



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 followers who has successfully sent *CEPOCH* to pleader in pleader.

VARIABLE *cepochRecv*

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* to represent the prefix entries of the history

that has applied to state machine, but we can tolerate that *applyIndex* = *commitIndex*.

This does not violate correctness.

VARIABLE *commitIndex*

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

VARIABLE *committedIndex*

Helper matrix for follower to stop sending *CEPOCH* to pleader in followers.

Because *CEPOCH* is the sole message which follower actively sends to pleader.

VARIABLE *cepochSent*

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*

VARIABLE *tempInitialHistory*

*serverVars*  $\triangleq$   $\langle \text{state}, \text{currentEpoch}, \text{leaderEpoch}, \text{leaderOracle}, \text{history}, \text{commitIndex} \rangle$

*leaderVars*  $\triangleq$   $\langle \text{cepocheRecv}, \text{ackRecv}, \text{ackldRecv}, \text{ackIndex}, \text{currentCounter}, \text{sendCounter}, \text{initialHistory}, \text{commitIndex} \rangle$

*tempVars*  $\triangleq$   $\langle \text{tempMaxEpoch}, \text{tempMaxLastEpoch}, \text{tempInitialHistory} \rangle$

*vars*  $\triangleq$   $\langle \text{serverVars}, \text{msgs}, \text{leaderVars}, \text{tempVars}, \text{cepocheSent} \rangle$

---

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

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

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

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

---

Define initial values for all variables

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

$$\begin{aligned}
\wedge \text{cephochSent} &= [s \in \text{Server} \mapsto \text{FALSE}] \\
\wedge \text{tempMaxEpoch} &= [s \in \text{Server} \mapsto 0] \\
\wedge \text{tempMaxLastEpoch} &= [s \in \text{Server} \mapsto 0] \\
\wedge \text{tempInitialHistory} &= [s \in \text{Server} \mapsto \langle \rangle]
\end{aligned}$$

---

A server becomes pleader and a quorum servers knows that.

$$\begin{aligned}
\text{Election}(i, Q) &\triangleq \\
&\wedge i \in Q \\
&\wedge \text{state}' = [s \in \text{Server} \mapsto \text{IF } s = i \text{ THEN } \text{ProspectiveLeader} \\
&\hspace{15em} \text{ELSE IF } s \in Q \text{ THEN } \text{Follower} \\
&\hspace{15em} \text{ELSE } \text{state}[s]] \\
&\wedge \text{cephochRecv}' = [\text{cephochRecv} \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge \text{ackRecv}' = [\text{ackRecv} \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge \text{ackldRecv}' = [\text{ackldRecv} \text{ EXCEPT } ![i] = \{i\}] \\
&\wedge \text{ackIndex}' = [ii \in \text{Server} \mapsto [ij \in \text{Server} \mapsto \\
&\hspace{10em} \text{IF } ii = i \text{ THEN } 0 \\
&\hspace{10em} \text{ELSE } \text{ackIndex}[ii][ij]]] \\
&\wedge \text{committedIndex}' = [\text{committedIndex} \text{ EXCEPT } ![i] = 0] \\
&\wedge \text{initialHistory}' = [\text{initialHistory} \text{ EXCEPT } ![i] = \langle \rangle] \\
&\wedge \text{tempMaxEpoch}' = [\text{tempMaxEpoch} \text{ EXCEPT } ![i] = \text{currentEpoch}[i]] \\
&\wedge \text{tempMaxLastEpoch}' = [\text{tempMaxLastEpoch} \text{ EXCEPT } ![i] = \text{currentEpoch}[i]] \\
&\wedge \text{tempInitialHistory}' = [\text{tempInitialHistory} \text{ EXCEPT } ![i] = \text{history}[i]] \\
&\wedge \text{leaderOracle}' = [s \in \text{Server} \mapsto \text{IF } s \in Q \text{ THEN } i \\
&\hspace{10em} \text{ELSE } \text{leaderOracle}[s]] \\
&\wedge \text{leaderEpoch}' = [s \in \text{Server} \mapsto \text{IF } s \in Q \text{ THEN } \text{currentEpoch}[s] \\
&\hspace{10em} \text{ELSE } \text{leaderEpoch}[s]] \\
&\wedge \text{cephochSent}' = [s \in \text{Server} \mapsto \text{IF } s \in Q \text{ THEN } \text{FALSE} \\
&\hspace{10em} \text{ELSE } \text{cephochSent}[s]] \\
&\wedge \text{msgs}' = [ii \in \text{Server} \mapsto [ij \in \text{Server} \mapsto \\
&\hspace{10em} \text{IF } ii \in Q \wedge ij \in Q \text{ THEN } \langle \rangle \\
&\hspace{10em} \text{ELSE } \text{msgs}[ii][ij]]] \\
&\wedge \text{UNCHANGED } \langle \text{currentEpoch}, \text{history}, \text{commitIndex}, \text{currentCounter}, \text{sendCounter} \rangle
\end{aligned}$$

