———————————————— MODULE *paxos* ————————————————

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/afonsonf/ceph-consensus-spec

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

External libraries used on:
*SequencesExt*: *SetToSeq*
*FiniteSetsExt*: *Min*, *Max*

## Constants

Set of *Monitors*.
CONSTANTS *Monitors*

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

Rank predicate, used to compute proposal numbers.
$rank(mon) \triangleq \text{CHOOSE } i \in 1 .. MonitorsLen : MonitorsSeq[i] = mon$

Set of possible values.
CONSTANTS *Value_set*

Predicate used in the cfg file to define the symmetry set.
$SYMM \triangleq Permutations(Monitors) \cup Permutations(Value\_set)$

Reserved value.
CONSTANTS *Nil*

CONSTANTS $STATE\_RECOVERING$, $STATE\_ACTIVE$,
$\qquad STATE\_UPDATING$, $STATE\_UPDATING\_PREVIOUS$,
$\qquad STATE\_WRITING$, $STATE\_WRITING\_PREVIOUS$,
$\qquad STATE\_REFRESH$, $STATE\_SHUTDOWN$

$state\_names \triangleq \{STATE\_RECOVERING, STATE\_ACTIVE,$
$\qquad STATE\_UPDATING, STATE\_UPDATING\_PREVIOUS,$
$\qquad STATE\_WRITING, STATE\_WRITING\_PREVIOUS,$
$\qquad STATE\_REFRESH, STATE\_SHUTDOWN\}$

CONSTANTS $PHASE\_ELECTION,$
$\qquad PHASE\_SEND\_COLLECT, PHASE\_COLLECT,$
$\qquad PHASE\_LEASE, PHASE\_LEASE\_DONE,$
$\qquad PHASE\_BEGIN,$
$\qquad PHASE\_COMMIT$

$phase\_names \triangleq \{PHASE\_ELECTION,$
$\qquad PHASE\_SEND\_COLLECT, PHASE\_COLLECT,$
$\qquad PHASE\_LEASE, PHASE\_LEASE\_DONE,$
$\qquad PHASE\_BEGIN,$
$\qquad PHASE\_COMMIT\}$

CONSTANTS $OP\_COLLECT$, $OP\_LAST$,
$\qquad OP\_BEGIN$, $OP\_ACCEPT$, $OP\_COMMIT$,
$\qquad OP\_LEASE$, $OP\_LEASE\_ACK$

$messages\_types \triangleq \{OP\_COLLECT, OP\_LAST,$
$\qquad OP\_BEGIN, OP\_ACCEPT, OP\_COMMIT,$
$\qquad OP\_LEASE, OP\_LEASE\_ACK\}$

**Global variables**

VARIABLE $epoch$

VARIABLE $messages$

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 proposal number when the commit phase starts.
This proposal number can be retrieved after a monitor crashes and restarts.

Type: $[Monitors \mapsto$ proposal number]

VARIABLE *uncommitted_pn*

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: $[Monitors \mapsto$ 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: $[Monitors \mapsto$ proposal number$]$
VARIABLE *accepted_pn*

A function that stores the first value version committed by each monitor.
Type: $[Monitors \mapsto$ value version$]$
VARIABLE *first_committed*

A function that stores the last value version committed by each monitor.
Type: $[Monitors \mapsto$ value version$]$
VARIABLE *last_committed*

## Collect phase variables

A function that stores the number of peers that accepted a collect request.
Type: $[Monitors \mapsto$ number of peers that accepted$]$
VARIABLE *num_last*

Used by leader when receiving responses in collect phase.
Type: $[Monitors \mapsto [Monitors \mapsto$ value version$]]$
VARIABLE *peer_first_committed*

Used by leader when receiving responses in collect phase.
Type: $[Monitors \mapsto [Monitors \mapsto$ 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: $[Monitord \mapsto [Monitors \mapsto Bool]]$
VARIABLE *accepted*

## Debug variables

Variables to help debug a behavior.
step is the diameter of a behavior/path.
*step_name* the current predicate being called.
VARIABLE *step_name*

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**

$global\_vars \triangleq \langle epoch,\ messages,\ message\_history,\ quorum,\ quorum\_sz \rangle$
$state\_vars \triangleq \langle isLeader,\ state,\ phase \rangle$
$restart\_vars \triangleq \langle uncommitted\_pn,\ uncommitted\_v,\ uncommitted\_value \rangle$
$data\_vars \triangleq \langle monitor\_store,\ values,\ accepted\_pn,\ first\_committed,\ last\_committed \rangle$
$collect\_vars \triangleq \langle num\_last,\ peer\_first\_committed,\ peer\_last\_committed \rangle$
$lease\_vars \triangleq \langle acked\_lease \rangle$
$commit\_vars \triangleq \langle pending\_proposal,\ new\_value,\ accepted \rangle$

$vars \triangleq \langle global\_vars,\ state\_vars,\ restart\_vars,\ data\_vars,\ collect\_vars,$
$\qquad\qquad lease\_vars,\ commit\_vars \rangle$

$Init\_global\_vars \triangleq$
$\qquad \wedge epoch = 1$
$\qquad \wedge messages = [mon1 \in Monitors \mapsto [mon2 \in Monitors \mapsto \langle\rangle]]$
$\qquad \wedge message\_history = \{\}$
$\qquad \wedge quorum = [mon \in Monitors \mapsto \text{TRUE}]$
$\qquad \wedge quorum\_sz = MonitorsLen$

$Init\_state\_vars \triangleq$
$\qquad \wedge isLeader = [mon \in Monitors \mapsto \text{FALSE}]$
$\qquad \wedge state\ = [mon \in Monitors \mapsto Nil]$
$\qquad \wedge phase = [mon \in Monitors \mapsto Nil]$

$Init\_restart\_vars \triangleq$
$\qquad \wedge uncommitted\_pn = [mon \in Monitors \mapsto 0]$
$\qquad \wedge uncommitted\_v = [mon \in Monitors \mapsto 0]$
$\qquad \wedge uncommitted\_value = [mon \in Monitors \mapsto Nil]$

$Init\_data\_vars \triangleq$
$\qquad \wedge monitor\_store = [mon \in Monitors \mapsto Nil]$
$\qquad \wedge values = [mon \in Monitors \mapsto [version \in \{\} \mapsto Nil]]$
$\qquad \wedge accepted\_pn = [mon \in Monitors \mapsto 0]$
$\qquad \wedge first\_committed = [mon \in Monitors \mapsto 0]$
$\qquad \wedge last\_committed = [mon \in Monitors \mapsto 0]$

