

MODULE *ceph*

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>

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

EXTENDS *Integers, FiniteSets, Sequences, TLC*

**Utils**

*Max* element from a set.

@type: *Set(Int)*  $\Rightarrow$  *Int*;

$Max(S) \triangleq \text{CHOOSE } x \in S : \forall y \in S : x \geq y$

*Min* element from a set.

@type: *Set(Int)*  $\Rightarrow$  *Int*;

$Min(S) \triangleq \text{CHOOSE } x \in S : \forall y \in S : x \leq y$

Set of monitors to a sequence.

RECURSIVE *SetToSeq*(-)

@type: *Set(MONITOR)*  $\Rightarrow$  *Seq(MONITOR)*;

*SetToSeq*(*S*)  $\triangleq$

IF *S* = {} THEN  $\langle \rangle$

ELSE LET *x*  $\triangleq$  CHOOSE *x*  $\in$  *S* : TRUE

IN  $\langle x \rangle \circ SetToSeq(S \setminus \{x\})$

**Constants**

Set of *Monitors*.

CONSTANTS @type: *Set(MONITOR)*; *Monitors*

Sequence of monitors.

@type: *Seq(MONITOR)*;

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

Number of monitors.

@type: *Int*;

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

Rank predicate, used to compute proposal numbers.

@type: *MONITOR*  $\Rightarrow$  *Int*;

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

Set of possible values.

CONSTANTS @type: *Set(VALUE)*; *Value\_set*

Predicate used in the cfg file to define the symmetry set.

Workaround for typechecker.  
 $\text{@typeAlias: } \text{MONITOR} = T;$   
 $\text{@typeAlias: } \text{VALUE} = T;$   
 $\text{SYMM} \triangleq \text{Permutations}(\text{Monitors}) \cup \text{Permutations}(\text{Value\_set})$

Reserved value.  
 $\text{CONSTANTS } \text{@type: VALUE; Nil}$

Paxos states.  
 $\text{CONSTANTS } \text{@type: STATE\_NAME; STATE\_RECOVERING, @type: STATE\_NAME; STATE\_ACTIVE, @type: STATE\_NAME; STATE\_UPDATING, @type: STATE\_NAME; STATE\_UPDATING\_PREVIOUS, @type: STATE\_NAME; STATE\_WRITING, @type: STATE\_NAME; STATE\_WRITING\_PREVIOUS, @type: STATE\_NAME; STATE\_REFRESH, @type: STATE\_NAME; STATE\_SHUTDOWN}$

Paxos auxiliary phase states.  
 They are used to force some sequence of steps.  
 $\text{CONSTANTS } \text{@type: PHASE\_NAME; PHASE\_ELECTION, @type: PHASE\_NAME; PHASE\_SEND\_COLLECT, @type: PHASE\_NAME; PHASE\_COLLECT, @type: PHASE\_NAME; PHASE\_LEASE, @type: PHASE\_NAME; PHASE\_LEASE\_DONE, @type: PHASE\_NAME; PHASE\_BEGIN, @type: PHASE\_NAME; PHASE\_COMMIT}$

Paxos message types.  
 $\text{CONSTANTS } \text{@type: MESSAGE\_OP; OP\_COLLECT, @type: MESSAGE\_OP; OP\_LAST, @type: MESSAGE\_OP; OP\_BEGIN, @type: MESSAGE\_OP; OP\_ACCEPT, @type: MESSAGE\_OP; OP\_COMMIT, @type: MESSAGE\_OP; OP\_LEASE, @type: MESSAGE\_OP; OP\_LEASE\_ACK}$

### Global variables

Integer representing the current epoch. If is odd trigger an election.  
 $\text{VARIABLE } \text{@type: Int; epoch}$

Store messages waiting to be handled.  
 $\text{VARIABLE } \text{@type: MONITOR} \rightarrow (\text{MONITOR} \rightarrow \text{Seq}(\text{MESSAGE})); \text{messages}$

Stores history of messages. Can be useful to find specific states.  
 $\text{VARIABLE } \text{@type: Set}(\text{MESSAGE}); \text{message\_history}$

