

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 deviations/abstractions done that may differ from the implementation are:

- The election logic. The leader is chosen randomly, and, for now, only one leader is chosen per epoch.
- The quorum of monitors. For now, the specification considers the quorum to be the set of all monitors and that the quorum does not change over time.
- The communication layer. The variable *messages* holds the messages waiting to be handled. For now, messages cannot be randomly duplicated nor lost, and some messages can be received out of order.
- The transactions. In this specification, transactions represent only a change of value in the variable *monitor_store*.
- Failure model. For now, if a monitor crashes it will instantly restart, resetting some variables and continuing to participate in the quorum.

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

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

Constants

Set of monitors.

CONSTANTS *Monitors*

Sequence of monitors and the rank predicate, used to compute proposal numbers.

$\text{ranks} \triangleq \text{SetToSeq}(\text{Monitors})$

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

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*

$\text{state_names} \triangleq \{\text{STATE_RECOVERING}, \text{STATE_ACTIVE},$
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_PRE_COLLECT, *PHASE_COLLECT*,
 PHASE_LEASE, *PHASE_LEASE_DONE*,
 PHASE_BEGIN, *PHASE_BEGIN_DONE*,
 PHASE_COMMIT, *PHASE_COMMIT_DONE*

phase_names \triangleq {*PHASE_ELECTION*,
 PHASE_PRE_COLLECT, *PHASE_COLLECT*,
 PHASE_LEASE, *PHASE_LEASE_DONE*,
 PHASE_BEGIN, *PHASE_BEGIN_DONE*,
 PHASE_COMMIT, *PHASE_COMMIT_DONE*}

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*

A function that stores messages.

Type: $\langle message \rangle$

VARIABLE *messages*

Stores history of message events. Can be useful to find specific states.

Type: {*messages*}

VARIABLE *message_history*

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: $[Monitors \mapsto \text{Value_set}]$

VARIABLE *uncommitted_value*

Data variables

A function that stores, for each monitor, the current store where the transactions are applied.
In this model, a transaction represents changing the value in the store.

Type: $[Monitors \mapsto \text{Value_set}]$

VARIABLE *monitor_store*

A function that stores the transaction log of each monitor.

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

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 for each monitor.

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

VARIABLE *first_committed*

A function that stores the last value version committed for 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*

Auxiliary variables

A function that stores, for each monitor, a queue of messages types to send.
 Type: $[Monitors \mapsto \langle Monitors \times messages_types \rangle]$
 VARIABLE *send_queue*

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_refreshes*

Variables initialization

global_vars $\triangleq \langle epoch, messages, message_history \rangle$
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 $\triangleq \langle acked_lease \rangle$

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

$\wedge \text{Init_collect_vars}$
 $\wedge \text{Init_lease_vars}$
 $\wedge \text{Init_auxiliary_vars}$
 $\wedge \text{Init_commit_vars}$
 $\wedge \text{step} = 0 \wedge \text{step_x} = \text{"init"} \wedge \text{number_refreshes} = 0$

Message manipulation

Add message m to the network $msgs$.
 $\text{WithMessage}(m, msgs) \triangleq$
 $\text{Append}(msgs, m)$

Remove message m from the network $msgs$.
 $\text{WithoutMessage}(m, msgs) \triangleq$
 $\text{Remove}(msgs, m)$

Adds the message m to the network.
 Variable $message_history$ has impact in performace, update only when debugging.
 Variables changed: messages.
 $\text{Send}(m) \triangleq$
 $\wedge \text{messages}' = \text{WithMessage}(m, \text{messages})$
 $\wedge \text{message_history}' = \text{message_history} \cup \{m\}$
 $\wedge \text{UNCHANGED } message_history$

Removes message m from the network.
 Variables changed: messages.
 $\text{Discard}(m) \triangleq$
 $\wedge \text{messages}' = \text{WithoutMessage}(m, \text{messages})$
 $\wedge \text{UNCHANGED } message_history$

Removes the request from network and adds the response.
 Variables changed: messages.
 $\text{Reply}(\text{response}, \text{request}) \triangleq$
 $\wedge \text{messages}' = \text{WithoutMessage}(\text{request}, \text{WithMessage}(\text{response}, \text{messages}))$
 $\wedge \text{message_history}' = \text{message_history} \cup \{\text{response}\}$
 $\wedge \text{UNCHANGED } message_history$

Helper predicates

Compute a new unique proposal number for a given monitor.
 Example: $\text{oldpn} = 305, \text{rank}(\text{mon}) = 5, \text{newpn} = 405$.
 $\text{get_new_proposal_number}(\text{mon}, \text{oldpn}) \triangleq$
 $((\text{oldpn} \div 100) + 1) * 100 + \text{rank}(\text{mon})$

Clear the variable *peer_first_committed*.
 Variables changed: *peer_first_committed*.
 $clear_peer_first_committed(mon) \triangleq$
 $peer_first_committed' = [peer_first_committed \text{ EXCEPT } ![mon] =$
 $[m \in Monitors \mapsto -1]]$

Clear the variable *peer_last_committed*.
 Variables changed: *peer_last_committed*.
 $clear_peer_last_committed(mon) \triangleq$
 $peer_last_committed' = [peer_last_committed \text{ EXCEPT } ![mon] =$
 $[m \in Monitors \mapsto -1]]$

