

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 both the messages waiting to be handled and the ones already received. 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/afonsof/ceph-paxos-tla>

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

$\text{state_names} \triangleq \{\text{STATE_RECOVERING}, \text{STATE_ACTIVE},$
 $\text{STATE_UPDATING}, \text{STATE_UPDATING_PREVIOUS},$
 $\text{STATE_WRITING}, \text{STATE_WRITING_PREVIOUS},$
 $\text{STATE_REFRESH}, \text{STATE_SHUTDOWN}\}$

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: [message 'm' \mapsto 1 if 'm' is in the network else 0]

VARIABLE *messages*

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

$$\begin{aligned}
 global_vars &\triangleq \langle epoch, messages \rangle \\
 state_vars &\triangleq \langle isLeader, state, phase \rangle \\
 restart_vars &\triangleq \langle uncommitted_v, uncommitted_value \rangle \\
 data_vars &\triangleq \langle monitor_store, values, accepted_pn, first_committed, last_committed \rangle \\
 collect_vars &\triangleq \langle num_last, peer_first_committed, peer_last_committed \rangle \\
 lease_vars &\triangleq \langle acked_lease \rangle \\
 commit_vars &\triangleq \langle pending_proposal, new_value, accepted \rangle \\
 auxiliary_vars &\triangleq \langle send_queue \rangle \\
 vars &\triangleq \langle global_vars, state_vars, restart_vars, data_vars, collect_vars, \\
 &\quad lease_vars, commit_vars, auxiliary_vars \rangle
 \end{aligned}$$

$$\begin{aligned}
& \textit{Init_global_vars} \triangleq \\
& \quad \wedge \textit{epoch} = 1 \\
& \quad \wedge \textit{messages} = [m \in \{\} \mapsto 0] \\
& \textit{Init_state_vars} \triangleq \\
& \quad \wedge \textit{isLeader} = [mon \in \textit{Monitors} \mapsto \text{FALSE}] \\
& \quad \wedge \textit{state} = [mon \in \textit{Monitors} \mapsto \text{Nil}] \\
& \quad \wedge \textit{phase} = [mon \in \textit{Monitors} \mapsto \text{Nil}] \\
& \textit{Init_restart_vars} \triangleq \\
& \quad \wedge \textit{uncommitted_v} = [mon \in \textit{Monitors} \mapsto 0] \\
& \quad \wedge \textit{uncommitted_value} = [mon \in \textit{Monitors} \mapsto \text{Nil}] \\
& \textit{Init_data_vars} \triangleq \\
& \quad \wedge \textit{monitor_store} = [mon \in \textit{Monitors} \mapsto \text{Nil}] \\
& \quad \wedge \textit{values} = [mon \in \textit{Monitors} \mapsto [version \in \{\} \mapsto \text{Nil}]] \\
& \quad \wedge \textit{accepted_pn} = [mon \in \textit{Monitors} \mapsto 0] \\
& \quad \wedge \textit{first_committed} = [mon \in \textit{Monitors} \mapsto 0] \\
& \quad \wedge \textit{last_committed} = [mon \in \textit{Monitors} \mapsto 0] \\
& \textit{Init_collect_vars} \triangleq \\
& \quad \wedge \textit{num_last} = [mon \in \textit{Monitors} \mapsto 0] \\
& \quad \wedge \textit{peer_first_committed} = [mon1 \in \textit{Monitors} \mapsto [mon2 \in \textit{Monitors} \mapsto -1]] \\
& \quad \wedge \textit{peer_last_committed} = [mon1 \in \textit{Monitors} \mapsto [mon2 \in \textit{Monitors} \mapsto -1]] \\
& \textit{Init_lease_vars} \triangleq \\
& \quad \wedge \textit{acked_lease} = [mon1 \in \textit{Monitors} \mapsto [mon2 \in \textit{Monitors} \mapsto \text{FALSE}]] \\
& \textit{Init_commit_vars} \triangleq \\
& \quad \wedge \textit{pending_proposal} = [mon \in \textit{Monitors} \mapsto \text{Nil}] \\
& \quad \wedge \textit{new_value} = [mon \in \textit{Monitors} \mapsto \text{Nil}] \\
& \quad \wedge \textit{accepted} = [mon1 \in \textit{Monitors} \mapsto [mon2 \in \textit{Monitors} \mapsto \text{FALSE}]] \\
& \textit{Init_auxiliary_vars} \triangleq \\
& \quad \wedge \textit{send_queue} = [mon \in \textit{Monitors} \mapsto \langle \rangle] \\
& \textit{Init} \triangleq \\
& \quad \wedge \textit{Init_global_vars} \\
& \quad \wedge \textit{Init_state_vars} \\
& \quad \wedge \textit{Init_restart_vars} \\
& \quad \wedge \textit{Init_data_vars} \\
& \quad \wedge \textit{Init_collect_vars} \\
& \quad \wedge \textit{Init_lease_vars} \\
& \quad \wedge \textit{Init_auxiliary_vars} \\
& \quad \wedge \textit{Init_commit_vars} \\
& \quad \wedge \textit{step} = 0 \wedge \textit{step_x} = \text{"init"} \wedge \textit{number_refreshes} = 0
\end{aligned}$$

