

This is a specification of the paxos algorithm implemented in Ceph. The specification is based on the following source file: <https://github.com/ceph/ceph/blob/master/src/mon/Paxos.cc>

The main mechanism abstracted that may differ from the version implemented in Ceph are:

- The election logic. The leader is chosen randomly, and, for now, only one leader is chosen per epoch. When a new epoch begins, the messages from the previous epoch are discarded.
- Monitor quorum. The quorum is defined in the election phase, using all monitors that are up. Different epochs can have different quorums.
- The communication layer. The variable `messages` represents connections between monitors (e.g. `messages[mon1][mon2]` holds the messages sent from `mon1` to `mon2`). Within a connection the messages are sent and received in order.
- The transactions. Transactions are simplified to represent only a change of a value in the variable `monitor_store`.
- Failure model. A monitor can crash if the remaining number of monitors is sufficient to form a quorum. When a monitor crashes, new elections are triggered and the monitor is marked to not be part of a quorum until he recovers.
- Timeouts. A timeout can occur at any point in the algorithm and it will trigger new elections.

For a more detailed overview of the specification: <https://github.com/afonsof/ceph-consensus-spec>

EXTENDS *Integers, FiniteSets, Sequences, TLC, SequencesExt, FiniteSetsExt*

Constants

Monitors.

CONSTANTS *Monitors*

$MonitorsSeq \triangleq TLCEval(SetToSeq(Monitors))$

$MonitorsLen \triangleq TLCEval(Len(MonitorsSeq))$

Rank predicate, used to compute proposal numbers.

$rank(mon) \triangleq \text{CHOOSE } i \in 1 \dots MonitorsLen : MonitorsSeq[i] = mon$

Set of possible values.

CONSTANTS *Value_set*

Reserved value.

CONSTANTS *Nil*

Paxos states:

CONSTANTS *STATE_RECOVERING, STATE_ACTIVE,*
STATE_UPDATING, STATE_UPDATING_PREVIOUS,
STATE_WRITING, STATE_WRITING_PREVIOUS,
STATE_REFRESH, STATE_SHUTDOWN

$state_names \triangleq \{STATE_RECOVERING, STATE_ACTIVE,$

*STATE_UPDATING, STATE_UPDATING_PREVIOUS,
STATE_WRITING, STATE_WRITING_PREVIOUS,
STATE_REFRESH, STATE_SHUTDOWN}*

Paxos auxiliary phase states:

They are used to force some sequence of steps.

CONSTANTS *PHASE_ELECTION,*
PHASE_SEND_COLLECT, PHASE_COLLECT,
PHASE_LEASE, PHASE_LEASE_DONE,
PHASE_BEGIN,
PHASE_COMMIT

phase_names \triangleq {*PHASE_ELECTION,*
PHASE_SEND_COLLECT, PHASE_COLLECT,
PHASE_LEASE, PHASE_LEASE_DONE,
PHASE_BEGIN,
PHASE_COMMIT}

Paxos message types:

CONSTANTS *OP_COLLECT, OP_LAST,*
OP_BEGIN, OP_ACCEPT, OP_COMMIT,
OP_LEASE, OP_LEASE_ACK

messages_types \triangleq {*OP_COLLECT, OP_LAST,*
OP_BEGIN, OP_ACCEPT, OP_COMMIT,
OP_LEASE, OP_LEASE_ACK}

Global variables

Integer representing the current epoch. If is odd trigger an election.

Type: Integer

VARIABLE *epoch*

Store messages waiting to be handled.

Type: [*Monitors* \mapsto [*Monitors* \mapsto *message*]]

VARIABLE *messages*

Stores history of messages. Can be useful to find specific states.

Type: {*messages*}

VARIABLE *message_history*

Stores if a monitor is up or down. All available monitors, in a given epoch, are part of the quorum.

Type: [*Monitors* \mapsto *Bool*]

VARIABLE *quorum*

Size of the current quorum.

Type: *Int*

VARIABLE *quorum_sz*

State variables

A function that stores the current leader. $isLeader[mon]$ is True iff mon is a leader, else False.

Type: $[Monitors \mapsto Bool]$

VARIABLE $isLeader$

A function that stores the state of each monitor.

Type: $[Monitors \mapsto state_names]$

VARIABLE $state$

A function that stores the phase of each monitor.

Type: $[Monitors \mapsto phase_names]$

VARIABLE $phase$

Restart variables

A function that stores, for each monitor, a value version when the commit phase starts.

This value version can be retrieved after a monitor crashes and restarts.

Type: $[Monitors \mapsto \text{value version}]$

VARIABLE $uncommitted_v$

A function that stores, for each monitor, a value when the commit phase starts.

This value can be retrieved after a monitor crashes and restarts.

Type: $[Monitors \mapsto Value_set]$

VARIABLE $uncommitted_value$

Data variables

A function that stores, for each monitor, the store where the transactions are applied.

In this model, a transaction represents changing the value in the store.

Type: $[Monitors \mapsto Value_set]$

VARIABLE $monitor_store$

A function that stores the transaction log of each monitor.

Type: $[Monitors \mapsto [value\ version \mapsto Value_set]]$

VARIABLE $values$

A function that stores the last proposal number accepted by each monitor.

Type: $[Monitors \mapsto \text{proposal number}]$

VARIABLE $accepted_pn$

A function that stores the first value version committed by each monitor.

Type: $[Monitors \mapsto \text{value version}]$

VARIABLE $first_committed$

A function that stores the last value version committed by each monitor.

Type: $[Monitors \mapsto \text{value version}]$

VARIABLE $last_committed$

Collect phase variables

A function that stores the number of peers that accepted a collect request.

Type: $[Monitors \mapsto \text{number of peers that accepted}]$

VARIABLE *num_last*

Used by leader when receiving responses in collect phase.

Type: $[Monitors \mapsto [Monitors \mapsto \text{value version}]]$

VARIABLE *peer_first_committed*

Used by leader when receiving responses in collect phase.

Type: $[Monitors \mapsto [Monitors \mapsto \text{value version}]]$

VARIABLE *peer_last_committed*

Lease phase variables

A function that stores, for each monitor, which of the peers have acked the lease request.

Type: $[Monitors \mapsto [Monitors \mapsto Bool]]$

VARIABLE *acked_lease*

Commit phase variables

A function that stores, for each monitor, the value proposed by a client.

Type: $[Monitors \mapsto Value_set \cup \{Nil\}]$

VARIABLE *pending_proposal*

A function that stores, for each monitor, the value to be committed in the begin phase.

Type: $[Monitors \mapsto Value_set \cup \{Nil\}]$

VARIABLE *new_value*

A function that stores, for each monitor, which of the peers have acked the begin request.

Type: $[Monitors \mapsto [Monitors \mapsto Bool]]$

VARIABLE *accepted*

Debug variables

Variables to help debug a behavior.

step is the diameter of a behavior/path.

step_x the current predicate being called.

VARIABLE *step*, *step_x*

Variables to limit the number of monitors crashes that can occur over a behavior.

This variable is used to limit the search space.

