $\wedge$  UNCHANGED  $\langle \text{commitIndex}, \text{leaderOracle} \rangle$   
 $\wedge$  UNCHANGED  $\langle \text{state}, \text{currentEpoch}, \text{leaderEpoch}, \text{history}, \text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{recovery}$

---

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

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

$\text{LeaderBroadcast1}(i) \triangleq$   
 $\wedge \text{state}[i] = \text{Leader}$   
 $\wedge \text{sendCounter}[i] < \text{currentCounter}[i]$   
 $\wedge \text{LET } \text{toBeSentCounter} \triangleq \text{sendCounter}[i] + 1$   
 $\text{toBeSentIndex} \triangleq \text{Len}(\text{initialHistory}[i]) + \text{toBeSentCounter}$   
 $\text{toBeSentEntry} \triangleq \text{history}[i][\text{toBeSentIndex}]$   
 IN  $\wedge \text{Broadcast}(i, [\text{mtype} \mapsto \text{PROPOSE},$

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

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

$$\begin{aligned}
& \text{FollowerBroadcast1}(i, j) \triangleq \\
& \quad \wedge \text{state}[i] = \text{Follower} \\
& \quad \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{PROPOSE} \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \text{IN } \quad \vee \quad \text{It should be that } \vee \text{msg.mproposal.counter} = 1 \\
& \quad \quad \vee \text{msg.mproposal.counter} = \text{history}[\text{Len}(\text{history})].\text{counter} + 1 \\
& \quad \quad \wedge \text{currentEpoch}[i] = \text{msg.mepoch} \\
& \quad \quad \wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{Append}(\text{history}[i], \text{msg.mproposal})] \\
& \quad \quad \wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = j] \\
& \quad \quad \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{ACK}, \\
& \quad \quad \quad \text{mepoch} \mapsto \text{currentEpoch}[i], \\
& \quad \quad \quad \text{mindex} \mapsto \text{Len}(\text{history}'[i])]) \\
& \quad \vee \quad \text{If happens, } \neq \text{ must be } >, \text{ namely a stale leader sends it.} \\
& \quad \quad \wedge \text{currentEpoch}[i] \neq \text{msg.mepoch} \\
& \quad \quad \wedge \text{Discard}(j, i) \\
& \quad \quad \wedge \text{UNCHANGED } \langle \text{history}, \text{leaderOracle} \rangle \\
& \quad \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{leaderEpoch}, \text{commitIndex}, \text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{recoveryIndex} \rangle
\end{aligned}$$

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

$$\begin{aligned}
& \text{LeaderHandleACK}(i, j) \triangleq \\
& \quad \wedge \text{state}[i] = \text{Leader} \\
& \quad \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{ACK} \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \text{IN } \quad \vee \quad \text{It should be that } \text{ackIndex}[i][j] + 1 \triangleq \text{msg.mindex} \\
& \quad \quad \wedge \text{currentEpoch}[i] = \text{msg.mepoch} \\
& \quad \quad \wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i][j] = \text{Maximum}(\{\text{ackIndex}[i][j], \text{msg.mindex}\})] \\
& \quad \vee \quad \text{If happens, } \neq \text{ must be } >, \text{ namely a stale follower sends it.} \\
& \quad \quad \wedge \text{currentEpoch}[i] \neq \text{msg.mepoch} \\
& \quad \quad \wedge \text{UNCHANGED } \text{ackIndex} \\
& \quad \wedge \text{Discard}(j, i) \\
& \quad \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{cluster}, \text{ceepochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{currentCounter}, \\
& \quad \quad \text{sendCounter}, \text{initialHistory}, \text{committedIndex}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, \text{recoveryIndex} \rangle
\end{aligned}$$