$Init\_collect\_vars \triangleq$
$\qquad \wedge num\_last = [mon \in Monitors \mapsto 0]$

5

$$\wedge \; peer\_first\_committed = [mon1 \in Monitors \mapsto [mon2 \in Monitors \mapsto \; -1]]$$
$$\wedge \; peer\_last\_committed = [mon1 \in Monitors \mapsto [mon2 \in Monitors \mapsto \; -1]]$$

$Init\_lease\_vars \;\triangleq$
$$\wedge \; acked\_lease = [mon1 \in Monitors \mapsto [mon2 \in Monitors \mapsto \text{FALSE}]]$$

$Init\_commit\_vars \;\triangleq$
$$\wedge \; pending\_proposal \;\; = [mon \in Monitors \mapsto Nil]$$
$$\wedge \; new\_value = [mon \in Monitors \mapsto Nil]$$
$$\wedge \; accepted = [mon1 \;\; \in Monitors \mapsto [mon2 \in Monitors \mapsto \text{FALSE}]]$$

$Init \;\triangleq$
$$\wedge \; Init\_global\_vars$$
$$\wedge \; Init\_state\_vars$$
$$\wedge \; Init\_restart\_vars$$
$$\wedge \; Init\_data\_vars$$
$$\wedge \; Init\_collect\_vars$$
$$\wedge \; Init\_lease\_vars$$
$$\wedge \; Init\_commit\_vars$$
$$\wedge \; step\_name = \text{"init"} \wedge number\_crashes = 0$$

## Message manipulation

Note: Variable $message\_history$ has impact in performace, update only when debugging.

Converts a set with at most one element to a sequence.
$SingleMessageSetToSeq(S) \;\triangleq$
   IF $\exists \, elem \in S : \text{TRUE}$ THEN LET $elem \;\triangleq$ CHOOSE $x \in S : \text{TRUE}$
   IN $\quad \langle elem \rangle$
   ELSE $\langle \rangle$

Add message $m$ to the network $msgs$.
$WithMessage(m, \, msgs) \;\triangleq$
   $[msgs$ EXCEPT $![m.from] =$
      $[msgs[m.from]$ EXCEPT $![m.dest] = Append(msgs[m.from][m.dest], \, m)]]$

Remove message $m$ from the network $msgs$.
$WithoutMessage(m, \, msgs) \;\triangleq$
   $[msgs$ EXCEPT $![m.from] =$
      $[msgs[m.from]$ EXCEPT $![m.dest] = Tail(msgs[m.from][m.dest])]]$

Adds the message $m$ to the network.
Variables changed: messages, $message\_history$.
$Send(m) \;\triangleq$
   $\wedge \; messages' = WithMessage(m, \, messages)$
   $\wedge \; message\_history' = message\_history \cup \{m\}$
   $\wedge$ UNCHANGED $message\_history$

Adds a set of messages to the network.
Variables changed: messages, $message\_history$.

$Send\_set(from, m\_set) \triangleq$
$\quad \land messages' = [messages \text{ EXCEPT } ![from] =$
$\quad\quad [mon \in Monitors \mapsto$
$\quad\quad\quad messages[from][mon] \circ SingleMessageSetToSeq(\{m \in m\_set : m.dest = mon\})]]$
$\quad\quad \land message\_history' = message\_history \cup m\_set$
$\quad \land \text{UNCHANGED } message\_history$

Removes the request from network and adds the response.
Variables changed: messages, $message\_history$.

$Reply(response, request) \triangleq$
$\quad \land messages' = WithoutMessage(request, WithMessage(response, messages))$
$\quad\quad \land message\_history' = message\_history \cup \{response\}$
$\quad \land \text{UNCHANGED } message\_history$

Removes the request from network and adds a set of messages.
Variables changed: messages, $message\_history$.

$Reply\_set(from, response\_set, request) \triangleq$
$\quad \land \text{LET } msgs \triangleq WithoutMessage(request, messages)$
$\quad\quad \text{IN} \quad messages' = [msgs \text{ EXCEPT } ![from] =$
$\quad\quad\quad [mon \in Monitors \mapsto$
$\quad\quad\quad\quad msgs[from][mon] \circ SingleMessageSetToSeq(\{m \in response\_set : m.dest = mon\})]]$
$\quad\quad \land message\_history' = message\_history \cup response\_set$
$\quad \land \text{UNCHANGED } message\_history$

Removes message $m$ from the network.
Variables changed: messages, $message\_history$.

$Discard(m) \triangleq$
$\quad \land \quad messages' = WithoutMessage(m, messages)$
$\quad \land \quad \text{UNCHANGED } message\_history$

**Helper predicates**

Computes a new unique proposal number for a given monitor.

Version A - Equal to the one in the source.
This version breaks the symmetry of the monitor set.
Example: $oldpn = 305$, $rank(mon) = 5$, $newpn = 405$.
$get\_new\_proposal\_number(mon, oldpn) \triangleq ((oldpn \div 100) + 1) * 100 + rank(mon)$

Version $B -$ Adapted to not break symmetry.
Example: $oldpn = 300$, $rank(mon) = 5$, $newpn = 400$.
$get\_new\_proposal\_number(mon, oldpn) \triangleq ((oldpn \div 100) + 1) * 100$

Clear the variable $peer\_first\_committed$.