Message manipulation

Add message m to the network $msgs$.

$$\begin{aligned} WithMessage(m, msgs) &\triangleq \\ (m :> 1) @@ msgs \end{aligned}$$

Remove message m from the network $msgs$.

$$\begin{aligned} WithoutMessage(m, msgs) &\triangleq \\ (m :> 0) @@ msgs \end{aligned}$$

Set of messages in network.

$$\begin{aligned} ValidMessage(msgs) &\triangleq \\ \{m \in \text{DOMAIN } messages : msgs[m] = 1\} \end{aligned}$$

Adds the message m to the network.

Variables changed: $messages$.

$$\begin{aligned} Send(m) &\triangleq \\ messages' = WithMessage(m, messages) \end{aligned}$$

Removes message m from the network.

Variables changed: $messages$.

$$\begin{aligned} Discard(m) &\triangleq \\ messages' = WithoutMessage(m, messages) \end{aligned}$$

Removes the request from network and adds the response.

Variables changed: $messages$.

$$\begin{aligned} Reply(response, request) &\triangleq \\ messages' = WithoutMessage(request, WithMessage(response, messages)) \end{aligned}$$

Helper predicates

Compute a new unique proposal number for a given monitor.

Example: $oldpn = 305$, $rank(mon) = 5$, $newpn = 405$.

$$\begin{aligned} get_new_proposal_number(mon, oldpn) &\triangleq \\ ((oldpn \div 100) + 1) * 100 + rank(mon) \end{aligned}$$

Clear the variable $peer_first_committed$.

Variables changed: $peer_first_committed$.

$$\begin{aligned} clear_peer_first_committed(mon) &\triangleq \\ peer_first_committed' = [peer_first_committed \text{ EXCEPT } ![mon] = \\ [m \in Monitors \mapsto -1]] \end{aligned}$$

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*.

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

Sends a lease message from the leader to a peer.

Variables changed: *messages*.

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

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*.

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

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*.

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

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.

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

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.

$$\begin{aligned}
& \text{begin}(mon, v) \triangleq \\
& \quad \wedge isLeader[mon] = \text{TRUE} \\
& \quad \wedge \vee state'[mon] = \text{STATE_UPDATING} \\
& \quad \quad \vee state'[mon] = \text{STATE_UPDATING_PREVIOUS} \\
& \quad \wedge Len(ranks) = 1 \vee num_last[mon] > Len(ranks) \div 2 \\
& \quad \wedge new_value[mon] = Nil \\
& \quad \wedge accepted' = [\text{accepted} \text{ EXCEPT } ![mon] = \\
& \quad \quad [m \in \text{Monitors} \mapsto \text{IF } m = mon \text{ THEN TRUE ELSE FALSE}]] \\
& \quad \wedge new_value' = [new_value \text{ EXCEPT } ![mon] = v] \\
& \quad \wedge phase' = [phase \text{ EXCEPT } ![mon] = \text{LEASE_BEGIN}] \\
& \quad \wedge send_queue' = [send_queue \text{ EXCEPT } ![mon] = \\
& \quad \quad send_queue[mon] \circ \text{SetToSeq}((\text{Monitors} \setminus \{mon\}) \times \{OP_BEGIN\})] \\
& \quad \wedge values' = [values \text{ EXCEPT } ![mon] = \\
& \quad \quad (values[mon] @@ ((last_committed[mon] + 1) :> new_value'[mon]))] \\
& \quad \wedge uncommitted_v' = [uncommitted_v \text{ EXCEPT } ![mon] = last_committed[mon] + 1] \\
& \quad \wedge uncommitted_value' = [uncommitted_value \text{ EXCEPT } ![mon] = v]
\end{aligned}$$

Sends a begin message from the leader to the peer.

Variables changed: messages.