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

$$\begin{aligned}
& \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}\} \\
& \quad \text{agreeIndexes} \triangleq \{\text{index} \in (\text{commitIndex}[i] + 1) \dots \text{Len}(\text{history}[i]) : \text{Agree}(\text{index}) \in \text{Quorum}\} \\
& \quad \text{newCommitIndex} \triangleq \text{IF } \text{agreeIndexes} \neq \{\} \text{ THEN } \text{Maximum}(\text{agreeIndexes}) \\
& \quad \quad \quad \text{ELSE } \text{commitIndex}[i] \\
& \text{IN } \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{newCommitIndex}] \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{leaderEpoch}, \text{leaderOracle}, \text{history}, \\
& \quad \text{msgs}, \text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, \text{proposalMsgsLog} \rangle \\
\text{LeaderBroadcast2}(i) & \triangleq \\
& \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}[i][\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 \\
\text{In phase } f32, \text{ follower receives } \text{COMMIT} \text{ and commits transaction.} \\
\text{FollowerBroadcast2}(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{COMMIT} \\
& \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \text{IN } \vee \wedge \text{currentEpoch}[i] = \text{msg.mepoch} \\
& \quad \quad \wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = j] \\
& \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 \wedge \text{history}[i][\text{msg.mindex}].\text{epoch} = \text{msg.mepoch} \\
& \quad \quad \quad \wedge \text{history}[i][\text{msg.mindex}].\text{counter} = \text{msg.mcounter} \\
& \quad \quad \quad \vee \text{msg.mindex} = 0 \\
& \quad \text{IN } \vee \text{new COMMIT} - \text{commit transaction in history} \\
& \quad \quad \wedge \text{infoOk} \\
& \quad \quad \wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{Maximum}(\{\text{commitIndex}[i], \text{msg.mindex}\})] \\
& \quad \quad \wedge \text{Discard}(j, i) \\
& \quad \vee \text{It may happen when the server is a new follower who joined in the cluster,} \\
& \quad \quad \text{and it misses the corresponding PROPOSE.} \\
& \quad \quad \wedge \neg \text{infoOk} \\
& \quad \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}
\end{aligned}$$

$$\begin{aligned}
& \vee \text{ stale } COMMIT - \text{discard} \\
& \wedge \text{currentEpoch}[i] \neq \text{msg.mepoch} \\
& \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED } \langle \text{commitIndex}, \text{leaderOracle} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{leaderEpoch}, \text{history}, \\
& \quad \text{leaderVars}, \text{tempVars}, \text{ceepochSent}, \text{recoveryVars}, \text{proposalMsgsLog} \rangle
\end{aligned}$$

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.

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

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

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

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

$\text{LeaderHandleACKLDinPhase3}(i, j) \triangleq$

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

ELSE  $cluster[i] \cup \{j\}$   
 $\wedge$  UNCHANGED  $\langle serverVars, cepochRecv, ackRecv, ackldRecv, currentCounter, sendCounter, initialHistory, committedIndex, tempVars, cepochSent, recoveryVars, proposalMsgsLo,$

Let me suppose three 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  $\wedge Maximum(\{currentEpoch[i], leaderEpoch[i]\}) < msg.mepoch$   
 $\wedge \vee msg.mtype = NEWEPOCH$   
 $\vee msg.mtype = NEWLEADER$   
 $\vee msg.mtype = COMMITLD$   
 $\vee msg.mtype = PROPOSE$   
 $\vee msg.mtype = COMMIT$   
 $\wedge state' = [state \text{ EXCEPT } ![i] = Follower]$   
 $\wedge currentEpoch' = [currentEpoch \text{ EXCEPT } ![i] = msg.mepoch]$   
 $\wedge leaderOracle' = [leaderOracle \text{ EXCEPT } ![i] = j]$   
 Here we should not use *Discard*.  
 $\wedge$  UNCHANGED  $\langle leaderEpoch, history, commitIndex, msgs, leaderVars, tempVars, cepochSent, recovery,$

---

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



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

---

Defines how the variables may transition.