$clear\_peer\_first\_committed(mon) \triangleq$
   $peer\_first\_committed' = [peer\_first\_committed$ EXCEPT $![mon] =$
                           $[m \in Monitors \mapsto -1]]$

$clear\_peer\_last\_committed(mon) \triangleq$
   $peer\_last\_committed' = [peer\_last\_committed$ EXCEPT $![mon] =$
                          $[m \in Monitors \mapsto -1]]$

$store\_state(mon, msg) \triangleq$

   $\wedge$ LET $logs \triangleq ($DOMAIN $msg.values) \cap (last\_committed[mon] + 1 \mathinner{\ldotp\ldotp} msg.last\_committed)$
     IN   $\wedge values' = [values$ EXCEPT $![mon] =$
             $[i \in$ DOMAIN $values[mon] \cup logs \mapsto$
                IF $i \in logs$
                THEN $msg.values[i]$
                ELSE $values[mon][i]]]$
          $\wedge last\_committed' = [last\_committed$ EXCEPT $![mon] = Max(logs \cup \{last\_committed[mon]\})]$
          $\wedge$ IF $logs \neq \{\} \wedge first\_committed[mon] = 0$
            THEN $first\_committed' =$
                       $[first\_committed$ EXCEPT $![mon] = Min(logs)]$
            ELSE $first\_committed' =$
                       $[first\_committed$ EXCEPT $![mon] = Min(logs \cup \{first\_committed[mon]\})]$
   $\wedge$ IF $last\_committed'[mon] = 0$
     THEN UNCHANGED $monitor\_store$
     ELSE $monitor\_store' = [monitor\_store$ EXCEPT $![mon] = values'[mon][last\_committed'[mon]]]$

$check\_and\_correct\_uncommitted(mon) \triangleq$
   IF $uncommitted\_v[mon] \leq last\_committed'[mon]$
     THEN $\wedge uncommitted\_v' = [uncommitted\_v$ EXCEPT $![mon] = 0]$
          $\wedge uncommitted\_pn' = [uncommitted\_pn$ EXCEPT $![mon] = 0]$
          $\wedge uncommitted\_value' = [uncommitted\_value$ EXCEPT $![mon] = Nil]$
     ELSE UNCHANGED $\langle uncommitted\_pn, uncommitted\_v, uncommitted\_value \rangle$

$bootstrap \triangleq$
   $\wedge epoch' = epoch + 1$

8

Changes $mon$ state to $STATE\_ACTIVE$.
Variables changed: state.
$finish\_round(mon) \stackrel{\Delta}{=}$
$\quad \land\ isLeader[mon] = \text{TRUE}$
$\quad \land\ state' = [state \text{ EXCEPT } ![mon] = STATE\_ACTIVE]$

Resets the variable acked lease and send lease messages to peers.
Variables changed: $acked\_lease$, messages, $message\_history$, phase.
$extend\_lease(mon) \stackrel{\Delta}{=}$
$\quad \land\ isLeader[mon] = \text{TRUE}$
$\quad \land\ acked\_lease' = [acked\_lease \text{ EXCEPT } ![mon] =$
$\quad\quad [m \in Monitors \mapsto \text{IF } m = mon \text{ THEN TRUE ELSE FALSE}]]$
$\quad \land\ Send\_set(mon,$
$\quad\quad \{[type \qquad\qquad \mapsto OP\_LEASE,$
$\quad\quad\ \ from \qquad\qquad \mapsto mon,$
$\quad\quad\ \ dest \qquad\qquad \mapsto dest,$
$\quad\quad\ \ last\_committed \mapsto last\_committed[mon]] : dest \in \{m \in Monitors \setminus \{mon\} : quorum[m]\}$
$\quad\quad\ \})$
$\quad \land\ phase' = [phase \text{ EXCEPT } ![mon] = PHASE\_LEASE]$

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.
$handle\_lease(mon,\ msg) \stackrel{\Delta}{=}$
$\quad \land\ $ discard if not peon or peon is behind
$\quad\quad \text{IF} \quad \lor\ isLeader[mon] = \text{TRUE}$
$\quad\quad\quad\quad \lor\ last\_committed[mon] \neq msg.last\_committed$
$\quad\quad \text{THEN} \quad \land\ Discard(msg)$
$\quad\quad\quad\quad\quad \land\ \text{UNCHANGED } state$
$\quad\quad \text{ELSE} \quad \land\ state' = [state \text{ EXCEPT } ![mon] = STATE\_ACTIVE]$
$\quad\quad\quad\quad \land\ Reply([type \qquad\quad \mapsto OP\_LEASE\_ACK,$
$\quad\quad\quad\quad\quad\ from \qquad\quad \mapsto mon,$
$\quad\quad\quad\quad\quad\ dest \qquad\quad \mapsto msg.from,$
$\quad\quad\quad\quad\quad\ first\_committed \mapsto first\_committed[mon],$
$\quad\quad\quad\quad\quad\ last\_committed \mapsto last\_committed[mon]],\ msg)$
$\quad\quad\quad\quad \land\ Discard(msg)$
$\quad \land\ \text{UNCHANGED } \langle epoch,\ quorum,\ quorum\_sz,\ isLeader,\ phase\rangle$
$\quad \land\ \text{UNCHANGED } \langle restart\_vars,\ data\_vars,\ collect\_vars,\ lease\_vars,\ commit\_vars\rangle$

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.

$handle\_lease\_ack(mon, msg) \triangleq$
  $\wedge phase[mon] = PHASE\_LEASE$
  $\wedge acked\_lease' = [acked\_lease$ EXCEPT $![mon] =$
   $[acked\_lease[mon]$ EXCEPT $![msg.from] =$ TRUE$]]$
  $\wedge Discard(msg)$
  $\wedge$ UNCHANGED $\langle epoch,\ quorum,\ quorum\_sz \rangle$
  $\wedge$ UNCHANGED $\langle state\_vars,\ restart\_vars,\ data\_vars,\ collect\_vars,\ commit\_vars \rangle$

$post\_lease\_ack(mon) \triangleq$
  $\wedge phase[mon] = PHASE\_LEASE$
  $\wedge phase' = [phase$ EXCEPT $![mon] = PHASE\_LEASE\_DONE]$
  $\wedge \forall m \in Monitors : quorum[m] \Rightarrow acked\_lease[mon][m] =$ TRUE
  $\wedge$ UNCHANGED $\langle isLeader,\ state \rangle$
  $\wedge$ UNCHANGED $\langle global\_vars,\ restart\_vars,\ data\_vars,\ collect\_vars,$
       $lease\_vars,\ commit\_vars \rangle$

**Commit phase predicates**

$begin(mon, v) \triangleq$
  $\wedge isLeader[mon] =$ 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 $![mon] =$
   $((last\_committed[mon] + 1) :> new\_value'[mon])$ @@ $values[mon]]$
  $\wedge Send\_set(mon,$
   $\{[type$      $\mapsto OP\_BEGIN,$
    $from$      $\mapsto mon,$
    $dest$      $\mapsto dest,$
    $last\_committed \mapsto last\_committed[mon],$
    $values$     $\mapsto values'[mon],$

$$pn \quad\quad\quad\quad\quad\quad \mapsto accepted\_pn[mon]] : dest \in \{m \in Monitors \setminus \{mon\} : quorum[m]\}$$
$$\})$$
$$\wedge uncommitted\_pn' = [uncommitted\_pn \text{ EXCEPT } ![mon] = accepted\_pn[mon]]$$
$$\wedge uncommitted\_v' = [uncommitted\_v \text{ EXCEPT } ![mon] = last\_committed[mon] + 1]$$
$$\wedge uncommitted\_value' = [uncommitted\_value \text{ EXCEPT } ![mon] = v]$$