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

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*.

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

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

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*

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

$\wedge \text{isLeader}[\text{mon}] = \text{TRUE}$

$\wedge \vee \text{state}[\text{mon}] = \text{STATE_UPDATING_PREVIOUS}$

$\vee \text{state}[\text{mon}] = \text{STATE_UPDATING}$

$\wedge \text{phase}[\text{mon}] = \text{PHASE_BEGIN}$

$\wedge \text{new_value}[\text{mon}] \neq \text{Nil}$

$\wedge \text{IF } \vee \text{msg.pn} \neq \text{accepted_pn}[\text{mon}]$

$\vee \wedge \text{last_committed}[\text{mon}] > 0$

$\wedge \text{msg.last_committed} < \text{last_committed}[\text{mon}] - 1$

THEN

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

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

ELSE

$\wedge \text{accepted}' = [\text{accepted} \text{ EXCEPT } ![\text{mon}] =$

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

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

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

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

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*

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

$\wedge \text{phase}[\text{mon}] = \text{PHASE_BEGIN}$

$\wedge \forall m \in \text{Monitors} : \text{accepted}[\text{mon}][m] = \text{TRUE}$

$\wedge \text{new_value}[\text{mon}] \neq \text{Nil}$

$\wedge \vee \text{state}[\text{mon}] = \text{STATE_UPDATING_PREVIOUS}$

$\vee \text{state}[\text{mon}] = \text{STATE_UPDATING}$

$\wedge \text{last_committed}' = [\text{last_committed} \text{ EXCEPT } ![\text{mon}] = \text{last_committed}[\text{mon}] + 1]$

$\wedge \text{IF } \text{first_committed}[\text{mon}] = 0$

THEN $\text{first_committed}' = [\text{first_committed} \text{ EXCEPT } ![\text{mon}] = \text{first_committed}[\text{mon}] + 1]$

ELSE UNCHANGED *first_committed*

$\wedge \text{monitor_store}' = [\text{monitor_store} \text{ EXCEPT } ![\text{mon}] = \text{values}[\text{mon}][\text{last_committed}[\text{mon}] + 1]]$

$\wedge \text{new_value}' = [\text{new_value} \text{ EXCEPT } ![\text{mon}] = \text{Nil}]$

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

$\text{send_queue}[\text{mon}] \circ \text{SetToSeq}((\text{Monitors} \setminus \{\text{mon}\}) \times \{\text{OP_COMMIT}\})]$

$\wedge \text{state}' = [\text{state} \text{ EXCEPT } ![\text{mon}] = \text{STATE_REFRESH}]$

$\wedge \text{phase}' = [\text{phase} \text{ EXCEPT } ![\text{mon}] = \text{PHASE_COMMIT}]$

$\wedge \text{UNCHANGED } \langle isLeader, values, accepted_pn, pending_proposal, accepted \rangle$
 $\wedge \text{UNCHANGED } \langle global_vars, restart_vars, collect_vars, lease_vars \rangle$

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

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

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

Send a commit message from the leader to a peer.

Variables changed: *messages*.

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

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*.

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

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*.

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

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*.

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

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*.

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

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*.

$pre_send_collect(mon) \triangleq$
 $\wedge state[mon] = STATE_RECOVERING$
 $\wedge isLeader[mon] = TRUE$
 $\wedge phase[mon] = PHASE_PRE_COLLECT$
 $\wedge clear_peer_first_committed(mon)$
 $\wedge clear_peer_last_committed(mon)$

\wedge IF $last_committed[mon] + 1 \in \text{DOMAIN } values[mon]$
 THEN $\wedge uncommitted_v' =$
 $[uncommitted_v \text{ EXCEPT } ![mon] = last_committed[mon] + 1]$
 $\wedge uncommitted_value' =$
 $[uncommitted_value \text{ EXCEPT } ![mon] = values[mon][last_committed[mon] + 1]]$
 ELSE UNCHANGED $\langle restart_vars \rangle$

 $\wedge num_last' = [num_last \text{ EXCEPT } ![mon] = 1]$
 $\wedge send_queue' = [send_queue \text{ EXCEPT } ![mon] =$
 $send_queue[mon] \circ SetToSeq((Monitors \setminus \{mon\}) \times \{OP_COLLECT\})]$
 $\wedge phase' = [phase \text{ EXCEPT } ![mon] = PHASE_COLLECT]$
 \wedge UNCHANGED $\langle isLeader, state \rangle$
 \wedge UNCHANGED $\langle global_vars, data_vars, lease_vars, commit_vars \rangle$

Send a collect message from the leader to a peer.