Store peer values and update *first_committed*, *last_committed* and *monitor_store* accordingly.
 Variables changed: *values*, *first_committed*, *last_committed*, *monitor_store*.
 $store_state(mon, msg) \triangleq$

Choose peer values from *mon* last committed + 1 to peer last committed.
 $\wedge \text{LET } logs \triangleq (\text{DOMAIN } msg.values) \cap (last_committed[mon] + 1 .. msg.last_committed)$
 $\text{IN } \wedge values' = [values \text{ EXCEPT } ![mon] =$
 $[i \in \text{DOMAIN } values[mon] \cup logs \mapsto$
 $\text{IF } i \notin \text{DOMAIN } values[mon]$
 $\text{THEN } msg.values[i]$
 $\text{ELSE } values[mon][i]]]$

Update last committed and first committed.
 $\wedge last_committed' = [last_committed \text{ EXCEPT } ![mon] = \text{Max}(logs \cup \{last_committed[mon]\})]$
 $\wedge \text{IF } logs \neq \{\} \wedge first_committed[mon] = 0$
 $\text{THEN } first_committed' =$
 $[first_committed \text{ EXCEPT } ![mon] = \text{Min}(logs)]$
 $\text{ELSE } first_committed' =$
 $[first_committed \text{ EXCEPT } ![mon] = \text{Min}(logs \cup \{first_committed[mon]\})]$

Update monitor store.
 $\wedge \text{IF } last_committed'[mon] = 0$
 $\text{THEN UNCHANGED } monitor_store$
 $\text{ELSE } monitor_store' = [monitor_store \text{ EXCEPT } ![mon] = values'[mon][last_committed'[mon]]]$

Check if uncommitted value version is still valid, else reset it.
 Variables changed: *uncommitted_v*, *uncommitted_value*.
 $check_and_correct_uncommitted(mon) \triangleq$
 $\text{IF } uncommitted_v[mon] \leq last_committed'[mon]$
 $\text{THEN } \wedge uncommitted_v' = [uncommitted_v \text{ EXCEPT } ![mon] = 0]$
 $\wedge uncommitted_value' = [uncommitted_value \text{ EXCEPT } ![mon] = Nil]$
 $\text{ELSE UNCHANGED } \langle uncommitted_v, uncommitted_value \rangle$

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*.

$finish_round(mon) \triangleq$
 $\wedge isLeader[mon] = \text{TRUE}$
 $\wedge state' = [state \text{ EXCEPT } ![mon] = \text{STATE_ACTIVE}]$

Resets the variable *acked_lease* and adds events to send lease messages to peers.

Variables changed: *acked_lease*, *send_queue*, *phase*.

$extend_lease(mon) \triangleq$
 $\wedge isLeader[mon] = \text{TRUE}$
 $\wedge acked_lease' = [acked_lease \text{ EXCEPT } ![mon] =$
 $\quad [m \in Monitors \mapsto \text{IF } m = mon \text{ THEN TRUE ELSE FALSE}]]$
 $\wedge send_queue' = [send_queue \text{ EXCEPT } ![mon] =$
 $\quad send_queue[mon] \circ SetToSeq((Monitors \setminus \{mon\}) \times \{OP_LEASE\})]$
 $\wedge phase' = [phase \text{ EXCEPT } ![mon] = \text{PHASE_LEASE}]$

Send a lease message from the leader to a peer.

Variables changed: *messages*.

$send_extend_lease(mon, dest) \triangleq$
 $\wedge isLeader[mon] = \text{TRUE}$
 $\wedge phase[mon] = \text{PHASE_LEASE}$
 $\wedge Send([type \mapsto OP_LEASE,$
 $\quad from \mapsto mon,$
 $\quad dest \mapsto dest,$
 $\quad last_committed \mapsto last_committed[mon]])$
 $\wedge \text{UNCHANGED } \langle epoch \rangle$
 $\wedge \text{UNCHANGED } \langle restart_vars, data_vars, state_vars, collect_vars, lease_vars, commit_vars \rangle$

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*, *state*.

$handle_lease(mon, msg) \triangleq$
 $\wedge \text{discard if not peon or peon is behind}$
 $\text{IF } \vee isLeader[mon] = \text{TRUE}$
 $\quad \vee last_committed[mon] \neq msg.last_committed$
 $\text{THEN } \wedge Discard(msg)$
 $\quad \wedge \text{UNCHANGED } state$
 $\text{ELSE } \wedge state' = [state \text{ EXCEPT } ![mon] = \text{STATE_ACTIVE}]$
 $\quad \wedge Reply([type \mapsto OP_LEASE_ACK,$
 $\quad from \mapsto mon,$

$$\begin{aligned}
& \text{dest} \mapsto \text{msg.from}, \\
& \text{first_committed} \mapsto \text{first_committed}[\text{mon}], \\
& \text{last_committed} \mapsto \text{last_committed}[\text{mon}], \text{msg}) \\
& \wedge \text{Discard}(\text{msg}) \\
& \wedge \text{UNCHANGED } \langle \text{epoch}, \text{isLeader}, \text{phase} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{restart_vars}, \text{data_vars}, \text{collect_vars}, \text{lease_vars}, \text{commit_vars}, \text{auxiliary_vars} \rangle
\end{aligned}$$

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

Once 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*.