$handle\_begin(mon, msg) \;\triangleq$
    $\wedge isLeader[mon] = \text{FALSE}$
    $\wedge \text{ IF } msg.pn < accepted\_pn[mon]$
        THEN
        $\wedge Discard(msg)$
        $\wedge \text{ UNCHANGED } \langle state, values, restart\_vars \rangle$
        ELSE
        $\wedge msg.pn = accepted\_pn[mon]$
        $\wedge msg.last\_committed = last\_committed[mon]$

        assign $values[mon][last\_committed[mon] + 1]$
        $\wedge values' = [values \text{ EXCEPT } ![mon] =$
          $((last\_committed[mon] + 1) :> msg.values[last\_committed[mon] + 1]) @@ values[mon]]$

        $\wedge state' = [state \text{ EXCEPT } ![mon] = STATE\_UPDATING]$
        $\wedge uncommitted\_pn' = [uncommitted\_pn \text{ EXCEPT } ![mon] = accepted\_pn[mon]]$
        $\wedge uncommitted\_v' = [uncommitted\_v \text{ EXCEPT } ![mon] = last\_committed[mon] + 1]$
        $\wedge uncommitted\_value' = [uncommitted\_value \text{ EXCEPT } ![mon] =$
          $values'[mon][last\_committed[mon] + 1]]$
        $\wedge Reply([type \quad\quad\quad \mapsto OP\_ACCEPT,$
               $from \quad\quad\quad\quad \mapsto mon,$
               $dest \quad\quad\quad\quad \mapsto msg.from,$
               $last\_committed \;\; \mapsto last\_committed[mon],$
               $pn \quad\quad\quad\quad\quad \mapsto accepted\_pn[mon]], msg)$
    $\wedge \text{ UNCHANGED } \langle epoch, quorum, quorum\_sz, isLeader, phase, monitor\_store,$
                $accepted\_pn, first\_committed, last\_committed \rangle$
    $\wedge \text{ UNCHANGED } \langle collect\_vars, lease\_vars, commit\_vars \rangle$

$handle\_accept(mon, msg) \;\triangleq$
    $\wedge isLeader[mon] = \text{TRUE}$
    $\wedge \vee state[mon] \;\; = STATE\_UPDATING\_PREVIOUS$
      $\vee state[mon] \;\; = STATE\_UPDATING$

$\land phase[mon] = PHASE\_BEGIN$
$\land new\_value[mon] \neq Nil$
$\land \text{IF} \quad \lor msg.pn \neq accepted\_pn[mon]$
$\qquad\quad \lor \land last\_committed[mon] > 0$
$\qquad\qquad\quad \land msg.last\_committed < last\_committed[mon] - 1$
$\quad \text{THEN} \text{ UNCHANGED } accepted$
$\quad \text{ELSE} \quad accepted' = [accepted \text{ EXCEPT } ![mon] =$
$\qquad\qquad\qquad [accepted[mon] \text{ EXCEPT } ![msg.from] = \text{TRUE}]]$
$\land Discard(msg)$
$\land \text{UNCHANGED } \langle epoch,\ quorum,\ quorum\_sz,\ pending\_proposal,\ new\_value \rangle$
$\land \text{UNCHANGED } \langle restart\_vars,\ state\_vars,\ data\_vars,\ collect\_vars,\ lease\_vars \rangle$

$post\_accept(mon) \triangleq$
$\quad \land phase[mon] = PHASE\_BEGIN$
$\quad \land \forall\, m \in Monitors : quorum[m] \Rightarrow accepted[mon][m] = \text{TRUE}$
$\quad \land new\_value[mon] \neq Nil$
$\quad \land \lor state[mon] = STATE\_UPDATING\_PREVIOUS$
$\qquad \lor state[mon] = STATE\_UPDATING$
$\quad \land last\_committed' = [last\_committed \text{ EXCEPT } ![mon] = last\_committed[mon] + 1]$

$\quad \land \text{IF } first\_committed[mon] = 0$
$\qquad \text{THEN } first\_committed' = [first\_committed \text{ EXCEPT } ![mon] = first\_committed[mon] + 1]$
$\qquad \text{ELSE} \text{ UNCHANGED } first\_committed$

$\quad \land monitor\_store' = [monitor\_store \text{ EXCEPT } ![mon] = values[mon][last\_committed[mon] + 1]]$
$\quad \land new\_value' = [new\_value \text{ EXCEPT } ![mon] = Nil]$
$\quad \land Send\_set(mon,$
$\qquad \{[type \qquad\qquad \mapsto OP\_COMMIT,$
$\qquad\quad from \qquad\qquad \mapsto mon,$
$\qquad\quad dest \qquad\qquad \mapsto dest,$
$\qquad\quad last\_committed \mapsto last\_committed'[mon],$
$\qquad\quad pn \qquad\qquad\quad \mapsto accepted\_pn[mon],$
$\qquad\quad values \qquad\qquad \mapsto values[mon]] : dest \in \{m \in Monitors \setminus \{mon\} : quorum[m]\}$
$\qquad \})$
$\quad \land state' = [state \text{ EXCEPT } ![mon] = STATE\_REFRESH]$
$\quad \land phase' = [phase \text{ EXCEPT } ![mon] = PHASE\_COMMIT]$
$\quad \land \text{UNCHANGED } \langle isLeader,\ values,\ accepted\_pn,\ pending\_proposal,\ accepted \rangle$
$\quad \land \text{UNCHANGED } \langle epoch,\ quorum,\ quorum\_sz,\ restart\_vars,\ collect\_vars,\ lease\_vars \rangle$