$\text{Next} \triangleq$

$$\begin{aligned}
& \vee \exists i \in \text{Server}, Q \in \text{Quorums} : \text{InitialElection}(i, Q) \\
& \vee \exists i \in \text{Server} : \text{Restart}(i) \\
& \vee \exists i \in \text{Server} : \text{RecoveryAfterRestart}(i) \\
& \vee \exists i, j \in \text{Server} : \text{HandleRecoveryRequest}(i, j) \\
& \vee \exists i, j \in \text{Server} : \text{HandleRecoveryResponse}(i, j) \\
& \vee \exists i, j \in \text{Server} : \text{FindCluster}(i) \\
& \vee \exists i, j \in \text{Server} : \text{LeaderTimeout}(i, j) \\
& \vee \exists i \in \text{Server} : \text{FollowerTimeout}(i) \\
& \vee \exists i \in \text{Server} : \text{FollowerDiscovery1}(i) \\
& \vee \exists i, j \in \text{Server} : \text{LeaderHandleCEPOCH}(i, j) \\
& \vee \exists i \in \text{Server} : \text{LeaderDiscovery1}(i) \\
& \vee \exists i, j \in \text{Server} : \text{FollowerDiscovery2}(i, j) \\
& \vee \exists i, j \in \text{Server} : \text{LeaderHandleACKE}(i, j) \\
& \vee \exists i \in \text{Server} : \text{LeaderDiscovery2Sync1}(i) \\
& \vee \exists i, j \in \text{Server} : \text{FollowerSync1}(i, j) \\
& \vee \exists i, j \in \text{Server} : \text{LeaderHandleACKLD}(i, j) \\
& \vee \exists i \in \text{Server} : \text{LeaderSync2}(i) \\
& \vee \exists i, j \in \text{Server} : \text{FollowerSync2}(i, j) \\
& \vee \exists i \in \text{Server}, v \in \text{Value} : \text{ClientRequest}(i, v) \\
& \vee \exists i \in \text{Server} : \text{LeaderBroadcast1}(i) \\
& \vee \exists i, j \in \text{Server} : \text{FollowerBroadcast1}(i, j) \\
& \vee \exists i, j \in \text{Server} : \text{LeaderHandleACK}(i, j) \\
& \vee \exists i \in \text{Server} : \text{LeaderAdvanceCommit}(i) \\
& \vee \exists i \in \text{Server} : \text{LeaderBroadcast2}(i) \\
& \vee \exists i, j \in \text{Server} : \text{FollowerBroadcast2}(i, j)
\end{aligned}$$

$$\begin{aligned}
& \forall \exists i, j \in \text{Server} : \text{LeaderHandleCEPOCHinPhase3}(i, j) \\
& \forall \exists i, j \in \text{Server} : \text{LeaderHandleACKLDinPhase3}(i, j) \\
& \forall \exists i \in \text{Server} : \text{DiscardStaleMessage}(i) \\
& \forall \exists i \in \text{Server} : \text{BecomeFollower}(i)
\end{aligned}$$

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

---

Safety properties of Zab consensus algorithm

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

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

$$\wedge \vee \text{state}[i] = \text{Leader}$$

$$\vee \wedge \text{state}[i] = \text{ProspectiveLeader}$$

$$\wedge \text{NullPoint} \in \text{ackRecv}[i] \quad \text{prospective leader determines its epoch after broadcasting NEWLE}$$

$$\wedge \vee \text{state}[j] = \text{Leader}$$

$$\vee \wedge \text{state}[j] = \text{ProspectiveLeader}$$

$$\wedge \text{NullPoint} \in \text{ackRecv}[j]$$

$$\wedge \text{currentEpoch}[i] = \text{currentEpoch}[j]$$

$$\Rightarrow i = j$$

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

So we can assume the set of delivered transactions is the prefix of history with index from 1 to  $\text{commitIndex}$ .

And 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.

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

$$\wedge idx1 < idx2$$

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

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

---

\ \* Modification History  
 \ \* Last modified *Thu Apr 29 17:16:17 CST 2021* by Dell  
 \ \* Created Sat *Dec 05 13:32:08 CST 2020* by Dell