A server halts and restarts.

$$\begin{aligned}
\text{Restart}(i) &\triangleq \\
&\wedge \text{state}' = [\text{state} \text{ EXCEPT } ![i] = \text{Follower}] \\
&\wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = \text{NullPoint}] \\
&\wedge \text{cephochSent}' = [\text{cephochSent} \text{ EXCEPT } ![i] = \text{FALSE}] \\
&\wedge \text{UNCHANGED } \langle \text{currentEpoch}, \text{leaderEpoch}, \text{history}, \text{commitIndex}, \text{leaderVars}, \text{tempVars}, \text{msgs} \rangle
\end{aligned}$$

---

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

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

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

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

$$\begin{aligned}
& \text{LeaderHandleCEPOCH}(i, j) \triangleq \\
& \quad \wedge \text{state}[i] = \text{ProspectiveLeader} \\
& \quad \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{CEPOCH} \\
& \quad \wedge \vee \text{redundant message - just discard} \\
& \quad \quad \wedge j \in \text{cepochRecv}[i] \\
& \quad \quad \wedge \text{UNCHANGED } \langle \text{tempMaxEpoch}, \text{cepochRecv} \rangle \\
& \quad \vee \text{new message - modify tempMaxEpoch and cepochRecv} \\
& \quad \quad \wedge j \notin \text{cepochRecv}[i] \\
& \quad \quad \wedge \text{LET } \text{newEpoch} \triangleq \text{Maximum}(\{\text{tempMaxEpoch}[i], \text{msgs}[j][i][1].\text{mepoch}\}) \\
& \quad \quad \quad \text{IN } \text{tempMaxEpoch}' = [\text{tempMaxEpoch} \text{ EXCEPT } ![i] = \text{newEpoch}] \\
& \quad \quad \quad \wedge \text{cepochRecv}' = [\text{cepochRecv} \text{ EXCEPT } ![i] = \text{cepochRecv}[i] \cup \{j\}] \\
& \quad \wedge \text{Discard}(j, i) \\
& \quad \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{ackRecv}, \text{ackldRecv}, \text{ackIndex}, \text{currentCounter}, \text{sendCounter}, \\
& \quad \quad \text{initialHistory}, \text{committedIndex}, \text{cepochSent}, \text{tempMaxLastEpoch}, \text{tempInitialHistory} \rangle
\end{aligned}$$