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

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

Once 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*.

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

Commit phase predicates

Start a commit phase by the leader. The variable *new_value* is assigned and the events to send begin messages to the peers are added to *send_queue*.

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

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

$$\begin{aligned}
\text{begin}(\text{mon}, v) & \triangleq \\
& \wedge \text{isLeader}[\text{mon}] = \text{TRUE} \\
& \wedge \vee \text{state}'[\text{mon}] = \text{STATE_UPDATING} \\
& \quad \vee \text{state}'[\text{mon}] = \text{STATE_UPDATING_PREVIOUS} \\
& \wedge \text{Len}(\text{ranks}) = 1 \vee \text{num_last}[\text{mon}] > \text{Len}(\text{ranks}) \div 2 \\
& \wedge \text{new_value}[\text{mon}] = \text{Nil}
\end{aligned}$$

$$\begin{aligned}
&\wedge \text{accepted}' = [\text{accepted} \text{ EXCEPT } ![mon] = \\
&\quad [m \in \text{Monitors} \mapsto \text{IF } m = mon \text{ THEN TRUE ELSE FALSE}]] \\
&\wedge \text{new_value}' = [\text{new_value} \text{ EXCEPT } ![mon] = v] \\
&\wedge \text{phase}' = [\text{phase} \text{ EXCEPT } ![mon] = \text{PHASE_BEGIN}] \\
&\wedge \text{send_queue}' = [\text{send_queue} \text{ EXCEPT } ![mon] = \\
&\quad \text{send_queue}[mon] \circ \text{SetToSeq}((\text{Monitors} \setminus \{mon\}) \times \{OP_BEGIN\})] \\
&\wedge \text{values}' = [\text{values} \text{ EXCEPT } ![mon] = \\
&\quad (\text{values}[mon] @@ ((\text{last_committed}[mon] + 1) :> \text{new_value}'[mon]))] \\
&\wedge \text{uncommitted_v}' = [\text{uncommitted_v} \text{ EXCEPT } ![mon] = \text{last_committed}[mon] + 1] \\
&\wedge \text{uncommitted_value}' = [\text{uncommitted_value} \text{ EXCEPT } ![mon] = v]
\end{aligned}$$

Sends a begin message from the leader to the peer.

Variables changed: messages.

$$\begin{aligned}
&\text{send_begin}(mon, dest) \triangleq \\
&\quad \wedge \text{isLeader}[mon] = \text{TRUE} \\
&\quad \wedge \text{phase}[mon] = \text{PHASE_BEGIN} \\
&\quad \wedge \text{Send}([type \mapsto OP_BEGIN, \\
&\quad \quad \quad from \mapsto mon, \\
&\quad \quad \quad dest \mapsto dest, \\
&\quad \quad \quad \text{last_committed} \mapsto \text{last_committed}[mon], \\
&\quad \quad \quad \text{values} \mapsto \text{values}[mon], \\
&\quad \quad \quad pn \mapsto \text{accepted_pn}[mon]]) \\
&\quad \wedge \text{UNCHANGED } \langle epoch \rangle \\
&\quad \wedge \text{UNCHANGED } \langle \text{state_vars}, \text{restart_vars}, \text{data_vars}, \text{collect_vars}, \text{lease_vars}, \text{commit_vars} \rangle
\end{aligned}$$

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/restart.

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

$$\begin{aligned}
&\text{handle_begin}(mon, msg) \triangleq \\
&\quad \wedge \text{isLeader}[mon] = \text{FALSE} \\
&\quad \wedge \text{IF } msg.pn < \text{accepted_pn}[mon] \\
&\quad \quad \text{THEN} \\
&\quad \quad \quad \wedge \text{Discard}(msg) \\
&\quad \quad \quad \wedge \text{UNCHANGED } \langle \text{state}, \text{restart_vars} \rangle \\
&\quad \quad \text{ELSE} \\
&\quad \quad \quad \wedge msg.pn = \text{accepted_pn}[mon] \\
&\quad \quad \quad \wedge msg.\text{last_committed} = \text{last_committed}[mon] \\
&\quad \quad \quad \text{assign } \text{values}[mon][\text{last_committed}[mon] + 1] \\
&\quad \quad \wedge \text{values}' = [\text{values} \text{ EXCEPT } ![mon] = \\
&\quad \quad \quad (\text{values}[mon] @@ ((\text{last_committed}[mon] + 1) :> msg.\text{values}[\text{last_committed}[mon] + 1]))] \\
&\quad \wedge \text{state}' = [\text{state} \text{ EXCEPT } ![mon] = \text{STATE_UPDATING}]
\end{aligned}$$

$$\begin{aligned} &\wedge \text{uncommitted_}v' = [\text{uncommitted_}v \text{ EXCEPT } ![mon] = \text{last_committed}[mon] + 1] \\ &\wedge \text{uncommitted_value}' = [\text{uncommitted_value} \text{ EXCEPT } ![mon] = \\ &\quad \text{values}'[mon][\text{last_committed}[mon] + 1]] \end{aligned}$$

$$\begin{aligned} \wedge \text{Reply}([type &\mapsto OP_ACCEPT, \\ &from \mapsto mon, \\ &dest \mapsto msg.from, \\ &last_committed \mapsto \text{last_committed}[mon], \\ &pn \mapsto \text{accepted_pn}[mon]], msg) \end{aligned}$$

$$\wedge \text{UNCHANGED } \langle epoch, isLeader, phase, monitor_store, \text{accepted_pn}, \text{first_committed}, \text{last_committed} \rangle$$

$$\wedge \text{UNCHANGED } \langle collect_vars, lease_vars, commit_vars, auxiliary_vars \rangle$$