VARIABLE *number_crashes*

Variables initialization

global_vars $\triangleq \langle \text{epoch}, \text{messages}, \text{message_history}, \text{quorum}, \text{quorum_sz} \rangle$

$$\begin{aligned}
state_vars &\triangleq \langle isLeader, state, phase \rangle \\
restart_vars &\triangleq \langle uncommitted_v, uncommitted_value \rangle \\
data_vars &\triangleq \langle monitor_store, values, accepted_pn, first_committed, last_committed \rangle \\
collect_vars &\triangleq \langle num_last, peer_first_committed, peer_last_committed \rangle \\
lease_vars &\triangleq \langle acked_lease \rangle \\
commit_vars &\triangleq \langle pending_proposal, new_value, accepted \rangle \\
\\
vars &\triangleq \langle global_vars, state_vars, restart_vars, data_vars, collect_vars, \\
&\quad lease_vars, commit_vars \rangle \\
\\
Init_global_vars &\triangleq \\
&\quad \wedge epoch = 1 \\
&\quad \wedge messages = [mon1 \in Monitors \mapsto [mon2 \in Monitors \mapsto \langle \rangle]] \\
&\quad \wedge message_history = \{ \} \\
&\quad \wedge quorum = [mon \in Monitors \mapsto \text{TRUE}] \\
&\quad \wedge quorum_sz = MonitorsLen \\
\\
Init_state_vars &\triangleq \\
&\quad \wedge isLeader = [mon \in Monitors \mapsto \text{FALSE}] \\
&\quad \wedge state = [mon \in Monitors \mapsto Nil] \\
&\quad \wedge phase = [mon \in Monitors \mapsto Nil] \\
\\
Init_restart_vars &\triangleq \\
&\quad \wedge uncommitted_v = [mon \in Monitors \mapsto 0] \\
&\quad \wedge uncommitted_value = [mon \in Monitors \mapsto Nil] \\
\\
Init_data_vars &\triangleq \\
&\quad \wedge monitor_store = [mon \in Monitors \mapsto Nil] \\
&\quad \wedge values = [mon \in Monitors \mapsto [version \in \{ \} \mapsto Nil]] \\
&\quad \wedge accepted_pn = [mon \in Monitors \mapsto 0] \\
&\quad \wedge first_committed = [mon \in Monitors \mapsto 0] \\
&\quad \wedge last_committed = [mon \in Monitors \mapsto 0] \\
\\
Init_collect_vars &\triangleq \\
&\quad \wedge num_last = [mon \in Monitors \mapsto 0] \\
&\quad \wedge peer_first_committed = [mon1 \in Monitors \mapsto [mon2 \in Monitors \mapsto -1]] \\
&\quad \wedge peer_last_committed = [mon1 \in Monitors \mapsto [mon2 \in Monitors \mapsto -1]] \\
\\
Init_lease_vars &\triangleq \\
&\quad \wedge acked_lease = [mon1 \in Monitors \mapsto [mon2 \in Monitors \mapsto \text{FALSE}]] \\
\\
Init_commit_vars &\triangleq \\
&\quad \wedge pending_proposal = [mon \in Monitors \mapsto Nil] \\
&\quad \wedge new_value = [mon \in Monitors \mapsto Nil] \\
&\quad \wedge accepted = [mon1 \in Monitors \mapsto [mon2 \in Monitors \mapsto \text{FALSE}]] \\
\\
Init &\triangleq \\
&\quad \wedge Init_global_vars
\end{aligned}$$

$\wedge \text{Init_state_vars}$
 $\wedge \text{Init_restart_vars}$
 $\wedge \text{Init_data_vars}$
 $\wedge \text{Init_collect_vars}$
 $\wedge \text{Init_lease_vars}$
 $\wedge \text{Init_commit_vars}$
 $\wedge \text{step} = 0 \wedge \text{step_x} = \text{"init"} \wedge \text{number_crashes} = 0$

Message manipulation

Note: Variable *message_history* has impact in performace, update only when debugging.

Add message *m* to the network *msgs*.

$\text{WithMessage}(m, \text{msgs}) \triangleq$
 $[\text{msgs} \text{ EXCEPT } ![m.\text{from}] =$
 $[\text{msgs}[m.\text{from}] \text{ EXCEPT } ![m.\text{dest}] = \text{Append}(\text{msgs}[m.\text{from}}][m.\text{dest}], m)]$

Remove message *m* from the network *msgs*.

$\text{WithoutMessage}(m, \text{msgs}) \triangleq$
 $[\text{msgs} \text{ EXCEPT } ![m.\text{from}] =$
 $[\text{msgs}[m.\text{from}] \text{ EXCEPT } ![m.\text{dest}] = \text{Remove}(\text{msgs}[m.\text{from}}][m.\text{dest}], m)]$

Adds the message *m* to the network.

Variables changed: *messages*, *message_history*.

$\text{Send}(m) \triangleq$
 $\wedge \text{messages}' = \text{WithMessage}(m, \text{messages})$
 $\wedge \text{message_history}' = \text{message_history} \cup \{m\}$
 $\wedge \text{UNCHANGED } \text{message_history}$

Adds a set of messages to the network.

Variables changed: *messages*, *message_history*.

$\text{Send_set}(\text{from}, m_set) \triangleq$
 $\wedge \text{messages}' = [\text{messages} \text{ EXCEPT } ![from] =$
 $[\text{mon} \in \text{Monitors} \mapsto$
 $\text{messages}[from][mon] \circ \text{SetToSeq}(\{m \in m_set : m.\text{dest} = mon\})]$
 $\wedge \text{message_history}' = \text{message_history} \cup m_set$
 $\wedge \text{UNCHANGED } \text{message_history}$

Removes the request from network and adds the response.

Variables changed: *messages*, *message_history*.

$\text{Reply}(\text{response}, \text{request}) \triangleq$
 $\wedge \text{messages}' = \text{WithoutMessage}(\text{request}, \text{WithMessage}(\text{response}, \text{messages}))$
 $\wedge \text{message_history}' = \text{message_history} \cup \{\text{response}\}$
 $\wedge \text{UNCHANGED } \text{message_history}$

Removes the request from network and adds a set of messages.

Variables changed: *messages*, *message_history*.
 $Reply_set(from, response_set, request) \triangleq$
 $\wedge LET\ msgs \triangleq WithoutMessage(request, messages)$
 $IN\ messages' = [msgs\ EXCEPT\ ![from] =$
 $\quad [mon \in Monitors \mapsto$
 $\quad\quad msgs[from][mon] \circ SetToSeq(\{m \in response_set : m.dest = mon\})]]$
 $\wedge message_history' = message_history \cup response_set$
 $\wedge UNCHANGED\ message_history$

Removes message *m* from the network.

Variables changed: *messages*, *message_history*.
 $Discard(m) \triangleq$
 $\wedge\ messages' = WithoutMessage(m, messages)$
 $\wedge\ UNCHANGED\ message_history$

Helper predicates

Computes a new unique proposal number for a given monitor.

Example: $oldpn = 305$, $rank(mon) = 5$, $newpn = 405$.
 $get_new_proposal_number(mon, oldpn) \triangleq$
 $((oldpn \div 100) + 1) * 100 + rank(mon)$

Clear the variable *peer_first_committed*.

Variables changed: *peer_first_committed*.
 $clear_peer_first_committed(mon) \triangleq$
 $peer_first_committed' = [peer_first_committed\ EXCEPT\ ![mon] =$
 $\quad [m \in Monitors \mapsto -1]]$

Clear the variable *peer_last_committed*.

Variables changed: *peer_last_committed*.
 $clear_peer_last_committed(mon) \triangleq$
 $peer_last_committed' = [peer_last_committed\ EXCEPT\ ![mon] =$
 $\quad [m \in Monitors \mapsto -1]]$

Store peer values and update *first_committed*, *last_committed* and *monitor_store* accordingly.