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

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

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

$$\begin{aligned}
& \text{FollowerDiscovery2}(i, j) \triangleq \\
& \quad \wedge \text{state}[i] = \text{Follower} \\
& \quad \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{NEWEPOCH} \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \text{IN } \vee \text{new NEWEPOCH - accept and reply} \\
& \quad \quad \wedge \text{currentEpoch}[i] \leq \text{msg.mepoch} \quad \text{Here use } \leq, \text{ because one follower may send CEPOCH more then o}
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{currentEpoch}' = [\text{currentEpoch} \text{ EXCEPT } ![i] = \text{msg.mepoch}] \\
& \wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = j] \\
& \wedge \text{Reply}(i, j, [\text{mtype} \mapsto \text{ACKE}, \\
& \quad \text{mepoch} \mapsto \text{msg.mepoch}, \\
& \quad \text{mlastEpoch} \mapsto \text{leaderEpoch}[i], \\
& \quad \text{mhf} \mapsto \text{history}[i]]) \\
& \vee \text{stale NEWEPOCH} - \text{diacard} \\
& \wedge \text{currentEpoch}[i] > \text{msg.mepoch} \\
& \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED} \langle \text{currentEpoch}, \text{leaderOracle} \rangle \\
& \wedge \text{UNCHANGED} \langle \text{state}, \text{leaderEpoch}, \text{history}, \text{leaderVars}, \text{commitIndex}, \text{ceepochSent}, \text{tempVars} \rangle
\end{aligned}$$

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

$$\begin{aligned}
& \text{LeaderHandleACKE}(i, j) \triangleq \\
& \quad \wedge \text{state}[i] = \text{ProspectiveLeader} \\
& \quad \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{msgs}[j][i][1].\text{mtype} = \text{ACKE} \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \quad \text{infoOk} \triangleq \vee \text{msg.mlastEpoch} > \text{tempMaxLastEpoch}[i] \\
& \quad \quad \quad \vee \wedge \text{msg.mlastEpoch} = \text{tempMaxLastEpoch}[i] \\
& \quad \quad \quad \quad \vee \wedge \text{LastZxid}(\text{msg.mhf})[1] > \text{LastZxid}(\text{tempInitialHistory}[i])[1] \\
& \quad \quad \quad \quad \vee \wedge \text{LastZxid}(\text{msg.mhf})[1] = \text{LastZxid}(\text{tempInitialHistory}[i])[1] \\
& \quad \quad \quad \quad \quad \wedge \text{LastZxid}(\text{msg.mhf})[2] \geq \text{LastZxid}(\text{tempInitialHistory}[i])[2] \\
& \quad \text{IN} \quad \vee \wedge \text{leaderEpoch}[i] = \text{msg.mepoch} \\
& \quad \quad \wedge \vee \wedge \text{infoOk} \\
& \quad \quad \quad \wedge \text{tempMaxLastEpoch}' = [\text{tempMaxLastEpoch} \text{ EXCEPT } ![i] = \text{msg.mlastEpoch}] \\
& \quad \quad \quad \wedge \text{tempInitialHistory}' = [\text{tempInitialHistory} \text{ EXCEPT } ![i] = \text{msg.mhf}] \\
& \quad \quad \vee \wedge \neg \text{infoOk} \\
& \quad \quad \quad \wedge \text{UNCHANGED} \langle \text{tempMaxLastEpoch}, \text{tempInitialHistory} \rangle \\
& \quad \quad \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 \quad \text{ELSE } \text{ackRecv}[i]] \\
& \quad \vee \wedge \text{leaderEpoch}[i] \neq \text{msg.mepoch} \\
& \quad \quad \wedge \text{UNCHANGED} \langle \text{tempMaxLastEpoch}, \text{tempInitialHistory}, \text{ackRecv} \rangle \\
& \quad \wedge \text{Discard}(j, i) \\
& \quad \wedge \text{UNCHANGED} \langle \text{serverVars}, \text{ceepochRecv}, \text{ackldRecv}, \text{ackIndex}, \text{currentCounter}, \\
& \quad \quad \text{sendCounter}, \text{initialHistory}, \text{committedIndex}, \text{ceepochSent}, \text{tempMaxEpoch} \rangle
\end{aligned}$$