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`, `accepted`

$\text{handle_accept}(mon, msg) \triangleq$

$$\begin{aligned} &\wedge isLeader[mon] = \text{TRUE} \\ &\wedge \vee state[mon] = STATE_UPDATING_PREVIOUS \\ &\quad \vee state[mon] = STATE_UPDATING \\ &\wedge phase[mon] = PHASE_BEGIN \\ &\wedge new_value[mon] \neq Nil \\ &\wedge \text{IF } \vee msg.pn \neq \text{accepted_pn}[mon] \\ &\quad \vee \wedge \text{last_committed}[mon] > 0 \\ &\quad \wedge msg.\text{last_committed} < \text{last_committed}[mon] - 1 \end{aligned}$$

THEN

$$\wedge \text{Discard}(msg)$$

$$\wedge \text{UNCHANGED } \text{accepted}$$

ELSE

$$\wedge \text{accepted}' = [\text{accepted} \text{ EXCEPT } ![mon] = \\ \quad [\text{accepted}[mon] \text{ EXCEPT } ![msg.from] = \text{TRUE}]]$$

$$\wedge \text{Discard}(msg)$$

$$\wedge \text{UNCHANGED } \langle epoch, pending_proposal, new_value \rangle$$

$$\wedge \text{UNCHANGED } \langle restart_vars, state_vars, data_vars, collect_vars, lease_vars, auxiliary_vars \rangle$$

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

The leader commits the transaction in `new_value` and adds events in `send_queue` to send commit messages to his peers.

Variables changed: `first_committed`, `last_committed`, `monitor_store`, `new_value`, `send_queue`, `state`, `phase`

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

$$\begin{aligned} &\wedge phase[mon] = PHASE_BEGIN \\ &\wedge \forall m \in \text{Monitors} : \text{accepted}[mon][m] = \text{TRUE} \\ &\wedge new_value[mon] \neq Nil \\ &\wedge \vee state[mon] = STATE_UPDATING_PREVIOUS \\ &\quad \vee state[mon] = STATE_UPDATING \end{aligned}$$

$$\begin{aligned}
& \wedge \text{last_committed}' = [\text{last_committed} \text{ EXCEPT } ![mon] = \text{last_committed}[mon] + 1] \\
& \wedge \text{IF } \text{first_committed}[mon] = 0 \\
& \quad \text{THEN } \text{first_committed}' = [\text{first_committed} \text{ EXCEPT } ![mon] = \text{first_committed}[mon] + 1] \\
& \quad \text{ELSE UNCHANGED } \text{first_committed} \\
& \wedge \text{monitor_store}' = [\text{monitor_store} \text{ EXCEPT } ![mon] = \text{values}[mon][\text{last_committed}[mon] + 1]] \\
& \wedge \text{new_value}' = [\text{new_value} \text{ EXCEPT } ![mon] = \text{Nil}] \\
& \wedge \text{send_queue}' = [\text{send_queue} \text{ EXCEPT } ![mon] = \\
& \quad \text{send_queue}[mon] \circ \text{SetToSeq}((\text{Monitors} \setminus \{mon\}) \times \{OP_COMMIT\})] \\
& \wedge \text{state}' = [\text{state} \text{ EXCEPT } ![mon] = \text{STATE_REFRESH}] \\
& \wedge \text{phase}' = [\text{phase} \text{ EXCEPT } ![mon] = \text{PHASE_COMMIT}] \\
& \wedge \text{UNCHANGED } \langle \text{isLeader}, \text{values}, \text{accepted_pn}, \text{pending_proposal}, \text{accepted} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{global_vars}, \text{restart_vars}, \text{collect_vars}, \text{lease_vars} \rangle
\end{aligned}$$

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*, *send_queue*.

$$\begin{aligned}
& \text{finish_commit}(mon) \triangleq \\
& \quad \wedge \text{state}[mon] = \text{STATE_REFRESH} \\
& \quad \wedge \text{phase}[mon] = \text{PHASE_COMMIT} \\
& \quad \wedge \text{finish_round}(mon) \\
& \quad \wedge \text{extend_lease}(mon) \\
& \quad \wedge \text{UNCHANGED } \langle \text{isLeader} \rangle \\
& \quad \wedge \text{UNCHANGED } \langle \text{global_vars}, \text{restart_vars}, \text{data_vars}, \text{collect_vars}, \text{commit_vars} \rangle
\end{aligned}$$

Send a commit message from the leader to a peer.