Variables changed: *values*, *first_committed*, *last_committed*, *monitor_store*.

$store_state(mon, msg) \triangleq$
 $\quad Choose\ peer\ values\ from\ mon\ last\ committed\ + 1\ to\ peer\ last\ committed.$
 $\wedge LET\ logs \triangleq (DOMAIN\ msg.values) \cap (last_committed[mon] + 1 .. msg.last_committed)$
 $IN\ \wedge\ values' = [values\ EXCEPT\ ![mon] =$
 $\quad [i \in DOMAIN\ values[mon] \cup logs \mapsto$
 $\quad\quad IF\ i \in logs$
 $\quad\quad\quad THEN\ msg.values[i]$
 $\quad\quad\quad ELSE\ values[mon][i]]]$
 $\quad Update\ last\ committed\ and\ first\ committed.$
 $\wedge last_committed' = [last_committed\ EXCEPT\ ![mon] = Max(logs \cup \{last_committed[mon]\})]$

\wedge IF $logs \neq \{\}$ \wedge $first_committed[mon] = 0$
 THEN $first_committed' =$
 $[first_committed \text{ EXCEPT } ![mon] = Min(logs)]$
 ELSE $first_committed' =$
 $[first_committed \text{ EXCEPT } ![mon] = Min(logs \cup \{first_committed[mon]\})]$
 Update monitor store.
 \wedge IF $last_committed'[mon] = 0$
 THEN UNCHANGED $monitor_store$
 ELSE $monitor_store' = [monitor_store \text{ EXCEPT } ![mon] = values'[mon][last_committed'[mon]]]$
 Check if uncommitted value version is still valid, else reset it.
 Variables changed: $uncommitted_v, uncommitted_value$.
 $check_and_correct_uncommitted(mon) \triangleq$
 IF $uncommitted_v[mon] \leq last_committed'[mon]$
 THEN \wedge $uncommitted_v' = [uncommitted_v \text{ EXCEPT } ![mon] = 0]$
 \wedge $uncommitted_value' = [uncommitted_value \text{ EXCEPT } ![mon] = Nil]$
 ELSE UNCHANGED $\langle uncommitted_v, uncommitted_value \rangle$
 Trigger new election by incrementing epoch.
 Variables changed: epoch.
 $bootstrap \triangleq$
 \wedge $epoch' = epoch + 1$

Lease phase predicates

Changes mon state to $STATE_ACTIVE$.
 Variables changed: state.
 $finish_round(mon) \triangleq$
 \wedge $isLeader[mon] = \text{TRUE}$
 \wedge $state' = [state \text{ EXCEPT } ![mon] = STATE_ACTIVE]$
 Resets the variable $acked_lease$ and send lease messages to peers.
 Variables changed: $acked_lease, messages, message_history, phase$.
 $extend_lease(mon) \triangleq$
 \wedge $isLeader[mon] = \text{TRUE}$
 \wedge $acked_lease' = [acked_lease \text{ EXCEPT } ![mon] =$
 $[m \in Monitors \mapsto \text{IF } m = mon \text{ THEN TRUE ELSE FALSE}]]$
 \wedge $Send_set(mon,$
 $\{[type \mapsto OP_LEASE,$
 $from \mapsto mon,$
 $dest \mapsto dest,$
 $last_committed \mapsto last_committed[mon]] : dest \in \{m \in Monitors \setminus \{mon\} : quorum[m]\}$
 $\})$
 \wedge $phase' = [phase \text{ EXCEPT } ![mon] = PHASE_LEASE]$

Handle a lease message. The peon changes his state and replies with a lease ack message.

The reply is commented because the lease ack is only used to check if all peers are up.

In the model this is done by “randomly” triggering the predicate *Timeout*. In this way, the search space is reduced.

Variables changed: *messages*, *message_history*, *state*.

$$\begin{aligned}
& \text{handle_lease}(mon, msg) \triangleq \\
& \wedge \text{discard if not peon or peon is behind} \\
& \text{IF } \vee isLeader[mon] = \text{TRUE} \\
& \quad \vee last_committed[mon] \neq msg.last_committed \\
& \text{THEN } \wedge Discard(msg) \\
& \quad \wedge \text{UNCHANGED } state \\
& \text{ELSE } \wedge state' = [state \text{ EXCEPT } ![mon] = STATE_ACTIVE] \\
& \quad \wedge \text{Reply}([type \mapsto OP_LEASE_ACK, \\
& \quad \quad from \mapsto mon, \\
& \quad \quad dest \mapsto msg.from, \\
& \quad \quad first_committed \mapsto first_committed[mon], \\
& \quad \quad last_committed \mapsto last_committed[mon]], msg) \\
& \quad \wedge Discard(msg) \\
& \wedge \text{UNCHANGED } \langle epoch, quorum, quorum_sz, isLeader, phase \rangle \\
& \wedge \text{UNCHANGED } \langle restart_vars, data_vars, collect_vars, lease_vars, commit_vars \rangle
\end{aligned}$$

Handle a lease ack message. The leader updates the *acked_lease* variable.

Because the *lease_ack* messages are not sent, this predicate is never called.

The reasoning for this is given in *handle_lease* comment.

Variables changed: *acked_lease*, *messages*, *message_history*.

$$\begin{aligned}
& \text{handle_lease_ack}(mon, msg) \triangleq \\
& \wedge phase[mon] = PHASE_LEASE \\
& \wedge acked_lease' = [acked_lease \text{ EXCEPT } ![mon] = \\
& \quad [acked_lease[mon] \text{ EXCEPT } ![msg.from] = \text{TRUE}]] \\
& \wedge Discard(msg) \\
& \wedge \text{UNCHANGED } \langle epoch, quorum, quorum_sz \rangle \\
& \wedge \text{UNCHANGED } \langle state_vars, restart_vars, data_vars, collect_vars, commit_vars \rangle
\end{aligned}$$

Predicate that is called when all peers ack the lease. The phase is changed to prevent loops.

Because the *lease_ack* messages are not sent, this predicate is never called.

The reasoning for this is given in *handle_lease* comment.

Variables changed: *phase*.

$$\begin{aligned}
& \text{post_lease_ack}(mon) \triangleq \\
& \wedge phase[mon] = PHASE_LEASE \\
& \wedge phase' = [phase \text{ EXCEPT } ![mon] = PHASE_LEASE_DONE] \\
& \wedge \forall m \in Monitors : quorum[m] \Rightarrow acked_lease[mon][m] = \text{TRUE} \\
& \wedge \text{UNCHANGED } \langle isLeader, state \rangle \\
& \wedge \text{UNCHANGED } \langle global_vars, restart_vars, data_vars, collect_vars, \\
& \quad lease_vars, commit_vars \rangle
\end{aligned}$$

Commit phase predicates

Start a commit phase by the leader. The variable *new_value* is assigned. Send begin messages to the peers.

The value of *uncommitted_v* and *uncommitted_value* are assigned in order for the leader to be able to recover from a crash.