$finish\_commit(mon) \triangleq$

$\land state[mon]\ = STATE\_REFRESH$
$\land phase[mon] = PHASE\_COMMIT$
$\land finish\_round(mon)$
$\land extend\_lease(mon)$
$\land$ UNCHANGED $\langle epoch,\ quorum,\ quorum\_sz,\ isLeader\rangle$
$\land$ 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, $first\_committed$, $last\_committed$, $monitor\_store$, $uncommitted\_v$, $uncommitted\_pn$, $uncommitted\_value$.
$handle\_commit(mon,\ msg)\ \overset{\Delta}{=}$
    $\land isLeader[mon] = $ FALSE
    $\land store\_state(mon,\ msg)$
    $\land check\_and\_correct\_uncommitted(mon)$
    $\land Discard(msg)$
    $\land$ UNCHANGED $\langle epoch,\ quorum,\ quorum\_sz,\ accepted\_pn\rangle$
    $\land$ 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: $pending\_proposal$.
$client\_request(mon,\ v)\ \overset{\Delta}{=}$
    $\land isLeader[mon] = $ TRUE
    $\land state[mon] = STATE\_ACTIVE$
    $\land pending\_proposal[mon] = Nil$
    $\land pending\_proposal' = [pending\_proposal$ EXCEPT $![mon] = v]$
    $\land$ UNCHANGED $\langle new\_value,\ accepted\rangle$
    $\land$ UNCHANGED $\langle global\_vars,\ state\_vars,\ restart\_vars,\ data\_vars,\ collect\_vars,\ lease\_vars\rangle$

Start a commit phase with the value on pending proposal.
Variables changed: state, $pending\_proposal$, accepted, $new\_value$, phase, messages, $message\_history$, values, $uncommitted\_pn$, $uncommitted\_v$, $uncommitted\_value$.
$propose\_pending(mon)\ \overset{\Delta}{=}$
    $\land phase[mon] = PHASE\_LEASE \lor phase[mon] = PHASE\_ELECTION$
    $\land state[mon]\ = STATE\_ACTIVE$
    $\land pending\_proposal[mon] \neq Nil$
    $\land pending\_proposal' = [pending\_proposal$ EXCEPT $![mon] = Nil]$
    $\land state' = [state$ EXCEPT $![mon] = STATE\_UPDATING]$
    $\land begin(mon,\ pending\_proposal[mon])$
    $\land$ UNCHANGED $\langle isLeader,\ monitor\_store,\ accepted\_pn,\ first\_committed,\ last\_committed\rangle$
    $\land$ UNCHANGED $\langle epoch,\ quorum,\ quorum\_sz,\ collect\_vars,\ lease\_vars\rangle$

**Collect phase predicates**

13

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.

$collect(mon, \; oldpn) \; \triangleq$
    $\land \; state[mon] = STATE\_RECOVERING$
    $\land \; isLeader[mon] = \text{TRUE}$
    $\land \; \text{LET} \; new\_pn \; \triangleq \; get\_new\_proposal\_number(mon, \; Max(\{oldpn, \; accepted\_pn[mon]\}))$
       $\text{IN} \quad \land \; accepted\_pn' = [accepted\_pn \; \text{EXCEPT} \; ![mon] = new\_pn]$
    $\land \; phase' = [phase \; \text{EXCEPT} \; ![mon] = PHASE\_SEND\_COLLECT]$

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_pn*, *uncommitted_v*, *uncommitted_value*, *num_las*
messages, *message_history*, phase.

$send\_collect(mon) \; \triangleq$
    $\land \; state[mon] = STATE\_RECOVERING$
    $\land \; isLeader[mon] = \text{TRUE}$
    $\land \; phase[mon] = PHASE\_SEND\_COLLECT$
    $\land \; clear\_peer\_first\_committed(mon)$
    $\land \; clear\_peer\_last\_committed(mon)$

    $\land \; \text{IF} \; last\_committed[mon] + 1 \in \text{DOMAIN} \; values[mon]$
      $\text{THEN} \quad \land \; uncommitted\_v' =$
            $[uncommitted\_v \; \text{EXCEPT} \; ![mon] = last\_committed[mon] + 1]$
          $\land \; uncommitted\_value' =$
            $[uncommitted\_value \; \text{EXCEPT} \; ![mon] = values[mon][last\_committed[mon] + 1]]$
          $\land \; uncommitted\_pn' = uncommitted\_pn$
      $\text{ELSE} \quad \text{UNCHANGED} \; \langle restart\_vars \rangle$

    $\land \; num\_last' = [num\_last \; \text{EXCEPT} \; ![mon] = 1]$
    $\land \; Send\_set(mon,$
      $\{[type \qquad\qquad \mapsto OP\_COLLECT,$
        $from \qquad\qquad \mapsto mon,$
        $dest \qquad\qquad \mapsto dest,$
        $first\_committed \mapsto first\_committed[mon],$
        $last\_committed \; \mapsto last\_committed[mon],$
        $pn \qquad\qquad\quad \mapsto accepted\_pn[mon]] : dest \in \{m \in Monitors \setminus \{mon\} : quorum[m]\}$
      $\})$
    $\land \; phase' = [phase \; \text{EXCEPT} \; ![mon] = PHASE\_COLLECT]$
    $\land \; \text{UNCHANGED} \; \langle isLeader, \; state \rangle$
    $\land \; \text{UNCHANGED} \; \langle epoch, \; quorum, \; quorum\_sz, \; data\_vars, \; lease\_vars, \; commit\_vars \rangle$

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*

$handle\_collect(mon,\ msg) \triangleq$
 $\land\ isLeader[mon] = \text{FALSE}$
 $\land\ state' = [state \text{ EXCEPT } ![mon] = STATE\_RECOVERING]$
 $\land\ \lor\ \land\ msg.first\_committed > last\_committed[mon] + 1$
   $\land\ bootstrap$
   $\land\ Discard(msg)$
   $\land\ \text{UNCHANGED } \langle accepted\_pn \rangle$
  $\lor\ \land\ msg.first\_committed \leq last\_committed[mon] + 1$
   $\land\ \text{IF } msg.pn > accepted\_pn[mon]$
    $\text{THEN } accepted\_pn' = [accepted\_pn \text{ EXCEPT } ![mon] = msg.pn]$
    $\text{ELSE } \text{UNCHANGED } accepted\_pn$
   $\land\ Reply([type\qquad\qquad\quad \mapsto OP\_LAST,$
      $from\qquad\qquad\quad \mapsto mon,$
      $dest\qquad\qquad\quad \mapsto msg.from,$
      $first\_committed\ \mapsto first\_committed[mon],$
      $last\_committed\ \ \mapsto last\_committed[mon],$
      $values\qquad\qquad \mapsto values[mon],$
      $uncommitted\_pn\ \mapsto uncommitted\_pn[mon],$
      $pn\qquad\qquad\quad \mapsto accepted\_pn'[mon]],\ msg)$
   $\land\ \text{UNCHANGED } epoch$
 $\land\ \text{UNCHANGED } \langle isLeader,\ phase,\ values,\ first\_committed,\ last\_committed,\ monitor\_store \rangle$
 $\land\ \text{UNCHANGED } \langle quorum,\ quorum\_sz,\ restart\_vars,\ collect\_vars,\ lease\_vars,\ commit\_vars \rangle$

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, $uncommitted\_pn$, $uncommitted\_v$, $uncommitted\_value$, $monitor\_s$
$accepted\_pn$, $first\_committed$, $last\_committed$, $num\_last$, $peer\_first\_committed$, $peer\_last\_committed$.