Variables changed: messages.

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

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, epoch, state, accepted_pn

$handle_collect(mon, msg) \triangleq$
 $\wedge isLeader[mon] = \text{FALSE}$
 $\wedge state' = [state \text{ EXCEPT } ![mon] = STATE_RECOVERING]$
 $\wedge \vee \wedge msg.first_committed > last_committed[mon] + 1$
 $\wedge bootstrap$
 $\wedge Discard(msg)$
 \wedge UNCHANGED $\langle accepted_pn \rangle$
 $\vee \wedge msg.first_committed \leq last_committed[mon] + 1$
 \wedge IF $msg.pn > accepted_pn[mon]$
 THEN $accepted_pn' = [accepted_pn \text{ EXCEPT } ![mon] = msg.pn]$
 ELSE UNCHANGED $accepted_pn$
 $\wedge Reply([type \mapsto OP_LAST,$
 $from \mapsto mon,$
 $dest \mapsto msg.from,$

$$\begin{aligned}
& first_committed \mapsto first_committed[mon], \\
& last_committed \mapsto last_committed[mon], \\
& values \mapsto values[mon], \\
& pn \mapsto accepted_pn'[mon], msg) \\
& \wedge \text{UNCHANGED } epoch \\
& \wedge \text{UNCHANGED } \langle isLeader, phase, values, first_committed, last_committed, monitor_store \rangle \\
& \wedge \text{UNCHANGED } \langle restart_vars, collect_vars, lease_vars, commit_vars, auxiliary_vars \rangle
\end{aligned}$$

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*.

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

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

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*.

$$\begin{aligned}
post_last(mon) & \triangleq \\
& \wedge isLeader[mon] = \text{TRUE} \\
& \wedge num_last[mon] = Len(ranks) \\
& \wedge phase[mon] = PHASE_COLLECT \\
& \wedge clear_peer_first_committed(mon) \\
& \wedge clear_peer_last_committed(mon) \\
& \wedge \text{IF } \wedge uncommitted_v[mon] = last_committed[mon] + 1 \\
& \quad \wedge uncommitted_value[mon] \neq Nil \\
& \quad \text{THEN } \wedge state' = [state \text{ EXCEPT } ![mon] = STATE_UPDATING_PREVIOUS] \\
& \quad \quad \wedge begin(mon, uncommitted_value) \\
& \quad \quad \wedge UNCHANGED \langle acked_lease \rangle \\
& \quad \text{ELSE } \wedge finish_round(mon)
\end{aligned}$$

$\wedge \text{extend_lease}(\text{mon})$
 $\wedge \text{UNCHANGED } \langle \text{accepted}, \text{new_value}, \text{values}, \text{restart_vars} \rangle$
 $\wedge \text{UNCHANGED } \langle \text{isLeader}, \text{monitor_store}, \text{accepted_pn}, \text{first_committed}, \text{last_committed} \rangle$
 $\wedge \text{UNCHANGED } \langle \text{num_last}, \text{pending_proposal} \rangle$
 $\wedge \text{UNCHANGED } \langle \text{global_vars} \rangle$

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*.

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

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

Variables changed: *accepted_pn*, *phase*.

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

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*.

$\text{restart_mon}(\text{mon}) \triangleq$
 $\wedge \text{messages}' = [m \in \text{DOMAIN } \text{messages} \mapsto \text{IF } m.\text{from} = \text{mon} \text{ THEN } 0 \text{ ELSE } \text{messages}[m]]$
 $\wedge \text{isLeader}' = [\text{isLeader} \text{ EXCEPT } ![\text{mon}] = \text{FALSE}]$
 $\wedge \text{phase}' = [\text{phase} \text{ EXCEPT } ![\text{mon}] = \text{PHASE_ELECTION}]$
 $\wedge \text{state}' = [\text{state} \text{ EXCEPT } ![\text{mon}] = \text{IF } \text{Len}(\text{ranks}) = 1$
 $\quad \text{THEN } \text{STATE_ACTIVE}$
 $\quad \text{ELSE } \text{STATE_RECOVERING}]$
 $\wedge \text{pending_proposal}' = [\text{pending_proposal} \text{ EXCEPT } ![\text{mon}] = \text{Nil}]$
 $\wedge \text{new_value}' = [\text{new_value} \text{ EXCEPT } ![\text{mon}] = \text{Nil}]$

$\wedge \text{number_refreshes}' = \text{number_refreshes} + 1$
 $\wedge \text{UNCHANGED } \langle \text{epoch}, \text{accepted} \rangle$
 $\wedge \text{UNCHANGED } \langle \text{restart_vars}, \text{data_vars}, \text{collect_vars}, \text{lease_vars}, \text{auxiliary_vars} \rangle$

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.