Variables changed: *accepted*, *new_value*, *phase*, *messages*, *message_history*, *values*, *uncommitted_v*, *uncommitted_value*.

```

begin(mon, v)  $\triangleq$ 
   $\wedge$  isLeader[mon] = TRUE
   $\wedge$   $\vee$  state'[mon] = STATE_UPDATING
     $\vee$  state'[mon] = STATE_UPDATING_PREVIOUS
   $\wedge$  quorum_sz = 1  $\vee$  num_last[mon] > MonitorsLen  $\div$  2
   $\wedge$  new_value[mon] = Nil
   $\wedge$  accepted' = [accepted EXCEPT ![mon] =
    [m  $\in$  Monitors  $\mapsto$  IF m = mon THEN TRUE ELSE FALSE]]
   $\wedge$  new_value' = [new_value EXCEPT ![mon] = v]
   $\wedge$  phase' = [phase EXCEPT ![mon] = PHASE_BEGIN]
   $\wedge$  values' = [values EXCEPT ![mon] =
    (values[mon] @@ ((last_committed[mon] + 1) :> new_value'[mon]))]
   $\wedge$  Send_set(mon,
    {[type  $\mapsto$  OP_BEGIN,
      from  $\mapsto$  mon,
      dest  $\mapsto$  dest,
      last_committed  $\mapsto$  last_committed[mon],
      values  $\mapsto$  values'[mon],
      pn  $\mapsto$  accepted_pn[mon]] : dest  $\in$  {m  $\in$  Monitors \ {mon} : quorum[m]}
    })
   $\wedge$  uncommitted_v' = [uncommitted_v EXCEPT ![mon] = last_committed[mon] + 1]
   $\wedge$  uncommitted_value' = [uncommitted_value EXCEPT ![mon] = v]

```

Handle a begin message. The monitor will accept if the proposal number in the message is greater or equal than the one he accepted.

Similar to what happens in begin, *uncommitted_v* and *uncommitted_value* are assigned in order for the monitor to recover in case of a crash.

Variables changed: *messages*, *message_history*, *state*, *values*, *uncommitted_v*, *uncommitted_value*.

```

handle_begin(mon, msg)  $\triangleq$ 
   $\wedge$  isLeader[mon] = FALSE
   $\wedge$  IF msg.pn < accepted_pn[mon]
    THEN
       $\wedge$  Discard(msg)
       $\wedge$  UNCHANGED <state, values, restart_vars>
    ELSE
       $\wedge$  msg.pn = accepted_pn[mon]
       $\wedge$  msg.last_committed = last_committed[mon]

      assign values[mon][last_committed[mon] + 1]
       $\wedge$  values' = [values EXCEPT ![mon] =
        (values[mon] @@ ((last_committed[mon] + 1) :> msg.values[last_committed[mon] + 1]))]

```

$$\begin{aligned}
& \wedge state' = [state \text{ EXCEPT } ![mon] = STATE_UPDATING] \\
& \wedge uncommitted_v' = [uncommitted_v \text{ EXCEPT } ![mon] = last_committed[mon] + 1] \\
& \wedge uncommitted_value' = [uncommitted_value \text{ EXCEPT } ![mon] = \\
& \quad values'[mon][last_committed[mon] + 1]] \\
& \wedge Reply([type \quad \mapsto OP_ACCEPT, \\
& \quad \quad from \quad \mapsto mon, \\
& \quad \quad dest \quad \mapsto msg.from, \\
& \quad \quad last_committed \mapsto last_committed[mon], \\
& \quad \quad pn \quad \mapsto accepted_pn[mon]], msg) \\
& \wedge UNCHANGED \langle epoch, quorum, quorum_sz, isLeader, phase, monitor_store, \\
& \quad \quad \quad accepted_pn, first_committed, last_committed \rangle \\
& \wedge UNCHANGED \langle collect_vars, lease_vars, commit_vars \rangle
\end{aligned}$$

Handle an accept message. If the leader receives a positive response from the peer, it will add it to the variable accepted.

Variables changed: messages, *message_history*, accepted

$$\begin{aligned}
& handle_accept(mon, msg) \triangleq \\
& \quad \wedge isLeader[mon] = TRUE \\
& \quad \wedge \vee state[mon] = STATE_UPDATING_PREVIOUS \\
& \quad \quad \vee state[mon] = STATE_UPDATING \\
& \quad \wedge phase[mon] = PHASE_BEGIN \\
& \quad \wedge new_value[mon] \neq Nil \\
& \quad \wedge IF \quad \vee msg.pn \neq accepted_pn[mon] \\
& \quad \quad \vee \wedge last_committed[mon] > 0 \\
& \quad \quad \quad \wedge msg.last_committed < last_committed[mon] - 1 \\
& \quad \quad THEN UNCHANGED accepted \\
& \quad \quad ELSE accepted' = [accepted EXCEPT ![mon] = \\
& \quad \quad \quad [accepted[mon] EXCEPT ![msg.from] = TRUE]] \\
& \quad \wedge Discard(msg) \\
& \quad \wedge UNCHANGED \langle epoch, quorum, quorum_sz, pending_proposal, new_value \rangle \\
& \quad \wedge UNCHANGED \langle restart_vars, state_vars, data_vars, collect_vars, lease_vars \rangle
\end{aligned}$$

Predicate that is enabled and called when all peers in the quorum accept begin request from leader.

The leader commits the transaction in *new_value* and sends commit messages to his peers.

Variables changed: *first_committed*, *last_committed*, *monitor_store*, *new_value*, messages, *message_history*, state, phase

$$\begin{aligned}
& post_accept(mon) \triangleq \\
& \quad \wedge phase[mon] = PHASE_BEGIN \\
& \quad \wedge \forall m \in Monitors : quorum[m] \Rightarrow accepted[mon][m] = TRUE \\
& \quad \wedge new_value[mon] \neq Nil \\
& \quad \wedge \vee state[mon] = STATE_UPDATING_PREVIOUS \\
& \quad \quad \vee state[mon] = STATE_UPDATING \\
& \quad \wedge last_committed' = [last_committed EXCEPT ![mon] = last_committed[mon] + 1] \\
& \quad \wedge IF first_committed[mon] = 0 \\
& \quad \quad THEN first_committed' = [first_committed EXCEPT ![mon] = first_committed[mon] + 1] \\
& \quad \quad ELSE UNCHANGED first_committed
\end{aligned}$$

$\wedge \text{monitor_store}' = [\text{monitor_store} \text{ EXCEPT } ![mon] = \text{values}[mon][\text{last_committed}[mon] + 1]]$
 $\wedge \text{new_value}' = [\text{new_value} \text{ EXCEPT } ![mon] = \text{Nil}]$
 $\wedge \text{Send_set}(mon,$
 $\quad \{[type \quad \mapsto OP_COMMIT,$
 $\quad \quad from \quad \mapsto mon,$
 $\quad \quad dest \quad \mapsto dest,$
 $\quad \quad last_committed \mapsto last_committed'[mon],$
 $\quad \quad pn \quad \mapsto accepted_pn[mon],$
 $\quad \quad values \quad \mapsto \text{values}[mon]] : dest \in \{m \in \text{Monitors} \setminus \{mon\} : \text{quorum}[m]\}$
 $\quad \})$
 $\wedge \text{state}' = [\text{state} \text{ EXCEPT } ![mon] = \text{STATE_REFRESH}]$
 $\wedge \text{phase}' = [\text{phase} \text{ EXCEPT } ![mon] = \text{PHASE_COMMIT}]$
 $\wedge \text{UNCHANGED } \langle isLeader, values, accepted_pn, pending_proposal, accepted \rangle$
 $\wedge \text{UNCHANGED } \langle epoch, quorum, quorum_sz, restart_vars, collect_vars, lease_vars \rangle$

Predicate that is called after *post_accept*. The leader finishes the commit phase by updating his state to *STATE_ACTIVE* and by extending the lease to his peers.