$handle\_last(mon,\ msg) \triangleq$
 $\land\ isLeader[mon] = \text{TRUE}$

 $\land\ peer\_first\_committed' = [peer\_first\_committed \text{ EXCEPT } ![mon] =$
  $[peer\_first\_committed[mon] \text{ EXCEPT } ![msg.from] = msg.first\_committed]]$
 $\land\ peer\_last\_committed' = [peer\_last\_committed \text{ EXCEPT } ![mon] =$
  $[peer\_last\_committed[mon] \text{ EXCEPT } ![msg.from] = msg.last\_committed]]$

 $\land\ \text{IF } msg.first\_committed > last\_committed[mon] + 1$
  $\text{THEN}$
  $\land\ bootstrap$
  $\land\ Discard(msg)$
  $\land\ \text{UNCHANGED } \langle num\_last,\ accepted\_pn,\ values,\ phase,\ monitor\_store \rangle$
  $\land\ \text{UNCHANGED } \langle first\_committed,\ last\_committed,\ restart\_vars \rangle$
  $\text{ELSE}$
  $\land\ store\_state(mon,\ msg)$
  $\land\ \text{IF } \exists\, peer \in Monitors :$

$\land$ *peer* $\neq$ *mon*
$\land$ *peer_last_committed$'$*[*mon*][*peer*] $\neq -1$
$\land$ *peer_last_committed$'$*[*mon*][*peer*] $+ 1 <$ *first_committed*[*mon*]
$\land$ *first_committed*[*mon*] $> 1$
THEN
$\land$ *bootstrap*
$\land$ *check_and_correct_uncommitted*(*mon*)
$\land$ *Discard*(*msg*)
$\land$ UNCHANGED $\langle$*phase*, *accepted_pn*, *num_last*$\rangle$
ELSE
$\land$ LET *monitors_behind* $\triangleq$ {*peer* $\in$ *Monitors* :
$\quad\land$ *peer* $\neq$ *mon*
$\quad\land$ *peer_last_committed$'$*[*mon*][*peer*] $\neq -1$
$\quad\land$ *peer_last_committed$'$*[*mon*][*peer*] $<$ *last_committed*[*mon*]
$\quad\land$ *quorum*[*peer*]}
IN  *Reply_set*(*mon*,
$\quad${[*type* $\qquad\qquad \mapsto OP\_COMMIT$,
$\quad\;\;$ *from* $\qquad\qquad \mapsto$ *mon*,
$\quad\;\;$ *dest* $\qquad\qquad \mapsto$ *dest*,
$\quad\;\;$ *last_committed* $\mapsto$ *last_committed$'$*[*mon*],
$\quad\;\;$ *pn* $\qquad\qquad \mapsto$ *accepted_pn*[*mon*],
$\quad\;\;$ *values* $\qquad\quad \mapsto$ *values*[*mon*]] : *dest* $\in$ *monitors_behind*
$\quad$}, *msg*)
$\land$ $\lor$ $\land$ *msg.pn* $>$ *accepted_pn*[*mon*]
$\quad\;\;\land$ *collect*(*mon*, *msg.pn*)
$\quad\;\;\land$ *check_and_correct_uncommitted*(*mon*)
$\quad\;\;\land$ UNCHANGED *num_last*

$\quad\lor$ $\land$ *msg.pn* $=$ *accepted_pn*[*mon*]
$\quad\;\;\land$ *num_last$'$* $=$ [*num_last* EXCEPT ![*mon*] $=$ *num_last*[*mon*] $+ 1$]
$\quad\;\;\land$ IF $\land$ *msg.last_committed* $+ 1 \in$ DOMAIN *msg.values*
$\qquad\quad\;\;\land$ *msg.last_committed* $\geq$ *last_committed$'$*[*mon*]
$\qquad\quad\;\;\land$ *msg.last_committed* $+ 1 \geq$ *uncommitted_v*[*mon*]
$\qquad\quad\;\;\land$ *msg.uncommitted_pn* $\geq$ *uncommitted_pn*[*mon*]
$\qquad$THEN $\land$ *uncommitted_v$'$* $=$
$\qquad\qquad\qquad$ [*uncommitted_v* EXCEPT ![*mon*] $=$ *msg.last_committed* $+ 1$]
$\qquad\qquad\;\;\land$ *uncommitted_pn$'$* $=$
$\qquad\qquad\qquad$ [*uncommitted_pn* EXCEPT ![*mon*] $=$ *msg.uncommitted_pn*]
$\qquad\qquad\;\;\land$ *uncommitted_value$'$* $=$
$\qquad\qquad\qquad$ [*uncommitted_value* EXCEPT ![*mon*] $=$ *msg.values*[*msg.last_committed* $+ 1$]]
$\qquad$ELSE *check_and_correct_uncommitted*(*mon*)
$\quad\;\;\land$ UNCHANGED $\langle$*phase*, *accepted_pn*$\rangle$

$\quad\lor$ $\land$ *msg.pn* $<$ *accepted_pn*[*mon*]
$\quad\;\;\land$ *check_and_correct_uncommitted*(*mon*)

16

$\qquad\qquad \land$ UNCHANGED $\langle phase,\ accepted\_pn,\ num\_last \rangle$
$\qquad\quad \land$ UNCHANGED $epoch$
$\qquad \land$ UNCHANGED $\langle epoch \rangle$