Stores if a monitor is up or down. All available monitors, in a given epoch, are part of the quorum.  
 $\text{VARIABLE } \text{@type: MONITOR} \rightarrow \text{Bool}; \text{quorum}$

Size of the current quorum.  
 $\text{VARIABLE } \text{@type: Int; quorum\_sz}$

### State variables

A function that stores the current leader.  $\text{isLeader}[mon]$  is True iff  $mon$  is a leader, else False.  
 $\text{VARIABLE } \text{@type: MONITOR} \rightarrow \text{Bool}; \text{isLeader}$

A function that stores the state of each monitor.  
VARIABLE @type: *MONITOR*  $\rightarrow$  *STATE\_NAME*; *state*

A function that stores the phase of each monitor.  
VARIABLE @type: *MONITOR*  $\rightarrow$  *PHASE\_NAME*; *phase*

#### Restart variables

A function that stores, for each monitor, a proposal number when the commit phase starts.  
This proposal number can be retrieved after a monitor crashes and restarts.  
VARIABLE @type: *MONITOR*  $\rightarrow$  *PN*; *pending\_pn*

A function that stores, for each monitor, a value version when the commit phase starts.  
This value version can be retrieved after a monitor crashes and restarts.  
VARIABLE @type: *MONITOR*  $\rightarrow$  *VALUE\_VERSION*; *pending\_v*

A function that stores, for each monitor, the best uncommitted *pn* received in the collect phase.  
VARIABLE @type: *MONITOR*  $\rightarrow$  *PN*; *uncommitted\_pn*

A function that stores, for each monitor, the best uncommitted value version received in the collect phase.  
VARIABLE @type: *MONITOR*  $\rightarrow$  *VALUE\_VERSION*; *uncommitted\_v*

A function that stores, for each monitor, the best uncommitted value received in the collect phase.  
VARIABLE @type: *MONITOR*  $\rightarrow$  *VALUE*; *uncommitted\_value*

#### Data variables

A function that stores, for each monitor, the store where the transactions are applied.  
In this model, a transaction represents changing the value in the store.  
VARIABLE @type: *MONITOR*  $\rightarrow$  *VALUE*; *monitor\_store*

A function that stores the transaction log of each monitor.  
VARIABLE @type: *MONITOR*  $\rightarrow$  (*VALUE\_VERSION*  $\rightarrow$  *VALUE*); *values*

A function that stores the last proposal number accepted by each monitor.  
VARIABLE @type: *MONITOR*  $\rightarrow$  *PN*; *accepted\_pn*

A function that stores the first value version committed by each monitor.  
VARIABLE @type: *MONITOR*  $\rightarrow$  *VALUE\_VERSION*; *first\_committed*

A function that stores the last value version committed by each monitor.  
VARIABLE @type: *MONITOR*  $\rightarrow$  *VALUE\_VERSION*; *last\_committed*

#### Collect phase variables

A function that stores the number of peers that accepted a collect request.  
VARIABLE @type: *MONITOR*  $\rightarrow$  *Int*; *num\_last*

Used by leader when receiving responses in collect phase.

VARIABLE @type:  $MONITOR \rightarrow (MONITOR \rightarrow VALUE\_VERSION)$ ; *peer\_first\_committed*

Used by leader when receiving responses in collect phase.

VARIABLE @type:  $MONITOR \rightarrow (MONITOR \rightarrow VALUE\_VERSION)$ ; *peer\_last\_committed*

#### Lease phase variables

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

VARIABLE @type:  $MONITOR \rightarrow (MONITOR \rightarrow Bool)$ ; *acked\_lease*

#### Commit phase variables

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

VARIABLE @type:  $MONITOR \rightarrow VALUE$ ; *pending\_proposal*

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

VARIABLE @type:  $MONITOR \rightarrow VALUE$ ; *new\_value*

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

VARIABLE @type:  $MONITOR \rightarrow (MONITOR \rightarrow Bool)$ ; *accepted*

#### Debug variables

Variables to help debug a behavior.

*step* is the diameter of a behavior/path.

*step\_name* the current predicate being called.