Variables changed: *messages*.

$$\begin{aligned}
& \text{send_commit}(mon, dest) \triangleq \\
& \quad \wedge \text{isLeader}[mon] = \text{TRUE} \\
& \quad \wedge \text{Send}([type \mapsto OP_COMMIT, \\
& \quad \quad from \mapsto mon, \\
& \quad \quad dest \mapsto dest, \\
& \quad \quad last_committed \mapsto \text{last_committed}[mon], \\
& \quad \quad pn \mapsto \text{accepted_pn}[mon], \\
& \quad \quad values \mapsto \text{values}[mon]]) \\
& \quad \wedge \text{UNCHANGED } \langle \text{epoch} \rangle \\
& \quad \wedge \text{UNCHANGED } \langle \text{state_vars}, \text{restart_vars}, \text{data_vars}, \text{collect_vars}, \text{lease_vars}, \text{commit_vars} \rangle
\end{aligned}$$

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

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

$$\begin{aligned}
& \text{handle_commit}(mon, msg) \triangleq \\
& \quad \wedge \text{isLeader}[mon] = \text{FALSE} \\
& \quad \wedge \text{store_state}(mon, msg) \\
& \quad \wedge \text{check_and_correct_uncommitted}(mon) \\
& \quad \wedge \text{Discard}(msg)
\end{aligned}$$

$\wedge \text{UNCHANGED } \langle \text{epoch}, \text{accepted_pn} \rangle$
 $\wedge \text{UNCHANGED } \langle \text{state_vars}, \text{collect_vars}, \text{lease_vars}, \text{commit_vars}, \text{auxiliary_vars} \rangle$

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.

This predicate has a big cost on performance, so there were some requirements added (monitor phase and state) to mitigate that.

Variables changed: *pending_proposal*.

$\text{client_request}(\text{mon}, v) \triangleq$
 $\wedge \text{phase}[\text{mon}] = \text{PHASE_LEASE} \vee \text{phase}[\text{mon}] = \text{PHASE_ELECTION}$
 $\wedge \text{isLeader}[\text{mon}] = \text{TRUE}$
 $\wedge \text{state}[\text{mon}] = \text{STATE_ACTIVE}$
 $\wedge \text{pending_proposal}[\text{mon}] = \text{Nil}$
 $\wedge \text{pending_proposal}' = [\text{pending_proposal} \text{ EXCEPT } ![\text{mon}] = v]$
 $\wedge \text{UNCHANGED } \langle \text{new_value}, \text{accepted} \rangle$
 $\wedge \text{UNCHANGED } \langle \text{global_vars}, \text{state_vars}, \text{restart_vars}, \text{data_vars}, \text{collect_vars}, \text{lease_vars}, \text{auxiliary_vars} \rangle$

Start a commit phase with the value on pending proposal.

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

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

Collect phase predicates

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

Variables changed: *accepted_pn*, *phase*.

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

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.

Add events to *send_queue* to send collect messages to the peers.

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

$$\begin{aligned}
pre_send_collect(mon) &\triangleq \\
&\wedge state[mon] = STATE_RECOVERING \\
&\wedge isLeader[mon] = TRUE \\
&\wedge phase[mon] = PHASE_PRE_COLLECT \\
&\wedge clear_peer_first_committed(mon) \\
&\wedge clear_peer_last_committed(mon) \\
&\wedge IF\ last_committed[mon] + 1 \in DOMAIN\ values[mon] \\
&\quad THEN\ \wedge\ uncommitted_v' = \\
&\quad\quad [uncommitted_v\ EXCEPT\ ![mon] = last_committed[mon] + 1] \\
&\quad\quad \wedge\ uncommitted_value' = \\
&\quad\quad\quad [uncommitted_value\ EXCEPT\ ![mon] = values[mon][last_committed[mon] + 1]] \\
&\quad ELSE\ UNCHANGED\ \langle restart_vars \rangle \\
&\wedge num_last' = [num_last\ EXCEPT\ ![mon] = 1] \\
&\wedge send_queue' = [send_queue\ EXCEPT\ ![mon] = \\
&\quad\quad send_queue[mon] \circ SetToSeq((Monitors \setminus \{mon\}) \times \{OP_COLLECT\})] \\
&\wedge phase' = [phase\ EXCEPT\ ![mon] = PHASE_COLLECT] \\
&\wedge UNCHANGED\ \langle isLeader, state \rangle \\
&\wedge UNCHANGED\ \langle global_vars, data_vars, lease_vars, commit_vars \rangle
\end{aligned}$$

Send a collect message from the leader to a peer.

Variables changed: *messages*.

$$\begin{aligned}
send_collect(mon, dest) &\triangleq \\
&\wedge state[mon] = STATE_RECOVERING \\
&\wedge isLeader[mon] = TRUE \\
&\wedge phase[mon] = PHASE_COLLECT \\
&\wedge Send([type \mapsto OP_COLLECT, \\
&\quad from \mapsto mon, \\
&\quad dest \mapsto dest, \\
&\quad first_committed \mapsto first_committed[mon], \\
&\quad last_committed \mapsto last_committed[mon], \\
&\quad pn \mapsto accepted_pn[mon]]) \\
&\wedge UNCHANGED\ epoch \\
&\wedge UNCHANGED\ \langle state_vars, restart_vars, data_vars, collect_vars, lease_vars, commit_vars \rangle
\end{aligned}$$

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*, *epoch*, *state*, *accepted_pn*

$$\begin{aligned}
handle_collect(mon, msg) &\triangleq \\
&\wedge isLeader[mon] = FALSE \\
&\wedge state' = [state\ EXCEPT\ ![mon] = STATE_RECOVERING]
\end{aligned}$$

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

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 with that.