Variables changed: *state, phase, acked_lease, messages, message_history*.

$\text{finish_commit}(mon) \triangleq$
 $\quad \wedge \text{state}[mon] = \text{STATE_REFRESH}$
 $\quad \wedge \text{phase}[mon] = \text{PHASE_COMMIT}$
 $\quad \wedge \text{finish_round}(mon)$
 $\quad \wedge \text{extend_lease}(mon)$
 $\quad \wedge \text{UNCHANGED } \langle epoch, quorum, quorum_sz, isLeader \rangle$
 $\quad \wedge \text{UNCHANGED } \langle restart_vars, data_vars, collect_vars, commit_vars \rangle$

Handle a commit message. The monitor stores the values sent by the leader commit message.

Variables changed: *messages, message_history, values, first_committed, last_committed, monitor_store, uncommitted_v, uncommitted_value*.

$\text{handle_commit}(mon, msg) \triangleq$
 $\quad \wedge isLeader[mon] = \text{FALSE}$
 $\quad \wedge \text{store_state}(mon, msg)$
 $\quad \wedge \text{check_and_correct_uncommitted}(mon)$
 $\quad \wedge \text{Discard}(msg)$
 $\quad \wedge \text{UNCHANGED } \langle epoch, quorum, quorum_sz, accepted_pn \rangle$
 $\quad \wedge \text{UNCHANGED } \langle state_vars, collect_vars, lease_vars, commit_vars \rangle$

Client Request

Request a transaction *v* to the monitor. The transaction is saved on pending proposal to be committed in the next available commit phase.

Variables changed: *pending_proposal*.

$\text{client_request}(mon, v) \triangleq$
 $\quad \wedge isLeader[mon] = \text{TRUE}$
 $\quad \wedge \text{state}[mon] = \text{STATE_ACTIVE}$

$\wedge \text{pending_proposal}[mon] = Nil$
 $\wedge \text{pending_proposal}' = [\text{pending_proposal} \text{ EXCEPT } ![mon] = v]$
 $\wedge \text{UNCHANGED } \langle \text{new_value}, \text{accepted} \rangle$
 $\wedge \text{UNCHANGED } \langle \text{global_vars}, \text{state_vars}, \text{restart_vars}, \text{data_vars}, \text{collect_vars}, \text{lease_vars} \rangle$

Start a commit phase with the value on pending proposal.

Variables changed: *state*, *pending_proposal*, *accepted*, *new_value*, *phase*, *messages*, *message_history*, *values*, *uncommitted_v*, *uncommitted_value*.

$\text{propose_pending}(mon) \triangleq$
 $\wedge \text{phase}[mon] = \text{PHASE_LEASE} \vee \text{phase}[mon] = \text{PHASE_ELECTION}$
 $\wedge \text{state}[mon] = \text{STATE_ACTIVE}$
 $\wedge \text{pending_proposal}[mon] \neq Nil$
 $\wedge \text{pending_proposal}' = [\text{pending_proposal} \text{ EXCEPT } ![mon] = Nil]$
 $\wedge \text{state}' = [\text{state} \text{ EXCEPT } ![mon] = \text{STATE_UPDATING}]$
 $\wedge \text{begin}(mon, \text{pending_proposal}[mon])$
 $\wedge \text{UNCHANGED } \langle \text{isLeader}, \text{monitor_store}, \text{accepted_pn}, \text{first_committed}, \text{last_committed} \rangle$
 $\wedge \text{UNCHANGED } \langle \text{epoch}, \text{quorum}, \text{quorum_sz}, \text{collect_vars}, \text{lease_vars} \rangle$

Collect phase predicates

Start collect phase. This first part of the collect phase is divided in two parts (*collect* and *send_collect*) in order to simplify variable changes (when collect is triggered from *handle_last*).

Variables changed: *accepted_pn*, *phase*.

$\text{collect}(mon, \text{oldpn}) \triangleq$
 $\wedge \text{state}[mon] = \text{STATE_RECOVERING}$
 $\wedge \text{isLeader}[mon] = \text{TRUE}$
 $\wedge \text{LET } \text{new_pn} \triangleq \text{get_new_proposal_number}(mon, \text{Max}(\{\text{oldpn}, \text{accepted_pn}[mon]\}))$
 $\quad \text{IN } \wedge \text{accepted_pn}' = [\text{accepted_pn} \text{ EXCEPT } ![mon] = \text{new_pn}]$
 $\wedge \text{phase}' = [\text{phase} \text{ EXCEPT } ![mon] = \text{PHASE_SEND_COLLECT}]$

Continue the start of the collect phase. Initialize the number of peers that accepted the proposal (*num_last*) and the variables with peers version numbers. Check if there is an uncommitted value.

Send collect messages to the peers.

Variables changed: *peer_first_committed*, *peer_last_committed*, *uncommitted_v*, *uncommitted_value*, *num_last*, *messages*, *message_history*, *phase*.

$\text{send_collect}(mon) \triangleq$
 $\wedge \text{state}[mon] = \text{STATE_RECOVERING}$
 $\wedge \text{isLeader}[mon] = \text{TRUE}$
 $\wedge \text{phase}[mon] = \text{PHASE_SEND_COLLECT}$
 $\wedge \text{clear_peer_first_committed}(mon)$
 $\wedge \text{clear_peer_last_committed}(mon)$
 $\wedge \text{IF } \text{last_committed}[mon] + 1 \in \text{DOMAIN } \text{values}[mon]$
 $\quad \text{THEN } \wedge \text{uncommitted_v}' =$
 $\quad \quad [\text{uncommitted_v} \text{ EXCEPT } ![mon] = \text{last_committed}[mon] + 1]$

$$\begin{aligned}
& \wedge \text{uncommitted_value}' = \\
& \quad [\text{uncommitted_value} \text{ EXCEPT } ![mon] = \text{values}[mon][\text{last_committed}[mon] + 1]] \\
& \text{ELSE UNCHANGED } \langle \text{restart_vars} \rangle \\
& \wedge \text{num_last}' = [\text{num_last} \text{ EXCEPT } ![mon] = 1] \\
& \wedge \text{Send_set}(mon, \\
& \quad \{[type \mapsto OP_COLLECT, \\
& \quad \quad from \mapsto mon, \\
& \quad \quad dest \mapsto dest, \\
& \quad \quad first_committed \mapsto first_committed[mon], \\
& \quad \quad last_committed \mapsto last_committed[mon], \\
& \quad \quad pn \mapsto \text{accepted_pn}[mon]] : dest \in \{m \in \text{Monitors} \setminus \{mon\} : \text{quorum}[m]\} \\
& \quad \}) \\
& \wedge \text{phase}' = [\text{phase} \text{ EXCEPT } ![mon] = PHASE_COLLECT] \\
& \wedge \text{UNCHANGED } \langle \text{isLeader}, \text{state} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{epoch}, \text{quorum}, \text{quorum_sz}, \text{data_vars}, \text{lease_vars}, \text{commit_vars} \rangle
\end{aligned}$$

Handle a collect message. The peer will accept the proposal number from the leader if it is bigger than the last proposal number he accepted.