VARIABLE @type: *Str*; *step\_name*

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

This variable is used to limit the search space.

VARIABLE @type: *Int*; *number\_crashes*

#### Variables initialization

@typeAlias:  $VALUE\_VERSION = Int$ ;

@typeAlias:  $PN = Int$ ;

*global\_vars*  $\triangleq \langle epoch, messages, message\_history, quorum, quorum\_sz \rangle$

*state\_vars*  $\triangleq \langle isLeader, state, phase \rangle$

*restart\_vars*  $\triangleq \langle pending\_pn, pending\_v, uncommitted\_pn, uncommitted\_v, uncommitted\_value \rangle$

*data\_vars*  $\triangleq \langle monitor\_store, values, accepted\_pn, first\_committed, last\_committed \rangle$

*collect\_vars*  $\triangleq \langle num\_last, peer\_first\_committed, peer\_last\_committed \rangle$

*lease\_vars*  $\triangleq \langle acked\_lease \rangle$

*commit\_vars*  $\triangleq \langle pending\_proposal, new\_value, accepted \rangle$

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

*Init\_global\_vars*  $\triangleq$

$$\begin{aligned}
& \wedge \text{epoch} = 1 \\
& \wedge \text{messages} = [\text{mon1} \in \text{Monitors} \mapsto [\text{mon2} \in \text{Monitors} \mapsto \langle \rangle]] \\
& \wedge \text{message\_history} = \{\} \\
& \wedge \text{quorum} = [\text{mon} \in \text{Monitors} \mapsto \text{TRUE}] \\
& \wedge \text{quorum\_sz} = \text{MonitorsLen} \\
\\
\text{Init\_state\_vars} & \triangleq \\
& \wedge \text{isLeader} = [\text{mon} \in \text{Monitors} \mapsto \text{FALSE}] \\
& \wedge \text{state} = [\text{mon} \in \text{Monitors} \mapsto \text{STATE\_RECOVERING}] \\
& \wedge \text{phase} = [\text{mon} \in \text{Monitors} \mapsto \text{PHASE\_ELECTION}] \\
\\
\text{Init\_restart\_vars} & \triangleq \\
& \wedge \text{pending\_pn} = [\text{mon} \in \text{Monitors} \mapsto 0] \\
& \wedge \text{pending\_v} = [\text{mon} \in \text{Monitors} \mapsto 0] \\
& \wedge \text{uncommitted\_pn} = [\text{mon} \in \text{Monitors} \mapsto 0] \\
& \wedge \text{uncommitted\_v} = [\text{mon} \in \text{Monitors} \mapsto 0] \\
& \wedge \text{uncommitted\_value} = [\text{mon} \in \text{Monitors} \mapsto \text{Nil}] \\
\\
\text{Init\_data\_vars} & \triangleq \\
& \wedge \text{monitor\_store} = [\text{mon} \in \text{Monitors} \mapsto \text{Nil}] \\
& \wedge \text{values} = [\text{mon} \in \text{Monitors} \mapsto [\text{version} \in \{\} \mapsto \text{Nil}]] \\
& \wedge \text{accepted\_pn} = [\text{mon} \in \text{Monitors} \mapsto 0] \\
& \wedge \text{first\_committed} = [\text{mon} \in \text{Monitors} \mapsto 0] \\
& \wedge \text{last\_committed} = [\text{mon} \in \text{Monitors} \mapsto 0] \\
\\
\text{Init\_collect\_vars} & \triangleq \\
& \wedge \text{num\_last} = [\text{mon} \in \text{Monitors} \mapsto 0] \\
& \wedge \text{peer\_first\_committed} = [\text{mon1} \in \text{Monitors} \mapsto [\text{mon2} \in \text{Monitors} \mapsto -1]] \\
& \wedge \text{peer\_last\_committed} = [\text{mon1} \in \text{Monitors} \mapsto [\text{mon2} \in \text{Monitors} \mapsto -1]] \\
\\
\text{Init\_lease\_vars} & \triangleq \\
& \wedge \text{acked\_lease} = [\text{mon1} \in \text{Monitors} \mapsto [\text{mon2} \in \text{Monitors} \mapsto \text{FALSE}]] \\
\\
\text{Init\_commit\_vars} & \triangleq \\
& \wedge \text{pending\_proposal} = [\text{mon} \in \text{Monitors} \mapsto \text{Nil}] \\
& \wedge \text{new\_value} = [\text{mon} \in \text{Monitors} \mapsto \text{Nil}] \\
& \wedge \text{accepted} = [\text{mon1} \in \text{Monitors} \mapsto [\text{mon2} \in \text{Monitors} \mapsto \text{FALSE}]] \\
\\
\text{Init} & \triangleq \\
& \wedge \text{Init\_global\_vars} \\
& \wedge \text{Init\_state\_vars} \\
& \wedge \text{Init\_restart\_vars} \\
& \wedge \text{Init\_data\_vars} \\
& \wedge \text{Init\_collect\_vars} \\
& \wedge \text{Init\_lease\_vars} \\
& \wedge \text{Init\_commit\_vars} \\
& \wedge \text{step\_name} = \text{"init"} \wedge \text{number\_crashes} = 0
\end{aligned}$$