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

$handle_last(mon, msg) \triangleq$

$$\begin{aligned}
& \wedge isLeader[mon] = \text{TRUE} \\
& \wedge peer_first_committed' = [peer_first_committed \text{ EXCEPT } ![mon] = \\
& \quad [peer_first_committed[mon] \text{ EXCEPT } ![msg.from] = msg.first_committed]] \\
& \wedge peer_last_committed' = [peer_last_committed \text{ EXCEPT } ![mon] = \\
& \quad [peer_last_committed[mon] \text{ EXCEPT } ![msg.from] = msg.last_committed]] \\
& \wedge IF msg.first_committed > last_committed[mon] + 1 \\
& \quad THEN \\
& \quad \quad \wedge bootstrap \\
& \quad \quad \wedge UNCHANGED \langle num_last, accepted_pn, values, phase, monitor_store \rangle \\
& \quad \quad \wedge UNCHANGED \langle first_committed, last_committed, restart_vars, auxiliary_vars \rangle \\
& \quad ELSE \\
& \quad \quad \wedge store_state(mon, msg) \\
& \quad \quad \wedge IF \exists peer \in Monitors : \\
& \quad \quad \quad \wedge peer \neq mon \\
& \quad \quad \quad \wedge peer_last_committed'[mon][peer] \neq -1 \\
& \quad \quad \quad \wedge peer_last_committed'[mon][peer] + 1 < first_committed[mon] \\
& \quad \quad \quad \wedge first_committed[mon] > 1 \\
& \quad THEN
\end{aligned}$$

```

 $\wedge$  bootstrap
 $\wedge$  check_and_correct_uncommitted(mon)
 $\wedge$  UNCHANGED  $\langle$ phase, accepted_pn, num_last, auxiliary_vars $\rangle$ 
ELSE
 $\wedge$  LET monitors_behind  $\triangleq$  {peer  $\in$  Monitors :
     $\wedge$  peer  $\neq$  mon
     $\wedge$  peer_last_committed'[mon][peer]  $\neq$  -1
     $\wedge$  peer_last_committed'[mon][peer] < last_committed[mon]}
IN send_queue' = [send_queue EXCEPT ![mon] =
    send_queue[mon]  $\circ$  SetToSeq((monitors_behind)  $\times$  {OP_COMMIT})]

 $\wedge$   $\vee$   $\wedge$  msg.pn > accepted_pn[mon]
     $\wedge$  collect(mon, msg.pn)
     $\wedge$  check_and_correct_uncommitted(mon)
     $\wedge$  UNCHANGED num_last

 $\vee$   $\wedge$  msg.pn = accepted_pn[mon]
     $\wedge$  num_last' = [num_last EXCEPT ![mon] = num_last[mon] + 1]

     $\wedge$  IF  $\wedge$  msg.last_committed + 1  $\in$  DOMAIN msg.values
         $\wedge$  msg.last_committed  $\geq$  last_committed'[mon]
         $\wedge$  msg.last_committed + 1  $\geq$  uncommitted_v[mon]
    THEN  $\wedge$  uncommitted_v' =
        [uncommitted_v EXCEPT ![mon] = msg.last_committed + 1]
         $\wedge$  uncommitted_value' =
        [uncommitted_value EXCEPT ![mon] = msg.values[msg.last_committed + 1]]
    ELSE check_and_correct_uncommitted(mon)

     $\wedge$  UNCHANGED  $\langle$ phase, accepted_pn $\rangle$ 

 $\vee$   $\wedge$  msg.pn < accepted_pn[mon]
     $\wedge$  check_and_correct_uncommitted(mon)
     $\wedge$  UNCHANGED  $\langle$ phase, accepted_pn, num_last $\rangle$ 
 $\wedge$  UNCHANGED epoch
 $\wedge$  UNCHANGED  $\langle$ epoch $\rangle$ 

 $\wedge$  Discard(msg)
 $\wedge$  UNCHANGED  $\langle$ isLeader, state $\rangle$ 
 $\wedge$  UNCHANGED  $\langle$ lease_vars, commit_vars $\rangle$ 

```

Predicate that is enabled and called when all peers 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*, *send_queue*, *values*, *uncommitted_v*, *uncommitted_value*, *acked_lease*.

```

post_last(mon)  $\triangleq$ 
 $\wedge$  isLeader[mon] = TRUE
 $\wedge$  num_last[mon] = Len(ranks)

```


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

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

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

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

Variables changed: *accepted_pn*, *phase*.

election_recover(mon) \triangleq

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

Timeouts and restart

Restart a monitor and wipe variables that are not persistent.