$\text{Timeout}(\text{mon}) \triangleq$
 $\wedge \text{phase}[\text{mon}] = \text{PHASE_COLLECT} \vee \text{phase}[\text{mon}] = \text{PHASE_BEGIN}$
 $\wedge \text{bootstrap}$
 $\wedge \text{send_queue}' = [m \in \text{Monitors} \mapsto \langle \rangle]$
 $\wedge \text{messages}' = [m \in \text{DOMAIN messages} \mapsto 0]$
 $\wedge \text{UNCHANGED } \langle \text{state_vars}, \text{restart_vars}, \text{data_vars}, \text{collect_vars}, \text{lease_vars}, \text{commit_vars} \rangle$

Dispatchers and next statement

Handle a message.

$\text{Receive}(\text{msg}) \triangleq$
 $\wedge \vee \text{phase}[\text{msg.dest}] = \text{PHASE_COLLECT}$
 $\vee \text{phase}[\text{msg.dest}] = \text{PHASE_BEGIN}$
 $\vee \text{phase}[\text{msg.dest}] = \text{PHASE_ELECTION}$
 \wedge
 $\vee \wedge \text{msg.type} = \text{OP_COLLECT}$
 $\wedge \text{handle_collect}(\text{msg.dest}, \text{msg})$
 $\wedge \text{step_x}' = \text{"receive collect"}$
 $\vee \wedge \text{msg.type} = \text{OP_LAST}$
 $\wedge \text{handle_last}(\text{msg.dest}, \text{msg})$
 $\wedge \text{step_x}' = \text{"receive last"}$
 $\vee \wedge \text{msg.type} = \text{OP_LEASE}$
 $\wedge \text{handle_lease}(\text{msg.dest}, \text{msg})$
 $\wedge \text{step_x}' = \text{"receive lease"}$
 $\vee \wedge \text{msg.type} = \text{OP_LEASE_ACK}$
 $\wedge \text{handle_lease_ack}(\text{msg.dest}, \text{msg})$
 $\wedge \text{step_x}' = \text{"receive lease_ack"}$
 $\vee \wedge \text{msg.type} = \text{OP_BEGIN}$
 $\wedge \text{handle_begin}(\text{msg.dest}, \text{msg})$
 $\wedge \text{step_x}' = \text{"receive begin"}$
 $\vee \wedge \text{msg.type} = \text{OP_ACCEPT}$
 $\wedge \text{handle_accept}(\text{msg.dest}, \text{msg})$
 $\wedge \text{step_x}' = \text{"receive accept"}$
 $\vee \wedge \text{msg.type} = \text{OP_COMMIT}$

$\wedge \text{handle_commit}(\text{msg.dest}, \text{msg})$
 $\wedge \text{step}_{x'} = \text{"receive commit"}$

Send a message from the event queue.

$\text{Send_from_queue}(\text{mon}) \triangleq$
 $\wedge \text{Len}(\text{send_queue}[\text{mon}]) > 0$
 \wedge
 $\vee \wedge \text{Head}(\text{send_queue}[\text{mon}])[2] = \text{OP_COLLECT}$
 $\wedge \text{send_collect}(\text{mon}, \text{Head}(\text{send_queue}[\text{mon}])[1])$
 $\wedge \text{step}_{x'} = \text{"send_collect"}$
 $\vee \wedge \text{Head}(\text{send_queue}[\text{mon}])[2] = \text{OP_LEASE}$
 $\wedge \text{send_extend_lease}(\text{mon}, \text{Head}(\text{send_queue}[\text{mon}])[1])$
 $\wedge \text{step}_{x'} = \text{"send_extend_lease"}$
 $\vee \wedge \text{Head}(\text{send_queue}[\text{mon}])[2] = \text{OP_BEGIN}$
 $\wedge \text{send_begin}(\text{mon}, \text{Head}(\text{send_queue}[\text{mon}])[1])$
 $\wedge \text{step}_{x'} = \text{"send_begin"}$
 $\vee \wedge \text{Head}(\text{send_queue}[\text{mon}])[2] = \text{OP_COMMIT}$
 $\wedge \text{send_commit}(\text{mon}, \text{Head}(\text{send_queue}[\text{mon}])[1])$
 $\wedge \text{step}_{x'} = \text{"send_commit"}$
 $\wedge \text{send_queue}' = [\text{send_queue} \text{ EXCEPT } ![\text{mon}] = \text{Tail}(\text{send_queue}[\text{mon}])]$

Limit some variables to reduce search space.