## Message manipulation

@typeAlias: MESSAGE = [type: MESSAGE\_OP, from: MONITOR, dest: MONITOR,  
first\_committed: VALUE\_VERSION, last\_committed:  
VALUE\_VERSION, values: (VALUE\_VERSION → VALUE),  
uncommitted\_pn: PN, pn: PN];

@typeAlias: MESSAGE\_QUEUE = MONITOR → (MONITOR → Seq(MESSAGE));

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

Converts a set with at most one element to a sequence.

@type: Set(MESSAGE) ⇒ Seq(MESSAGE);

SingleMessageSetToSeq(*S*)  $\triangleq$   
IF  $\exists elem \in S : \text{TRUE}$  THEN LET *elem*  $\triangleq$  CHOOSE  $x \in S : \text{TRUE}$   
IN  $\langle elem \rangle$   
ELSE  $\langle \rangle$

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

@type: (MESSAGE, MESSAGE\_QUEUE) ⇒ MESSAGE\_QUEUE;

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

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

@type: (MESSAGE, MESSAGE\_QUEUE) ⇒ MESSAGE\_QUEUE;

WithoutMessage(*m*, *msgs*)  $\triangleq$   
[*msgs* EXCEPT ![*m.from*] =  
[*msgs*[*m.from*] EXCEPT ![*m.dest*] = Tail(*msgs*[*m.from*][*m.dest*))]]

Adds the message *m* to the network.

Variables changed: *messages*, *message\_history*.

@type: MESSAGE ⇒ Bool;

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

Adds a set of messages to the network.

Variables changed: *messages*, *message\_history*.

@type: (MONITOR, Set(MESSAGE)) ⇒ Bool;

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

Removes the request from network and adds the response.

Variables changed: *messages*, *message\_history*.

@type: (MESSAGE, MESSAGE)  $\Rightarrow$  Bool;

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

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

Variables changed: *messages*, *message\_history*.

@type: (MONITOR, Set(MESSAGE), MESSAGE)  $\Rightarrow$  Bool;

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

Removes message *m* from the network.

Variables changed: *messages*, *message\_history*.

@type: MESSAGE  $\Rightarrow$  Bool;

$Discard(m) \triangleq$   
 $\wedge\ messages' = WithoutMessage(m, messages)$   
 $\wedge\ UNCHANGED\ message\_history$

### Helper predicates

Computes a new unique proposal number for a given monitor.

Version A - Equal to the one in the source.

This version breaks the symmetry of the monitor set.

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

@type: (MONITOR, Int)  $\Rightarrow$  Int;

$get\_new\_proposal\_number(mon, oldpn) \triangleq ((oldpn \div 100) + 1) * 100 + rank(mon)$

Version B - Adapted to not break symmetry.

Example:  $oldpn = 300$ ,  $rank(mon) = 5$ ,  $newpn = 400$ .