Variables changed: messages, message_history, epoch, state, accepted_pn

$$\begin{aligned}
& \text{handle_collect}(mon, msg) \triangleq \\
& \quad \wedge \text{isLeader}[mon] = \text{FALSE} \\
& \quad \wedge \text{state}' = [\text{state} \text{ EXCEPT } ![mon] = \text{STATE_RECOVERING}] \\
& \quad \wedge \vee \wedge \text{msg.first_committed} > \text{last_committed}[mon] + 1 \\
& \quad \quad \wedge \text{bootstrap} \\
& \quad \quad \wedge \text{Discard}(msg) \\
& \quad \quad \wedge \text{UNCHANGED } \langle \text{accepted_pn} \rangle \\
& \quad \vee \wedge \text{msg.first_committed} \leq \text{last_committed}[mon] + 1 \\
& \quad \quad \wedge \text{IF } \text{msg.pn} > \text{accepted_pn}[mon] \\
& \quad \quad \quad \text{THEN } \text{accepted_pn}' = [\text{accepted_pn} \text{ EXCEPT } ![mon] = \text{msg.pn}] \\
& \quad \quad \quad \text{ELSE UNCHANGED } \text{accepted_pn} \\
& \quad \wedge \text{Reply}([type \mapsto OP_LAST, \\
& \quad \quad from \mapsto mon, \\
& \quad \quad dest \mapsto \text{msg.from}, \\
& \quad \quad first_committed \mapsto first_committed[mon], \\
& \quad \quad last_committed \mapsto last_committed[mon], \\
& \quad \quad values \mapsto \text{values}[mon], \\
& \quad \quad pn \mapsto \text{accepted_pn}'[mon]], msg) \\
& \quad \wedge \text{UNCHANGED } \text{epoch} \\
& \quad \wedge \text{UNCHANGED } \langle \text{isLeader}, \text{phase}, \text{values}, \text{first_committed}, \text{last_committed}, \text{monitor_store} \rangle \\
& \quad \wedge \text{UNCHANGED } \langle \text{quorum}, \text{quorum_sz}, \text{restart_vars}, \text{collect_vars}, \text{lease_vars}, \text{commit_vars} \rangle
\end{aligned}$$

Handle a last message (response from a peer to the leader collect message).

The peers first and last committed version are stored. If the leader is behind, bootstraps. Stores any value that the peer may have committed (store_state). If peer is behind send commit message with leader values.

If peer accepted proposal number increase num last, if he sent a bigger proposal number start a new collect phase.

Variables changed: *messages*, *message_history*, *epoch*, *phase*, *uncommitted_v*, *uncommitted_value*, *monitor_store*, *values*, *accepted_pn*, *first_committed*, *last_committed*, *num_last*, *peer_first_committed*, *peer_last_committed*.

$$\begin{aligned}
& \text{handle_last}(\text{mon}, \text{msg}) \triangleq \\
& \quad \wedge \text{isLeader}[\text{mon}] = \text{TRUE} \\
& \quad \wedge \text{peer_first_committed}' = [\text{peer_first_committed} \text{ EXCEPT } ![\text{mon}] = \\
& \quad \quad [\text{peer_first_committed}[\text{mon}] \text{ EXCEPT } ![\text{msg.from}] = \text{msg.first_committed}]] \\
& \quad \wedge \text{peer_last_committed}' = [\text{peer_last_committed} \text{ EXCEPT } ![\text{mon}] = \\
& \quad \quad [\text{peer_last_committed}[\text{mon}] \text{ EXCEPT } ![\text{msg.from}] = \text{msg.last_committed}]] \\
& \quad \wedge \text{IF } \text{msg.first_committed} > \text{last_committed}[\text{mon}] + 1 \\
& \quad \quad \text{THEN} \\
& \quad \quad \quad \wedge \text{bootstrap} \\
& \quad \quad \quad \wedge \text{Discard}(\text{msg}) \\
& \quad \quad \quad \wedge \text{UNCHANGED } \langle \text{num_last}, \text{accepted_pn}, \text{values}, \text{phase}, \text{monitor_store} \rangle \\
& \quad \quad \quad \wedge \text{UNCHANGED } \langle \text{first_committed}, \text{last_committed}, \text{restart_vars} \rangle \\
& \quad \quad \text{ELSE} \\
& \quad \quad \quad \wedge \text{store_state}(\text{mon}, \text{msg}) \\
& \quad \quad \quad \wedge \text{IF } \exists \text{peer} \in \text{Monitors} : \\
& \quad \quad \quad \quad \wedge \text{peer} \neq \text{mon} \\
& \quad \quad \quad \quad \wedge \text{peer_last_committed}'[\text{mon}][\text{peer}] \neq -1 \\
& \quad \quad \quad \quad \wedge \text{peer_last_committed}'[\text{mon}][\text{peer}] + 1 < \text{first_committed}[\text{mon}] \\
& \quad \quad \quad \quad \wedge \text{first_committed}[\text{mon}] > 1 \\
& \quad \quad \quad \text{THEN} \\
& \quad \quad \quad \quad \wedge \text{bootstrap} \\
& \quad \quad \quad \quad \wedge \text{check_and_correct_uncommitted}(\text{mon}) \\
& \quad \quad \quad \quad \wedge \text{Discard}(\text{msg}) \\
& \quad \quad \quad \quad \wedge \text{UNCHANGED } \langle \text{phase}, \text{accepted_pn}, \text{num_last} \rangle \\
& \quad \quad \quad \text{ELSE} \\
& \quad \quad \quad \text{LET } \text{monitors_behind} \triangleq \{ \text{peer} \in \text{Monitors} : \\
& \quad \quad \quad \quad \wedge \text{peer} \neq \text{mon} \\
& \quad \quad \quad \quad \wedge \text{peer_last_committed}'[\text{mon}][\text{peer}] \neq -1 \\
& \quad \quad \quad \quad \wedge \text{peer_last_committed}'[\text{mon}][\text{peer}] < \text{last_committed}[\text{mon}] \\
& \quad \quad \quad \quad \wedge \text{quorum}[\text{peer}] \} \\
& \quad \quad \text{IN } \text{Reply_set}(\text{mon}, \\
& \quad \quad \quad \{ [\text{type} \mapsto \text{OP_COMMIT}, \\
& \quad \quad \quad \text{from} \mapsto \text{mon}, \\
& \quad \quad \quad \text{dest} \mapsto \text{dest}, \\
& \quad \quad \quad \text{last_committed} \mapsto \text{last_committed}'[\text{mon}], \\
& \quad \quad \quad \text{pn} \mapsto \text{accepted_pn}[\text{mon}], \\
& \quad \quad \quad \text{values} \mapsto \text{values}[\text{mon}]] : \text{dest} \in \text{monitors_behind} \\
& \quad \quad \quad \}, \text{msg}) \\
& \quad \wedge \vee \wedge \text{msg.pn} > \text{accepted_pn}[\text{mon}] \\
& \quad \quad \wedge \text{collect}(\text{mon}, \text{msg.pn}) \\
& \quad \quad \wedge \text{check_and_correct_uncommitted}(\text{mon})
\end{aligned}$$