$\quad \land$ UNCHANGED $\langle quorum,\ quorum\_sz,\ isLeader,\ state \rangle$
$\quad \land$ 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, $uncommitted\_pn$, $uncommitted\_v$, $uncommitted\_value$, $acked\_lease$.
$post\_last(mon) \triangleq$
$\quad \land isLeader[mon] =$ TRUE
$\quad \land num\_last[mon] = quorum\_sz$
$\quad \land phase[mon] = PHASE\_COLLECT$

$\quad \land clear\_peer\_first\_committed(mon)$
$\quad \land clear\_peer\_last\_committed(mon)$

$\quad \land$ IF $\ \land uncommitted\_v[mon] = last\_committed[mon] + 1$
$\qquad\quad\ \land uncommitted\_value[mon] \neq Nil$
$\qquad$ THEN $\ \land state' = [state$ EXCEPT $![mon] = STATE\_UPDATING\_PREVIOUS]$
$\qquad\qquad\quad \land begin(mon,\ uncommitted\_value[mon])$
$\qquad\qquad\quad \land$ UNCHANGED $\langle acked\_lease \rangle$
$\qquad$ ELSE $\ \land finish\_round(mon)$
$\qquad\qquad\quad \land extend\_lease(mon)$
$\qquad\qquad\quad \land$ UNCHANGED $\langle accepted,\ new\_value,\ values,\ restart\_vars \rangle$

$\quad \land$ UNCHANGED $\langle isLeader,\ monitor\_store,\ accepted\_pn,\ first\_committed,\ last\_committed \rangle$
$\quad \land$ 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.
$leader\_election \triangleq$
$\quad \land \exists\, mon \in Monitors :$
$\qquad \land quorum[mon]$
$\qquad \land isLeader' = [m \in Monitors \mapsto$ IF $m = mon$ THEN TRUE ELSE FALSE$]$
$\qquad \land state' = [m \in Monitors \mapsto$
$\qquad\quad$ IF $quorum\_sz = 1$ THEN $STATE\_ACTIVE$ ELSE $STATE\_RECOVERING]$
$\quad \land phase' = [m \in Monitors \mapsto PHASE\_ELECTION]$
$\quad \land new\_value' = [m \in Monitors \mapsto Nil]$
$\quad \land pending\_proposal' = [m \in Monitors \mapsto Nil]$
$\quad \land epoch' = epoch + 1$

$\wedge\ messages' = [mon1 \in Monitors \mapsto [mon2 \in Monitors \mapsto \langle\rangle]]$

$\wedge$ UNCHANGED $\langle quorum,\ quorum\_sz,\ accepted,\ message\_history\rangle$

$\wedge$ 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.

$election\_recover(mon) \stackrel{\Delta}{=}$

$\quad \wedge\ quorum\_sz > 1$

$\quad \wedge\ phase[mon] = PHASE\_ELECTION$

$\quad \wedge\ collect(mon, 0)$

$\quad \wedge$ UNCHANGED $\langle isLeader,\ state,\ values,\ first\_committed,\ last\_committed,\ monitor\_store\rangle$

$\quad \wedge$ UNCHANGED $\langle global\_vars,\ restart\_vars,\ collect\_vars,\ lease\_vars,\ commit\_vars\rangle$

**Timeouts and restart**

$crash\_mon(mon) \stackrel{\Delta}{=}$

$\quad \wedge\ quorum\_sz > (MonitorsLen \div 2) + 1$

$\quad \wedge\ quorum[mon] = \text{TRUE}$

$\quad \wedge\ quorum' = [quorum\ \text{EXCEPT}\ ![mon] = \text{FALSE}]$

$\quad \wedge\ quorum\_sz' = quorum\_sz - 1$

$\quad \wedge\ bootstrap$

$\quad \wedge\ number\_crashes' = number\_crashes + 1$

$\quad \wedge$ UNCHANGED $\langle messages,\ message\_history\rangle$

$\quad \wedge$ UNCHANGED $\langle state\_vars,\ restart\_vars,\ data\_vars,\ collect\_vars,\ lease\_vars,\ commit\_vars\rangle$

$restore\_mon(mon) \stackrel{\Delta}{=}$

$\quad \wedge\ quorum[mon] = \text{FALSE}$

$\quad \wedge\ quorum' = [quorum\ \text{EXCEPT}\ ![mon] = \text{TRUE}]$

$\quad \wedge\ quorum\_sz' = quorum\_sz + 1$

$\quad \wedge\ bootstrap$

$\quad \wedge$ UNCHANGED $\langle messages,\ message\_history\rangle$

$\quad \wedge$ UNCHANGED $\langle state\_vars,\ restart\_vars,\ data\_vars,\ collect\_vars,\ lease\_vars,\ commit\_vars\rangle$

Monitor timeout (simulate the various timeouts that can occur). Triggers new elections.
Variables changed: epoch.

$Timeout(mon) \stackrel{\Delta}{=}$

$\quad \wedge\quad bootstrap$

$\quad \wedge\quad$ UNCHANGED $\langle messages,\ quorum,\ quorum\_sz,\ message\_history,\ state\_vars,\ restart\_vars,$
$\qquad\qquad\qquad data\_vars,\ collect\_vars,\ lease\_vars,\ commit\_vars\rangle$

**Dispatchers and next statement**

Handle a message.

$Receive(msg) \stackrel{\Delta}{=}$

$\quad \wedge\quad \vee\ \wedge\ msg.type = OP\_COLLECT$

18

$\land handle\_collect(msg.dest, \; msg)$
$\land step\_name' = \text{``receive collect''}$

$\lor \; \land msg.type = OP\_LAST$
$\quad \land handle\_last(msg.dest, \; msg)$
$\quad \land step\_name' = \text{``receive last''}$

$\lor \; \land msg.type = OP\_LEASE$
$\quad \land handle\_lease(msg.dest, \; msg)$
$\quad \land step\_name' = \text{``receive lease''}$

$\lor \; \land msg.type = OP\_LEASE\_ACK$
$\quad \land handle\_lease\_ack(msg.dest, \; msg)$
$\quad \land step\_name' = \text{``receive lease\_ack''}$

$\lor \; \land msg.type = OP\_BEGIN$
$\quad \land handle\_begin(msg.dest, \; msg)$
$\quad \land step\_name' = \text{``receive begin''}$

$\lor \; \land msg.type = OP\_ACCEPT$
$\quad \land handle\_accept(msg.dest, \; msg)$
$\quad \land step\_name' = \text{``receive accept''}$

$\lor \; \land msg.type = OP\_COMMIT$
$\quad \land handle\_commit(msg.dest, \; msg)$
$\quad \land step\_name' = \text{``receive commit''}$