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

State transitions.

$\text{Next} \triangleq$
 $\wedge \text{reduce_search_space}$
 $\wedge \text{IF } \text{epoch} \% 2 = 1 \text{ THEN}$
 $\wedge \text{leader_election}$
 $\wedge \text{step}_{x'} = \text{"election"} \wedge \text{step}' = \text{step} + 1$
 $\wedge \text{UNCHANGED } \text{number_refreshes}$
 ELSE
 $\text{IF } \exists \text{mon} \in \text{Monitors} : \text{Len}(\text{send_queue}[\text{mon}]) > 0$
 $\text{THEN } \wedge \exists \text{mon} \in \text{Monitors} : \text{Send_from_queue}(\text{mon})$
 $\wedge \text{step}' = \text{step} + 1$
 $\wedge \text{UNCHANGED } \text{number_refreshes}$
 ELSE
 \vee
 $\wedge \exists \text{mon} \in \text{Monitors} : \text{election_recover}(\text{mon})$

$\wedge \text{step_x}' = \text{"election_recover"} \wedge \text{step}' = \text{step} + 1$
 $\wedge \text{UNCHANGED } \text{number_refreshes}$

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

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

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

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

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

$\vee \wedge \exists \text{mon} \in \text{Monitors} : \exists v \in \text{Value_set} : \text{client_request}(\text{mon}, v)$
 $\wedge \text{step_x}' = \text{"client_request"} \wedge \text{step}' = \text{step} + 1$
 $\wedge \text{UNCHANGED } \text{number_refreshes}$

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

$\vee \wedge \exists m \in \text{ValidMessage}(\text{messages}) : \text{Receive}(m)$
 $\wedge \text{step}' = \text{step} + 1$
 $\wedge \text{UNCHANGED } \text{number_refreshes}$

$\vee \wedge \exists \text{mon} \in \text{Monitors} : \text{restart_mon}(\text{mon})$
 $\wedge \text{step_x}' = \text{"restart mon"} \wedge \text{step}' = \text{step} + 1$

$\vee \wedge \exists \text{mon} \in \text{Monitors} : \text{Timeout}(\text{mon})$
 $\wedge \text{step_x}' = \text{"timeout and restart"} \wedge \text{step}' = \text{step} + 1$
 $\wedge \text{UNCHANGED } \text{number_refreshes}$

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'.

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

Invariants to test in model check

$$Inv \triangleq \wedge \text{TRUE}$$

$$\wedge Inv_diam(20)$$

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(
 $\exists msg1, msg2 \in \text{DOMAIN } messages :$
 $\wedge msg1.type = OP_COLLECT \wedge msg2.type = OP_COLLECT$
 $\wedge msg1.from \neq msg2.from$
)

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

Inv_find_state(
 $\wedge step_x = \text{"restart mon"}$
 $\setminus * \text{The system is in collect phase and no } OP_LAST \text{ message has been received.}$
 $\setminus * isLeader[mon] = \text{TRUE}$ assures that the leader was not the one that crashed.
 $\wedge \exists mon \in \text{Monitors} :$
 $\wedge isLeader[mon] = \text{TRUE}$
 $\wedge phase[mon] = PHASE_COLLECT$
 $\wedge num_last[mon] = 1$
 $\setminus * \text{All the collect requests have been handled by the peers.}$
 $\wedge \forall m \in \text{ValidMessage}(messages) :$
 $m.type \neq OP_COLLECT$
 $\wedge number_refreshes = 2$
 $\wedge epoch = 2$
)

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

Inv_find_state(
 $\wedge step_x = \text{"restart mon"}$
 $\wedge \exists msg \in \text{ValidMessage}(messages) :$
 $msg.type = OP_ACCEPT$
 $\wedge \forall mon \in \text{Monitors} :$
 $isLeader[mon] = \text{FALSE}$
 $\wedge number_refreshes = 2$
 $\wedge epoch = 2$
)

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

$\setminus * \text{Modification History}$
 $\setminus * \text{Last modified Mon Mar 01 11:27:57 WET 2021 by } afonsonf$

\ * Created *Mon Jan 11 16:15:26 WET 2021* by *afonsof*