$\wedge \text{UNCHANGED } num_last$
 $\vee \wedge msg.pn = accepted_pn[mon]$
 $\wedge num_last' = [num_last \text{ EXCEPT } ![mon] = num_last[mon] + 1]$
 $\wedge \text{IF } \wedge msg.last_committed + 1 \in \text{DOMAIN } msg.values$
 $\wedge msg.last_committed \geq last_committed'[mon]$
 $\wedge msg.last_committed + 1 \geq uncommitted_v[mon]$
 $\text{THEN } \wedge uncommitted_v' =$
 $\quad [uncommitted_v \text{ EXCEPT } ![mon] = msg.last_committed + 1]$
 $\wedge uncommitted_value' =$
 $\quad [uncommitted_value \text{ EXCEPT } ![mon] = msg.values[msg.last_committed + 1]]$
 $\text{ELSE } check_and_correct_uncommitted(mon)$
 $\wedge \text{UNCHANGED } \langle phase, accepted_pn \rangle$
 $\vee \wedge msg.pn < accepted_pn[mon]$
 $\wedge check_and_correct_uncommitted(mon)$
 $\wedge \text{UNCHANGED } \langle phase, accepted_pn, num_last \rangle$
 $\wedge \text{UNCHANGED } epoch$
 $\wedge \text{UNCHANGED } \langle epoch \rangle$
 $\wedge \text{UNCHANGED } \langle quorum, quorum_sz, isLeader, state \rangle$
 $\wedge \text{UNCHANGED } \langle lease_vars, commit_vars \rangle$

Predicate that is enabled and called when all peers in quorum accept collect request from leader. If there is an uncommitted value, a commit phase is started with that value, else the leader changes to *ACTIVE_STATE* and extends the lease to his peers.

Variables changed: *peer_first_committed*, *peer_last_committed*, *state*, *accepted*, *new_value*, *phase*, *messages*, *message_history*, *values*, *uncommitted_v*, *uncommitted_value*, *acked_lease*.

$post_last(mon) \triangleq$
 $\wedge isLeader[mon] = \text{TRUE}$
 $\wedge num_last[mon] = quorum_sz$
 $\wedge phase[mon] = PHASE_COLLECT$
 $\wedge clear_peer_first_committed(mon)$
 $\wedge clear_peer_last_committed(mon)$
 $\wedge \text{IF } \wedge uncommitted_v[mon] = last_committed[mon] + 1$
 $\quad \wedge uncommitted_value[mon] \neq Nil$
 $\text{THEN } \wedge state' = [state \text{ EXCEPT } ![mon] = STATE_UPDATING_PREVIOUS]$
 $\quad \wedge begin(mon, uncommitted_value[mon])$
 $\quad \wedge \text{UNCHANGED } \langle acked_lease \rangle$
 $\text{ELSE } \wedge finish_round(mon)$
 $\quad \wedge extend_lease(mon)$
 $\quad \wedge \text{UNCHANGED } \langle accepted, new_value, values, restart_vars \rangle$
 $\wedge \text{UNCHANGED } \langle isLeader, monitor_store, accepted_pn, first_committed, last_committed \rangle$
 $\wedge \text{UNCHANGED } \langle epoch, quorum, quorum_sz, num_last, pending_proposal \rangle$

Leader election

Elect one monitor as a leader and initialize the remaining ones as peons.

Variables changed: *isLeader*, *state*, *phase*, *new_value*, *pending_proposal*, *epoch*.

leader_election \triangleq

- $\wedge \exists mon \in Monitors :$
 - $\wedge quorum[mon]$
 - $\wedge isLeader' = [m \in Monitors \mapsto \text{IF } m = mon \text{ THEN TRUE ELSE FALSE}]$
 - $\wedge state' = [m \in Monitors \mapsto$
 - $\text{IF } quorum_sz = 1 \text{ THEN } STATE_ACTIVE \text{ ELSE } STATE_RECOVERING]$
 - $\wedge phase' = [m \in Monitors \mapsto PHASE_ELECTION]$
 - $\wedge new_value' = [m \in Monitors \mapsto Nil]$
 - $\wedge pending_proposal' = [m \in Monitors \mapsto Nil]$
 - $\wedge epoch' = epoch + 1$
 - $\wedge messages' = [mon1 \in Monitors \mapsto [mon2 \in Monitors \mapsto \langle \rangle]]$
 - $\wedge \text{UNCHANGED } \langle quorum, quorum_sz, accepted, message_history \rangle$
 - $\wedge \text{UNCHANGED } \langle data_vars, restart_vars, collect_vars, lease_vars \rangle$

Start recovery phase if number of monitors in quorum is greater than 1.

Variables changed: *accepted_pn*, *phase*.

election_recover(mon) \triangleq

- $\wedge quorum_sz > 1$
- $\wedge phase[mon] = PHASE_ELECTION$
- $\wedge collect(mon, 0)$
- $\wedge \text{UNCHANGED } \langle isLeader, state, values, first_committed, last_committed, monitor_store \rangle$
- $\wedge \text{UNCHANGED } \langle global_vars, restart_vars, collect_vars, lease_vars, commit_vars \rangle$

Timeouts and restart

crash_mon(mon) \triangleq

- $\wedge quorum_sz > (MonitorsLen \div 2) + 1$
- $\wedge quorum[mon] = \text{TRUE}$
- $\wedge quorum' = [quorum \text{ EXCEPT } ![mon] = \text{FALSE}]$
- $\wedge quorum_sz' = quorum_sz - 1$
- $\wedge bootstrap$
 - $\wedge number_crashes' = number_crashes + 1$
- $\wedge \text{UNCHANGED } \langle messages, message_history \rangle$
- $\wedge \text{UNCHANGED } \langle state_vars, restart_vars, data_vars, collect_vars, lease_vars, commit_vars \rangle$

restore_mon(mon) \triangleq

- $\wedge quorum[mon] = \text{FALSE}$
- $\wedge quorum' = [quorum \text{ EXCEPT } ![mon] = \text{TRUE}]$
- $\wedge quorum_sz' = quorum_sz + 1$
- $\wedge bootstrap$

$\wedge \text{UNCHANGED } \langle \text{messages}, \text{message_history} \rangle$
 $\wedge \text{UNCHANGED } \langle \text{state_vars}, \text{restart_vars}, \text{data_vars}, \text{collect_vars}, \text{lease_vars}, \text{commit_vars} \rangle$

Monitor timeout (simulate the various timeouts that can occur). Triggers new elections.

Variables changed: epoch.

$\text{Timeout}(\text{mon}) \triangleq$

$\wedge \text{bootstrap}$
 $\wedge \text{UNCHANGED } \langle \text{messages}, \text{quorum}, \text{quorum_sz}, \text{message_history}, \text{state_vars}, \text{restart_vars}, \text{data_vars}, \text{collect_vars}, \text{lease_vars}, \text{commit_vars} \rangle$

Dispatchers and next statement

Handle a message.