Variables changed: *messages*, *isLeader*, *phase*, *state*, *pending_proposal*, *new_value*, *number_refreshes*.

$$\begin{aligned}
\text{restart_mon}(\text{mon}) &\triangleq \\
&\wedge \text{messages}' = \text{SelectSeq}(\text{messages}, \text{LAMBDA } t : t.\text{from} \neq \text{mon}) \\
&\wedge \text{isLeader}' = [\text{isLeader} \text{ EXCEPT } ![\text{mon}] = \text{FALSE}] \\
&\wedge \text{phase}' = [\text{phase} \text{ EXCEPT } ![\text{mon}] = \text{PHASE_ELECTION}] \\
&\wedge \text{state}' = [\text{state} \text{ EXCEPT } ![\text{mon}] = \text{IF } \text{Len}(\text{ranks}) = 1 \\
&\quad \text{THEN } \text{STATE_ACTIVE} \\
&\quad \text{ELSE } \text{STATE_RECOVERING}] \\
&\wedge \text{pending_proposal}' = [\text{pending_proposal} \text{ EXCEPT } ![\text{mon}] = \text{Nil}] \\
&\wedge \text{new_value}' = [\text{new_value} \text{ EXCEPT } ![\text{mon}] = \text{Nil}] \\
&\wedge \text{number_refreshes}' = \text{number_refreshes} + 1 \\
&\wedge \text{UNCHANGED } \langle \text{epoch}, \text{message_history}, \text{accepted} \rangle \\
&\wedge \text{UNCHANGED } \langle \text{restart_vars}, \text{data_vars}, \text{collect_vars}, \text{lease_vars}, \text{auxiliary_vars} \rangle
\end{aligned}$$

Monitor timeout (simulate message not received). Triggers new elections.

Messages in network and events in *send_queue* are cleared.

Variables changed: *epoch*, *send_queue*, *messages*.

$$\begin{aligned}
\text{Timeout}(\text{mon}) &\triangleq \\
&\wedge \text{phase}[\text{mon}] = \text{PHASE_COLLECT} \vee \text{phase}[\text{mon}] = \text{PHASE_BEGIN} \\
&\wedge \text{bootstrap} \\
&\wedge \text{send_queue}' = [m \in \text{Monitors} \mapsto \langle \rangle] \\
&\wedge \text{messages}' = \langle \rangle \\
&\wedge \text{UNCHANGED } \langle \text{message_history}, \text{state_vars}, \text{restart_vars}, \text{data_vars}, \text{collect_vars}, \text{lease_vars}, \text{commit_vars} \rangle
\end{aligned}$$

Dispatchers and next statement

Handle a message.

$$\begin{aligned}
\text{Receive}(\text{msg}) &\triangleq \\
&\wedge \vee \text{phase}[\text{msg.dest}] = \text{PHASE_COLLECT} \\
&\quad \vee \text{phase}[\text{msg.dest}] = \text{PHASE_BEGIN} \\
&\quad \vee \text{phase}[\text{msg.dest}] = \text{PHASE_ELECTION} \\
&\wedge \\
&\quad \vee \wedge \text{msg.type} = \text{OP_COLLECT} \\
&\quad \quad \wedge \text{handle_collect}(\text{msg.dest}, \text{msg}) \\
&\quad \quad \wedge \text{step_x}' = \text{"receive collect"} \\
&\quad \vee \wedge \text{msg.type} = \text{OP_LAST} \\
&\quad \quad \wedge \text{handle_last}(\text{msg.dest}, \text{msg}) \\
&\quad \quad \wedge \text{step_x}' = \text{"receive last"} \\
&\quad \vee \wedge \text{msg.type} = \text{OP_LEASE} \\
&\quad \quad \wedge \text{handle_lease}(\text{msg.dest}, \text{msg}) \\
&\quad \quad \wedge \text{step_x}' = \text{"receive lease"}
\end{aligned}$$

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

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

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

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

Send a message from the event queue.

$Send_from_queue(mon) \triangleq$
 $\wedge Len(send_queue[mon]) > 0$
 \wedge
 $\vee \wedge Head(send_queue[mon])[2] = OP_COLLECT$
 $\wedge send_collect(mon, Head(send_queue[mon])[1])$
 $\wedge step_x' = \text{"send_collect"}$

$\vee \wedge Head(send_queue[mon])[2] = OP_LEASE$
 $\wedge send_extend_lease(mon, Head(send_queue[mon])[1])$
 $\wedge step_x' = \text{"send_extend_lease"}$

$\vee \wedge Head(send_queue[mon])[2] = OP_BEGIN$
 $\wedge send_begin(mon, Head(send_queue[mon])[1])$
 $\wedge step_x' = \text{"send_begin"}$

$\vee \wedge Head(send_queue[mon])[2] = OP_COMMIT$
 $\wedge send_commit(mon, Head(send_queue[mon])[1])$
 $\wedge step_x' = \text{"send_commit"}$

$\wedge send_queue' = [send_queue \text{ EXCEPT } ![mon] = Tail(send_queue[mon])]$

Limit some variables to reduce search space.