$$\begin{aligned}
& \text{LeaderDiscovery2Sync1}(i) \triangleq \\
& \quad \wedge \text{state}[i] = \text{ProspectiveLeader} \\
& \quad \wedge \text{ackRecv}[i] \in \text{Quorums} \\
& \quad \wedge \text{currentEpoch}' = [\text{currentEpoch} \text{ EXCEPT } ![i] = \text{leaderEpoch}[i]] \\
& \quad \wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{tempInitialHistory}[i]] \\
& \quad \wedge \text{initialHistory}' = [\text{initialHistory} \text{ EXCEPT } ![i] = \text{tempInitialHistory}[i]] \\
& \quad \wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = 0]
\end{aligned}$$

$$\begin{aligned}
& \wedge ackeRecv' = [ackeRecv \text{ EXCEPT } ![i] = \{\}] \\
& \wedge ackIndex' = [ackIndex \text{ EXCEPT } ![i] = Len(tempInitialHistory[i])] \\
& \text{until now, phase1(Discovery) ends} \\
& \wedge Broadcast(i, [mtype \mapsto NEWLEADER, \\
& \quad mepoch \mapsto currentEpoch[i], \\
& \quad minitialHistory \mapsto history'[i]]) \\
& \wedge UNCHANGED \langle state, leaderEpoch, leaderOracle, cepochRecv, ackldRecv, \\
& \quad currentCounter, sendCounter, committedIndex, cepochSent, tempVars \rangle
\end{aligned}$$


---

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

$$\begin{aligned}
FollowerSync1(i, j) & \triangleq \\
& \wedge state[i] = Follower \\
& \wedge msgs[j][i] \neq \langle \rangle \\
& \wedge msgs[j][i][1].mtype = NEWLEADER \\
& \wedge LET \text{ msg } \triangleq msgs[j][i][1] \\
& \text{IN } \vee \text{ new NEWLEADER - accept and reply} \\
& \quad \wedge currentEpoch[i] \leq msg.mepoch \\
& \quad \wedge currentEpoch' = [currentEpoch \text{ EXCEPT } ![i] = msg.mepoch] \\
& \quad \wedge leaderEpoch' = [leaderEpoch \text{ EXCEPT } ![i] = msg.mepoch] \\
& \quad \wedge leaderOracle' = [leaderOracle \text{ EXCEPT } ![i] = j] \\
& \quad \wedge history' = [history \text{ EXCEPT } ![i] = msg.minitialHistory] \\
& \quad \wedge commitIndex' = [commitIndex \text{ EXCEPT } ![i] = 0] \\
& \quad \wedge Reply(i, j, [mtype \mapsto ACKLD, \\
& \quad \quad mepoch \mapsto msg.mepoch, \\
& \quad \quad mhistory \mapsto msg.minitialHistory]) \\
& \vee \text{ stale NEWLEADER - discard} \\
& \quad \wedge currentEpoch[i] > msg.mepoch \\
& \quad \wedge Discard(j, i) \\
& \quad \wedge UNCHANGED \langle currentEpoch, leaderEpoch, leaderOracle, history, commitIndex \rangle \\
& \wedge UNCHANGED \langle state, leaderVars, tempVars, cepochSent \rangle
\end{aligned}$$