Limit some variables to reduce search space.
$reduce\_search\_space \;\triangleq$
$\quad \land epoch \neq 8$
$\quad \land \lor \forall mon \in Monitors : last\_committed[mon] < 2$
$\qquad\quad \lor \forall mon2 \in Monitors\colon new\_value[mon2] = Nil$
$\quad \land \forall mon \in Monitors : accepted\_pn[mon] < 300$
$\quad \land number\_crashes \neq 4$

State transitions.
$Next \;\triangleq$
$\quad \land reduce\_search\_space$
$\quad \land \text{IF } epoch\%2 = 1 \text{ THEN}$
$\qquad \land leader\_election$
$\qquad \land step\_name' = \text{``election''}$
$\qquad \land \text{UNCHANGED } number\_crashes$
$\quad\;\; \text{ELSE}$
$\qquad \lor \; \land \exists mon \in Monitors : election\_recover(mon)$
$\qquad\quad\;\; \land step\_name' = \text{``election\_recover''}$
$\qquad\quad\;\; \land \text{UNCHANGED } number\_crashes$

$\qquad \lor \; \land \exists mon \in Monitors : send\_collect(mon)$

19

$\qquad \land step\_name' =$ "send_collect"
$\qquad \land$ UNCHANGED $number\_crashes$

$\lor \ \land \exists\, mon \in Monitors : post\_last(mon)$
$\qquad \land step\_name' =$ "post_last"
$\qquad \land$ UNCHANGED $number\_crashes$

$\lor \ \land \exists\, mon \in Monitors : post\_lease\_ack(mon)$
$\qquad \land step\_name' =$ "post_lease_ack"
$\qquad \land$ UNCHANGED $number\_crashes$

$\lor \ \land \exists\, mon \in Monitors : post\_accept(mon)$
$\qquad \land step\_name' =$ "post_accept"
$\qquad \land$ UNCHANGED $number\_crashes$

$\lor \ \land \exists\, mon \in Monitors : finish\_commit(mon)$
$\qquad \land step\_name' =$ "finish_commit"
$\qquad \land$ UNCHANGED $number\_crashes$

$\lor \ \land \exists\, mon \in Monitors : \exists\, v \in Value\_set : client\_request(mon,\, v)$
$\qquad \land step\_name' =$ "client_request"
$\qquad \land$ UNCHANGED $number\_crashes$

$\lor \ \land \exists\, mon \in Monitors : propose\_pending(mon)$
$\qquad \land step\_name' =$ "propose_pending"
$\qquad \land$ UNCHANGED $number\_crashes$

$\lor \ \land \exists\, mon1,\, mon2 \in Monitors :$
$\qquad\quad \land mon1 \neq mon2$
$\qquad\quad \land Len(messages[mon1][mon2]) > 0$
$\qquad\quad \land Receive(messages[mon1][mon2][1])$
$\qquad \land$ UNCHANGED $number\_crashes$

$\lor \ \land \exists\, mon \in Monitors : crash\_mon(mon)$
$\qquad \land step\_name' =$ "crash_mon"
$\qquad \land$ UNCHANGED $number\_crashes$

$\lor \ \land \exists\, mon \in Monitors : restore\_mon(mon)$
$\qquad \land step\_name' =$ "restore_mon"
$\qquad \land$ UNCHANGED $number\_crashes$

$\lor \ \land \exists\, mon \in Monitors : Timeout(mon)$
$\qquad \land step\_name' =$ "timeout_and_restart"
$\qquad \land$ UNCHANGED $number\_crashes$

**Safety invariants**

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

$same\_monitor\_store \;\triangleq\; \forall\, mon1,\, mon2 \in Monitors:$
$\quad state[mon1] = STATE\_ACTIVE \wedge state[mon2] = STATE\_ACTIVE$
$\quad\quad \Rightarrow monitor\_store[mon1] = monitor\_store[mon2]$

$Inv \;\triangleq\; \wedge same\_monitor\_store$

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'.
$Inv\_diam(size) \;\triangleq\; TLCGet(\text{``level''}) \neq size - 1$

Invariants to test in model check
$DEBUG\_Inv \;\triangleq\; \wedge \text{TRUE}$
$\quad\quad\quad\quad\quad\quad \wedge Inv\_diam(20)$

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.
$Inv\_find\_state($
$\quad \exists\, msg1,\, msg2 \in message\_history:$
$\quad\quad \wedge msg1.type = OP\_COLLECT \wedge msg2.type = OP\_COLLECT$
$\quad\quad \wedge msg1.from \neq msg2.from$
$)$

Find a state where a monitor crashed during the collect phase and fails to send a $OP\_LAST$ message.
$Inv\_find\_state($
$\quad \wedge step\_name = \text{``crash mon''}$

$\quad \setminus * \textit{The system is in collect phase and no OP\_LAST message has been received.}$
$\quad \setminus * isLeader[mon] = \text{TRUE } \textit{assures that the leader was not the one that crashed.}$
$\quad \wedge \exists\, mon \in Monitors:$
$\quad\quad \wedge isLeader[mon] = \text{TRUE}$
$\quad\quad \wedge phase[mon] = PHASE\_COLLECT$
$\quad\quad \wedge num\_last[mon] = 1$

$\quad \setminus * \textit{All the collect requests have been handled by the peers.}$
$\quad \wedge \forall\, mon1,\, mon2 \in Monitors:$
$\quad\quad \forall\, i \in 1\,..\,Len(messages[mon1][mon2]): messages[mon1][mon2][i].type \neq OP\_COLLECT$

$\quad \wedge epoch = 2$
$)$

Find a state where the leader crashes during the commit phase, failing to complete the commit.
$Inv\_find\_state($
$\quad \wedge step\_name = \text{``crash mon''}$
$\quad \wedge \exists\, mon1,\, mon2 \in Monitors:$

$\exists\, i \in 1 \mathinner{.\,.} Len(messages[mon1][mon2]) : messages[mon1][mon2][i].type = OP\_ACCEPT$
$\land\, \forall\, mon \in Monitors :$
$\quad isLeader[mon] = \text{FALSE}$
$\land\, epoch = 2$
$)$

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 Apr* 14 14:21:13 WEST 2021 by *afonsonf*
\ * Created *Mon Jan* 11 16:15:26 WET 2021 by *afonsonf*