$\text{Receive}(\text{msg}) \triangleq$

$\wedge \vee \wedge \text{msg.type} = \text{OP_COLLECT}$
 $\wedge \text{handle_collect}(\text{msg.dest}, \text{msg})$
 $\wedge \text{step_x}' = \text{"receive collect"}$

 $\vee \wedge \text{msg.type} = \text{OP_LAST}$
 $\wedge \text{handle_last}(\text{msg.dest}, \text{msg})$
 $\wedge \text{step_x}' = \text{"receive last"}$

 $\vee \wedge \text{msg.type} = \text{OP_LEASE}$
 $\wedge \text{handle_lease}(\text{msg.dest}, \text{msg})$
 $\wedge \text{step_x}' = \text{"receive lease"}$

 $\vee \wedge \text{msg.type} = \text{OP_LEASE_ACK}$
 $\wedge \text{handle_lease_ack}(\text{msg.dest}, \text{msg})$
 $\wedge \text{step_x}' = \text{"receive lease_ack"}$

 $\vee \wedge \text{msg.type} = \text{OP_BEGIN}$
 $\wedge \text{handle_begin}(\text{msg.dest}, \text{msg})$
 $\wedge \text{step_x}' = \text{"receive begin"}$

 $\vee \wedge \text{msg.type} = \text{OP_ACCEPT}$
 $\wedge \text{handle_accept}(\text{msg.dest}, \text{msg})$
 $\wedge \text{step_x}' = \text{"receive accept"}$

 $\vee \wedge \text{msg.type} = \text{OP_COMMIT}$
 $\wedge \text{handle_commit}(\text{msg.dest}, \text{msg})$
 $\wedge \text{step_x}' = \text{"receive commit"}$

Limit some variables to reduce search space.

$\text{reduce_search_space} \triangleq$

$\wedge \text{epoch} \neq 8$
 $\wedge \text{number_crashes} \neq 20$
 $\wedge \forall \text{mon} \in \text{Monitors} : \text{last_committed}[\text{mon}] < 2$

$\Rightarrow \forall mon2 \in Monitors : new_value[mon2] = Nil$
 $\wedge \forall mon \in Monitors : accepted_pn[mon] < 300$

State transitions.

$Next \triangleq$

$\wedge reduce_search_space$
 \wedge IF $epoch \% 2 = 1$ THEN
 $\wedge leader_election$
 $\wedge step_x' = \text{"election"} \wedge step' = step + 1$
 \wedge UNCHANGED $number_crashes$
ELSE
 $\vee \wedge \exists mon \in Monitors : election_recover(mon)$
 $\wedge step_x' = \text{"election_recover"} \wedge step' = step + 1$
 \wedge UNCHANGED $number_crashes$
 $\vee \wedge \exists mon \in Monitors : send_collect(mon)$
 $\wedge step_x' = \text{"pre_send_collect"} \wedge step' = step + 1$
 \wedge UNCHANGED $number_crashes$
 $\vee \wedge \exists mon \in Monitors : post_last(mon)$
 $\wedge step_x' = \text{"post_last"} \wedge step' = step + 1$
 \wedge UNCHANGED $number_crashes$
 $\vee \wedge \exists mon \in Monitors : post_lease_ack(mon)$
 $\wedge step_x' = \text{"post_lease_ack"} \wedge step' = step + 1$
 \wedge UNCHANGED $number_crashes$
 $\vee \wedge \exists mon \in Monitors : post_accept(mon)$
 $\wedge step_x' = \text{"post_accept"} \wedge step' = step + 1$
 \wedge UNCHANGED $number_crashes$
 $\vee \wedge \exists mon \in Monitors : finish_commit(mon)$
 $\wedge step_x' = \text{"finish_commit"} \wedge step' = step + 1$
 \wedge UNCHANGED $number_crashes$
 $\vee \wedge \exists mon \in Monitors : \exists v \in Value_set : client_request(mon, v)$
 $\wedge step_x' = \text{"client_request"} \wedge step' = step + 1$
 \wedge UNCHANGED $number_crashes$
 $\vee \wedge \exists mon \in Monitors : propose_pending(mon)$
 $\wedge step_x' = \text{"propose_pending"} \wedge step' = step + 1$
 \wedge UNCHANGED $number_crashes$
 $\vee \wedge \exists mon1, mon2 \in Monitors :$
 $\wedge mon1 \neq mon2$
 $\wedge Len(messages[mon1][mon2]) > 0$
 $\wedge Receive(messages[mon1][mon2][1])$
 $\wedge step' = step + 1$

$$\begin{aligned}
& \wedge \text{UNCHANGED } number_crashes \\
& \vee \wedge \exists mon \in Monitors : crash_mon(mon) \\
& \quad \wedge step_x' = \text{"crash mon"} \wedge step' = step + 1 \\
& \quad \wedge \text{UNCHANGED } number_crashes \\
& \vee \wedge \exists mon \in Monitors : restore_mon(mon) \\
& \quad \wedge step_x' = \text{"restore mon"} \wedge step' = step + 1 \\
& \quad \wedge \text{UNCHANGED } number_crashes \\
& \vee \wedge \exists mon \in Monitors : Timeout(mon) \\
& \quad \wedge step_x' = \text{"timeout and restart"} \wedge step' = step + 1 \\
& \quad \wedge \text{UNCHANGED } number_crashes
\end{aligned}$$

Safety invariants

If two monitors are in state active then their *monitor_store* must have the same value.

$$\begin{aligned}
same_monitor_store & \triangleq \forall mon1, mon2 \in Monitors : \\
& state[mon1] = STATE_ACTIVE \wedge state[mon2] = STATE_ACTIVE \\
& \Rightarrow monitor_store[mon1] = monitor_store[mon2] \\
Inv & \triangleq \wedge same_monitor_store
\end{aligned}$$

Test/Debug invariants

Invariant used to search for a state where 'x' happens.

$$Inv_find_state(x) \triangleq \neg x$$

Invariant used to search for a behavior of diameter equal to 'size'.

$$Inv_diam(size) \triangleq step \neq size - 1$$

Invariants to test in model check

$$\begin{aligned}
DEBUG_Inv & \triangleq \wedge \text{TRUE} \\
& \wedge Inv_diam(20)
\end{aligned}$$

Examples:

Find a behavior with a diameter of size 60.

$$Inv_diam(60)$$

Find a behavior where two different monitors assume the role of a leader.

$$\begin{aligned}
& Inv_find_state(\\
& \quad \exists msg1, msg2 \in message_history : \\
& \quad \quad \wedge msg1.type = OP_COLLECT \wedge msg2.type = OP_COLLECT \\
& \quad \quad \wedge msg1.from \neq msg2.from \\
&)
\end{aligned}$$

Find a state where a monitor crashed during the collect phase and fails to send a *OP_LAST* message.

```

Inv_find_state(
  ∧ step_x = "crash mon"

  \ * The system is in collect phase and no OP_LAST message has been received.
  \ * isLeader[mon] = TRUE assures that the leader was not the one that crashed.
  ∧ ∃ mon ∈ Monitors :
    ∧ isLeader[mon] = TRUE
    ∧ phase[mon] = PHASE_COLLECT
    ∧ num_last[mon] = 1

  \ * All the collect requests have been handled by the peers.
  ∧ ∀ mon1, mon2 ∈ Monitors :
    ∀ i ∈ 1 .. Len(messages[mon1][mon2]) : messages[mon1][mon2][i].type ≠ OP_COLLECT

  ∧ epoch = 2
)

```

Find a state where the leader crashes during the commit phase, failing to complete the commit.

```

Inv_find_state(
  ∧ step_x = "crash mon"
  ∧ ∃ mon1, mon2 ∈ Monitors :
    ∃ i ∈ 1 .. Len(messages[mon1][mon2]) : messages[mon1][mon2][i].type = OP_ACCEPT
  ∧ ∀ mon ∈ Monitors :
    isLeader[mon] = FALSE
  ∧ epoch = 2
)

```

Note: After finding a state, that complete state can be used as an initial state to analyze behaviors from there.

```

\ * Modification History
\ * Last modified Wed Mar 17 14:56:41 WET 2021 by afonsonf
\ * Created Mon Jan 11 16:15:26 WET 2021 by afonsonf

```