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

$$\begin{aligned}
LeaderHandleACKLD(i, j) & \triangleq \\
& \wedge state[i] = ProspectiveLeader \\
& \wedge msgs[j][i] \neq \langle \rangle \\
& \wedge msgs[j][i][1].mtype = ACKLD \\
& \wedge LET \text{ msg } \triangleq msgs[j][i][1] \\
& \text{IN } \vee \text{ new ACK-LD - accept} \\
& \quad \wedge currentEpoch[i] = msg.mepoch \\
& \quad \wedge ackIndex' = [ackIndex \text{ EXCEPT } ![i][j] = Len(initialHistory[i])] \\
& \quad \wedge ackldRecv' = [ackldRecv \text{ EXCEPT } ![i] = \text{IF } j \notin ackldRecv[i] \text{ THEN } ackldRecv[i] \cup \{j\} \\
& \quad \quad \quad \text{ELSE } ackldRecv[i}]] \\
& \vee \text{ stale ACK-LD - impossible}
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{currentEpoch}[i] \neq \text{msg.mepoch} \\
& \wedge \text{UNCHANGED } \langle \text{ackldRecv}, \text{ackIndex} \rangle \\
& \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{ceepochRecv}, \text{ackRecv}, \text{currentCounter}, \\
& \quad \text{sendCounter}, \text{initialHistory}, \text{committedIndex}, \text{tempVars}, \text{ceepochSent} \rangle \\
\text{LeaderSync2}(i) & \triangleq \\
& \wedge \text{state}[i] = \text{ProspectiveLeader} \\
& \wedge \text{ackldRecv}[i] \in \text{Quorums} \\
& \wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{Len}(\text{history}[i])] \\
& \wedge \text{committedIndex}' = [\text{committedIndex} \text{ EXCEPT } ![i] = \text{Len}(\text{history}[i])] \\
& \wedge \text{state}' = [\text{state} \text{ EXCEPT } ![i] = \text{Leader}] \\
& \wedge \text{currentCounter}' = [\text{currentCounter} \text{ EXCEPT } ![i] = 0] \\
& \wedge \text{sendCounter}' = [\text{sendCounter} \text{ EXCEPT } ![i] = 0] \\
& \wedge \text{ackldRecv}' = [\text{ackldRecv} \text{ EXCEPT } ![i] = \{\}] \\
& \wedge \text{Broadcast}(i, [\text{mtype} \mapsto \text{COMMITLD}, \\
& \quad \text{mepoch} \mapsto \text{currentEpoch}[i]]) \\
& \wedge \text{UNCHANGED } \langle \text{currentEpoch}, \text{leaderEpoch}, \text{leaderOracle}, \text{history}, \text{ceepochRecv}, \\
& \quad \text{ackRecv}, \text{ackIndex}, \text{initialHistory}, \text{tempVars}, \text{ceepochSent} \rangle
\end{aligned}$$

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

$$\begin{aligned}
\text{FollowerSync2}(i, j) & \triangleq \\
& \wedge \text{state}[i] = \text{Follower} \\
& \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \wedge \text{msgs}[j][i][1].\text{mtype} = \text{COMMITLD} \\
& \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \text{IN } \vee \text{new COMMIT-LD - commit all transactions in initial history} \\
& \quad \wedge \text{currentEpoch}[i] = \text{msg.mepoch} \\
& \quad \wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{Len}(\text{history}[i])] \\
& \quad \wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = j] \\
& \vee \text{stale COMMIT-LD - discard} \\
& \quad \wedge \text{currentEpoch}[i] \neq \text{msg.mepoch} \\
& \quad \wedge \text{UNCHANGED } \langle \text{commitIndex}, \text{leaderOracle} \rangle \\
& \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{leaderOracle}, \text{history}, \text{leaderVars}, \text{tempVars}, \text{ceepochSent} \rangle
\end{aligned}$$


---

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

$$\begin{aligned}
\text{ClientRequest}(i, v) & \triangleq \\
& \wedge \text{state}[i] = \text{Leader} \\
& \wedge \text{currentCounter}' = [\text{currentCounter} \text{ EXCEPT } ![i] = \text{currentCounter}[i] + 1] \\
& \wedge \text{LET } \text{newTransaction} \triangleq [\text{epoch} \mapsto \text{currentEpoch}[i], \\
& \quad \text{counter} \mapsto \text{currentCounter}'[i], \\
& \quad \text{value} \mapsto v] \\
& \text{IN } \wedge \text{history}' = [\text{history} \text{ EXCEPT } ![i] = \text{Append}(\text{history}[i], \text{newTransaction})] \\
& \wedge \text{ackIndex}' = [\text{ackIndex} \text{ EXCEPT } ![i] = \text{Len}(\text{history}'[i])]
\end{aligned}$$



$\wedge$  UNCHANGED  $\langle msgs, state, currentEpoch, leaderEpoch, leaderOracle, commitIndex, cepochRecv, ackeRecv, ackldRecv, sendCounter, initialHistory, committedIndex, tempVars, cepoch$