@type: (MONITOR, Int)  $\Rightarrow$  Int;

$get\_new\_proposal\_number(mon, oldpn) \triangleq ((oldpn \div 100) + 1) * 100 + epoch$

Clear the variable *peer\_first\_committed*.

Variables changed: *peer\_first\_committed*.

@type: MONITOR  $\Rightarrow$  Bool;

$clear\_peer\_first\_committed(mon) \triangleq$   
 $peer\_first\_committed' = [peer\_first\_committed\ EXCEPT\ ![mon] =$

$[m \in \text{Monitors} \mapsto -1]$

Clear the variable *peer\_last\_committed*.

Variables changed: *peer\_last\_committed*.

@type: *MONITOR*  $\Rightarrow$  *Bool*;

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

Store peer values and update *first\_committed*, *last\_committed* and *monitor\_store* accordingly.

Variables changed: *values*, *first\_committed*, *last\_committed*, *monitor\_store*.

@type: (*MONITOR*, *MESSAGE*)  $\Rightarrow$  *Bool*;

$\text{store\_state}(\text{mon}, \text{msg}) \triangleq$   
 Choose peer values from *mon* last committed + 1 to peer last committed.  
 $\wedge \text{LET } \text{logs} \triangleq (\text{DOMAIN } \text{msg.values}) \cap (\text{last\_committed}[mon] + 1 .. \text{msg.last\_committed})$   
 IN  $\wedge \text{values}' = [\text{values} \text{ EXCEPT } ![mon] =$   
 $[i \in \text{DOMAIN } \text{values}[mon] \cup \text{logs} \mapsto$   
 $\text{IF } i \in \text{logs}$   
 $\text{THEN } \text{msg.values}[i]$   
 $\text{ELSE } \text{values}[mon][i]]]$   
 Update last committed and first committed.  
 $\wedge \text{last\_committed}' = [\text{last\_committed} \text{ EXCEPT } ![mon] = \text{Max}(\text{logs} \cup \{\text{last\_committed}[mon]\})]$   
 $\wedge \text{IF } \text{logs} \neq \{\}$   $\wedge \text{first\_committed}[mon] = 0$   
 $\text{THEN } \text{first\_committed}' =$   
 $[\text{first\_committed} \text{ EXCEPT } ![mon] = \text{Min}(\text{logs})]$   
 $\text{ELSE } \text{first\_committed}' =$   
 $[\text{first\_committed} \text{ EXCEPT } ![mon] = \text{Min}(\text{logs} \cup \{\text{first\_committed}[mon]\})]$   
 Update monitor store.  
 $\wedge \text{IF } \text{last\_committed}'[mon] = 0$   
 $\text{THEN UNCHANGED } \text{monitor\_store}$   
 $\text{ELSE } \text{monitor\_store}' = [\text{monitor\_store} \text{ EXCEPT } ![mon] = \text{values}'[mon][\text{last\_committed}'[mon]]]$

Check if uncommitted value version is still valid, else reset it.

Variables changed: *uncommitted\_pn*, *uncommitted\_v*, *uncommitted\_value*.

@type: *MONITOR*  $\Rightarrow$  *Bool*;

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

Trigger new election by incrementing epoch.

Variables changed: *epoch*.

@type: *Bool*;

$\text{bootstrap} \triangleq$



$\wedge epoch' = epoch + 1$

### Lease phase predicates

Changes *mon* state to *STATE\_ACTIVE*.

Variables changed: *state*.

@type: *MONITOR*  $\Rightarrow$  *Bool*;

$finish\_round(mon) \triangleq$   
 $\wedge isLeader[mon] = \text{TRUE}$   
 $\wedge state' = [state \text{ EXCEPT } ![mon] = \text{STATE\_ACTIVE}]$

Resets the variable *acked\_lease* and send lease messages to peers.

Variables changed: *acked\_lease*, *messages*, *message\_history*, *phase*.