$reduce_search_space \triangleq$
 $\wedge epoch < 5$
 $\wedge \forall mon \in Monitors : accepted_pn[mon] < 300$
 $\wedge \forall mon \in Monitors : last_committed[mon] < 2$
 $\wedge number_refreshes < 3$
 $\wedge step < 100$

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_refreshes
ELSE
  IF  $\exists mon \in Monitors : Len(send\_queue[mon]) > 0$ 
    THEN  $\wedge \exists mon \in Monitors : Send\_from\_queue(mon)$ 
       $\wedge$   $step' = step + 1$ 
       $\wedge$  UNCHANGED number_refreshes
    ELSE
       $\vee$ 
       $\wedge \exists mon \in Monitors : election\_recover(mon)$ 
       $\wedge$   $step\_x' = \text{"election\_recover"} \wedge step' = step + 1$ 
       $\wedge$  UNCHANGED number_refreshes

       $\vee$   $\wedge \exists mon \in Monitors : pre\_send\_collect(mon)$ 
       $\wedge$   $step\_x' = \text{"pre\_send\_collect"} \wedge step' = step + 1$ 
       $\wedge$  UNCHANGED number_refreshes

       $\vee$   $\wedge \exists mon \in Monitors : post\_last(mon)$ 
       $\wedge$   $step\_x' = \text{"post\_last"} \wedge step' = step + 1$ 
       $\wedge$  UNCHANGED number_refreshes

       $\vee$   $\wedge \exists mon \in Monitors : post\_lease\_ack(mon)$ 
       $\wedge$   $step\_x' = \text{"post\_lease\_ack"} \wedge step' = step + 1$ 
       $\wedge$  UNCHANGED number_refreshes

       $\vee$   $\wedge \exists mon \in Monitors : post\_accept(mon)$ 
       $\wedge$   $step\_x' = \text{"post\_accept"} \wedge step' = step + 1$ 
       $\wedge$  UNCHANGED number_refreshes

       $\vee$   $\wedge \exists mon \in Monitors : finish\_commit(mon)$ 
       $\wedge$   $step\_x' = \text{"finish\_commit"} \wedge step' = step + 1$ 
       $\wedge$  UNCHANGED number_refreshes

       $\vee$   $\wedge \exists mon \in Monitors : \exists v \in Value\_set : client\_request(mon, v)$ 
       $\wedge$   $step\_x' = \text{"client\_request"} \wedge step' = step + 1$ 
       $\wedge$  UNCHANGED number_refreshes

       $\vee$   $\wedge \exists mon \in Monitors : propose\_pending(mon)$ 
       $\wedge$   $step\_x' = \text{"propose\_pending"} \wedge step' = step + 1$ 
       $\wedge$  UNCHANGED number_refreshes

       $\vee$   $\wedge \exists i \in 1 \dots Len(messages) : Receive(messages[i])$ 
       $\wedge$   $step' = step + 1$ 
       $\wedge$  UNCHANGED number_refreshes

       $\vee$   $\wedge \exists mon \in Monitors : restart\_mon(mon)$ 

```

$$\begin{aligned}
& \wedge \text{step_}x' = \text{"restart mon"} \wedge \text{step}' = \text{step} + 1 \\
& \vee \wedge \exists \text{mon} \in \text{Monitors} : \text{Timeout}(\text{mon}) \\
& \quad \wedge \text{step_}x' = \text{"timeout and restart"} \wedge \text{step}' = \text{step} + 1 \\
& \quad \wedge \text{UNCHANGED } \text{number_refreshes}
\end{aligned}$$

Test/Debug invariants

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

$$\text{Inv_find_state}(x) \triangleq \neg x$$

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

$$\text{Inv_diam}(\text{size}) \triangleq \text{step} \neq \text{size} - 1$$

Invariants to test in model check

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

Examples:

Find a behavior with a diameter of size 60.

$$\text{Inv_diam}(60)$$

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

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

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

$$\begin{aligned}
& \text{Inv_find_state}(\\
& \quad \wedge \text{step_}x = \text{"restart mon"} \\
& \quad \backslash * \text{The system is in collect phase and no } \text{OP_LAST} \text{ message has been received.} \\
& \quad \backslash * \text{isLeader}[\text{mon}] = \text{TRUE} \text{ assures that the leader was not the one that crashed.} \\
& \quad \wedge \exists \text{mon} \in \text{Monitors} : \\
& \quad \quad \wedge \text{isLeader}[\text{mon}] = \text{TRUE} \\
& \quad \quad \wedge \text{phase}[\text{mon}] = \text{PHASE_COLLECT} \\
& \quad \quad \wedge \text{num_last}[\text{mon}] = 1 \\
& \quad \backslash * \text{All the collect requests have been handled by the peers.} \\
& \quad \wedge \forall i \in 1 \dots \text{Len}(\text{messages}) : \\
& \quad \quad \text{messages}[i].\text{type} \neq \text{OP_COLLECT} \\
& \quad \wedge \text{epoch} = 2 \\
&)
\end{aligned}$$

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

$$\begin{aligned}
& \text{Inv_find_state}(\\
& \quad \wedge \text{step_}x = \text{"restart mon"}
\end{aligned}$$

```

 $\wedge \exists i \in 1 \dots Len(messages) :$ 
     $messages[i].type = OP\_ACCEPT$ 
 $\wedge \forall mon \in Monitors :$ 
     $isLeader[mon] = FALSE$ 
 $\wedge epoch = 2$ 
)

```

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

```

\ * Modification History
\ * Last modified Fri Mar 05 16:07:45 WET 2021 by afonsonf
\ * Created Mon Jan 11 16:15:26 WET 2021 by afonsonf

```