$LeaderBroadcast1(i) \triangleq$

$\wedge state[i] = Leader$

$\wedge sendCounter[i] < currentCounter[i]$

$\wedge$  LET  $toBeSentCounter \triangleq sendCounter[i] + 1$

$toBeSentIndex \triangleq Len(initialHistory[i]) + toBeSentCounter$

$toBeSentEntry \triangleq history[i][toBeSentIndex]$

IN  $\wedge Broadcast(i, [mtype \mapsto PROPOSE,$   
 $mepoch \mapsto currentEpoch[i],$   
 $mproposal \mapsto toBeSentEntry])$

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

$\wedge$  UNCHANGED  $\langle serverVars, cepochRecv, ackeRecv, ackldRecv, ackIndex,$   
 $currentCounter, initialHistory, committedIndex, tempVars, cepochSent \rangle$

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

$FollowerBroadcast1(i, j) \triangleq$

$\wedge state[i] = Follower$

$\wedge msgs[j][i] \neq \langle \rangle$

$\wedge msgs[j][i][1].mtype = PROPOSE$

$\wedge$  LET  $msg \triangleq msgs[j][i][1]$

IN  $\vee$  It should be that  $msg.mproposal.counter = 1 \vee msg.mproposal.counter = history[Len(history)].counter + 1$

$\wedge currentEpoch[i] = msg.mepoch$

$\wedge history' = [history \text{ EXCEPT } ![i] = Append(history[i], msg.mproposal)]$

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

$\wedge Reply(i, j, [mtype \mapsto ACK,$   
 $mepoch \mapsto currentEpoch[i],$   
 $mindex \mapsto Len(history'[i]))$

$\vee$  If happens,  $\neq$  must be  $>$ , namely a stale leader sends it.

$\wedge currentEpoch[i] \neq msg.mepoch$

$\wedge Discard(j, i)$

$\wedge$  UNCHANGED  $\langle history, leaderOracle \rangle$

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

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

$LeaderHandleACK(i, j) \triangleq$

$\wedge state[i] = Leader$

$\wedge msgs[j][i] \neq \langle \rangle$

$\wedge msgs[j][i][1].mtype = ACK$

$\wedge$  LET  $msg \triangleq msgs[j][i][1]$

IN  $\vee$  There should be  $ackIndex[i][j] + 1 \triangleq msg.mindex$

$\wedge currentEpoch[i] = msg.mepoch$

$\wedge ackIndex' = [ackIndex \text{ EXCEPT } ![i][j] = Maximum(\{ackIndex[i][j], msg.mindex\})]$

$\vee$  If happens,  $\neq$  must be  $>$ , namely a stale follower sends it.

$$\begin{aligned}
& \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{ceepochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{currentCounter}, \\
& \quad \text{sendCounter}, \text{initialHistory}, \text{committedIndex}, \text{tempVars}, \text{ceepochSent} \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} : \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} \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 \\
& \text{IN } \wedge \text{Broadcast}(i, [\text{mtype} \mapsto \text{COMMIT}, \\
& \quad \text{mepoch} \mapsto \text{currentEpoch}[i], \\
& \quad \text{mindex} \mapsto \text{newCommittedIndex}, \\
& \quad \text{mcounter} \mapsto \text{history}[\text{newCommittedIndex}].\text{counter}]) \\
& \wedge \text{committedIndex}' = [\text{committedIndex} \text{ EXCEPT } ![i] = \text{committedIndex}[i] + 1] \\
& \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{ceepochRecv}, \text{ackRecv}, \text{ackldRecv}, \text{ackIndex}, \text{currentCounter}, \\
& \quad \text{sendCounter}, \text{initialHistory}, \text{tempVars}, \text{ceepochSent} \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] \\
& \text{IN } \vee \text{new } \text{COMMIT} - \text{commit transaction in history} \\
& \quad \wedge \text{currentEpoch}[i] = \text{msg.mepoch} \\
& \quad \wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = \text{Maximum}(\{\text{commitIndex}[i], \text{msg.mindex}\})] \\
& \quad \wedge \text{leaderOracle}' = [\text{leaderOracle} \text{ EXCEPT } ![i] = j] \\
& \vee \text{stale } \text{COMMIT} - \text{discard} \\
& \quad \wedge \text{currentEpoch}[i] \neq \text{msg.mepoch} \\
& \quad \wedge \text{UNCHANGED } \langle \text{commitIndex}, \text{leaderOracle} \rangle \\
& \wedge \text{Discard}(j, i) \\
& \wedge \text{UNCHANGED } \langle \text{state}, \text{currentEpoch}, \text{leaderEpoch}, \text{history}, \\
& \quad \text{leaderVars}, \text{tempVars}, \text{ceepochSent} \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 *NEWLEADER* and *COMMIT*.

In phase *l33*, upon receiving *CEPOCH*, leader *l* proposes back *NEWEPOCH* and *NEWLEADER*.

$LeaderHandleCEPOCHinPhase3(i, j) \triangleq$

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

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

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

$LeaderHandleACKLDinPhase3(i, j) \triangleq$

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

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.

$$\begin{aligned}
& \text{BecomeFollower}(i) \triangleq \\
& \quad \wedge \exists j \in \text{Server} \setminus \{i\} : \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \text{IN } \wedge \text{currentEpoch}[i] < \text{msg.mepoch} \\
& \quad \quad \wedge \vee \text{msg.mtype} = \text{NEWEPOCH} \\
& \quad \quad \vee \text{msg.mtype} = \text{NEWLEADER} \\
& \quad \quad \vee \text{msg.mtype} = \text{COMMITLD} \\
& \quad \quad \vee \text{msg.mtype} = \text{PROPOSE} \\
& \quad \quad \vee \text{msg.mtype} = \text{COMMIT} \\
& \quad \quad \wedge \text{state}' = [\text{state} \quad \text{EXCEPT } ![i] = \text{Follower}] \\
& \quad \quad \wedge \text{currentEpoch}' = [\text{currentEpoch} \quad \text{EXCEPT } ![i] = \text{msg.mepoch}] \\
& \quad \quad \wedge \text{leaderOracle}' = [\text{leaderOracle} \quad \text{EXCEPT } ![i] = j] \\
& \quad \quad \text{Here we should not use Discard.} \\
& \quad \wedge \text{UNCHANGED } \langle \text{leaderEpoch}, \text{history}, \text{commitIndex}, \text{msgs}, \text{leaderVars}, \text{tempVars}, \text{ceepochSent} \rangle
\end{aligned}$$

---


$$\begin{aligned}
& \text{DiscardStaleMessage}(i) \triangleq \\
& \quad \wedge \exists j \in \text{Server} \setminus \{i\} : \wedge \text{msgs}[j][i] \neq \langle \rangle \\
& \quad \wedge \text{LET } \text{msg} \triangleq \text{msgs}[j][i][1] \\
& \quad \text{IN } \vee \wedge \text{state}[i] = \text{Follower} \\
& \quad \quad \wedge \vee \text{msg.mepoch} < \text{currentEpoch}[i] \\
& \quad \quad \vee \text{msg.mtype} = \text{CEPOCH} \\
& \quad \quad \vee \text{msg.mtype} = \text{ACKE} \\
& \quad \quad \vee \text{msg.mtype} = \text{ACKLD} \\
& \quad \quad \vee \text{msg.mtype} = \text{ACK} \\
& \quad \vee \wedge \text{state}[i] = \text{Leader} \\
& \quad \quad \wedge \text{msg.mtype} \neq \text{CEPOCH} \\
& \quad \quad \wedge \vee \text{msg.mepoch} < \text{currentEpoch}[i] \\
& \quad \quad \quad \vee \text{msg.mtype} = \text{ACKE} \quad \text{response of NEWEPOCH} \\
& \quad \vee \wedge \text{state}[i] = \text{ProspectiveLeader} \\
& \quad \quad \wedge \text{msg.mtype} \neq \text{CEPOCH} \\
& \quad \quad \wedge \vee \text{msg.mepoch} < \text{currentEpoch}[i] \\
& \quad \quad \quad \vee \text{msg.mtype} = \text{ACK} \\
& \quad \wedge \text{Discard}(j, i) \\
& \quad \wedge \text{UNCHANGED } \langle \text{serverVars}, \text{leaderVars}, \text{tempVars}, \text{ceepochSent} \rangle
\end{aligned}$$

---

Defines how the variables may transition.

$Next \triangleq$

$$\begin{aligned}
& \vee \exists i \in Server : Restart(i) \\
& \vee \exists i \in Server, Q \in Quorums : Election(i, Q) \\
& \vee \exists i \in Server : FollowerDiscovery1(i) \\
& \vee \exists i, j \in Server : LeaderHandleCEPOCH(i, j) \\
& \vee \exists i \in Server : LeaderDiscovery1(i) \\
& \vee \exists i, j \in Server : FollowerDiscovery2(i, j) \\
& \vee \exists i, j \in Server : LeaderHandleACKE(i, j) \\
& \vee \exists i \in Server : LeaderDiscovery2Sync1(i) \\
& \vee \exists i, j \in Server : FollowerSync1(i, j) \\
& \vee \exists i, j \in Server : LeaderHandleACKLD(i, j) \\
& \vee \exists i \in Server : LeaderSync2(i) \\
& \vee \exists i, j \in Server : FollowerSync2(i, j) \\
& \vee \exists i \in Server, v \in Value : ClientRequest(i, v) \\
& \vee \exists i \in Server : LeaderBroadcast1(i) \\
& \vee \exists i, j \in Server : FollowerBroadcast1(i, j) \\
& \vee \exists i, j \in Server : LeaderHandleACK(i, j) \\
& \vee \exists i \in Server : LeaderAdvanceCommit(i) \\
& \vee \exists i \in Server : LeaderBroadcast2(i) \\
& \vee \exists i, j \in Server : FollowerBroadcast2(i, j) \\
& \vee \exists i, j \in Server : LeaderHandleCEPOCHinPhase3(i, j) \\
& \vee \exists i, j \in Server : LeaderHandleACKLDinPhase3(i, j) \\
& \vee \exists i \in Server : DiscardStaleMessage(i) \\
& \vee \exists i \in Server : BecomeFollower(i)
\end{aligned}$$

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

---

Defines some variants, safety propoties, and liveness propoties of zab consensus algorithm.

Safety properties

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

$Consistency \triangleq$

$$\begin{aligned}
& \exists i, j \in Server : \\
& \quad \wedge state[i] = Leader \\
& \quad \wedge state[j] = Leader \\
& \quad \wedge currentEpoch[i] = currentEpoch[j] \\
& \Rightarrow i = j
\end{aligned}$$

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

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

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

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

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.

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.

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.

Integrity  $\triangleq \forall l, f \in \text{Server}, \text{msg} \in \text{msgs}:$   
 $\wedge \text{state}[l] = \text{Leader} \wedge \text{state}[f] = \text{Follower}$   
 $\wedge \text{msg.type} = \text{COMMIT} \wedge \text{msg} \in \text{history}[f]$   
 $\Rightarrow \text{msg} \in \text{history}[l]$

LivenessProperty1  $\triangleq \forall i, j \in \text{Server}, \text{msg} \in \text{msgs}: (\text{state}[i] = \text{Leader}) \wedge (\text{msg.type} = \text{COMMIT}) \leadsto (\text{msg} \in \text{history}[j]) \wedge (\text{state}[j] = \text{Follower})$

\ \* Modification History  
 \ \* Last modified *Thu Apr 15 21:25:04 CST 2021* by Dell  
 \ \* Created Sat *Dec 05 13:32:08 CST 2020* by Dell