@type: *MONITOR*  $\Rightarrow$  *Bool*;

$extend\_lease(mon) \triangleq$   
 $\wedge isLeader[mon] = \text{TRUE}$   
 $\wedge acked\_lease' = [acked\_lease \text{ EXCEPT } ![mon] =$   
 $\quad [m \in Monitors \mapsto \text{IF } m = mon \text{ THEN TRUE ELSE FALSE}]]$   
 $\wedge Send\_set(mon,$   
 $\quad \{[type \mapsto OP\_LEASE,$   
 $\quad \quad from \mapsto mon,$   
 $\quad \quad dest \mapsto dest,$   
 $\quad \quad last\_committed \mapsto last\_committed[mon]] : dest \in \{m \in Monitors \setminus \{mon\} : quorum[m]\}$   
 $\quad \})$   
 $\wedge phase' = [phase \text{ EXCEPT } ![mon] = \text{PHASE\_LEASE}]$

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

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

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

Variables changed: *messages*, *message\_history*, *state*.

@type: (*MONITOR*, *MESSAGE*)  $\Rightarrow$  *Bool*;

$handle\_lease(mon, msg) \triangleq$   
 $\wedge \text{discard if not peon or peon is behind}$   
 $\text{IF } \vee isLeader[mon] = \text{TRUE}$   
 $\quad \vee last\_committed[mon] \neq msg.last\_committed$   
 $\text{THEN } \wedge Discard(msg)$   
 $\quad \wedge \text{UNCHANGED } state$   
 $\text{ELSE } \wedge state' = [state \text{ EXCEPT } ![mon] = \text{STATE\_ACTIVE}]$   
 $\quad \wedge Reply([type \mapsto OP\_LEASE\_ACK,$   
 $\quad \quad from \mapsto mon,$   
 $\quad \quad dest \mapsto msg.from,$   
 $\quad \quad first\_committed \mapsto first\_committed[mon],$   
 $\quad \quad last\_committed \mapsto last\_committed[mon]], msg)$   
 $\quad \wedge Discard(msg)$   
 $\wedge \text{UNCHANGED } \langle epoch, quorum, quorum\_sz, isLeader, phase \rangle$

$\wedge \text{UNCHANGED } \langle \text{restart\_vars}, \text{data\_vars}, \text{collect\_vars}, \text{lease\_vars}, \text{commit\_vars} \rangle$

Handle a lease ack message. The leader updates the *acked\_lease* variable.

Because the *lease\_ack* messages are not sent, this predicate is never called.

The reasoning for this is given in *handle\_lease* comment.

Variables changed: *acked\_lease*, *messages*, *message\_history*.

@type: (MONITOR, MESSAGE)  $\Rightarrow$  Bool;

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

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

Because the *lease\_ack* messages are not sent, this predicate is never called.

The reasoning for this is given in *handle\_lease* comment.

Variables changed: *phase*.

@type: MONITOR  $\Rightarrow$  Bool;

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

### Commit phase predicates

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

The new value is stored in *values* and *pending\_pn* is assigned in order for the leader to be able to recover from a crash.

Variables changed: *accepted*, *new\_value*, *phase*, *messages*, *message\_history*, *values*, *pending\_pn*, *pending\_v*.

@type: (MONITOR, VALUE)  $\Rightarrow$  Bool;

$\text{begin}(\text{mon}, v) \triangleq$   
 $\wedge \text{isLeader}[\text{mon}] = \text{TRUE}$   
 $\wedge \vee \text{state}'[\text{mon}] = \text{STATE\_UPDATING}$   
 $\quad \vee \text{state}'[\text{mon}] = \text{STATE\_UPDATING\_PREVIOUS}$   
 $\wedge \text{quorum\_sz} = 1 \vee \text{num\_last}[\text{mon}] > \text{MonitorsLen} \div 2$   
 $\wedge \text{new\_value}[\text{mon}] = \text{Nil}$   
 $\wedge \text{accepted}' = [\text{accepted} \text{ EXCEPT } ![\text{mon}] =$   
 $\quad [m \in \text{Monitors} \mapsto \text{IF } m = \text{mon} \text{ THEN TRUE ELSE FALSE}]]$   
 $\wedge \text{new\_value}' = [\text{new\_value} \text{ EXCEPT } ![\text{mon}] = v]$   
 $\wedge \text{phase}' = [\text{phase} \text{ EXCEPT } ![\text{mon}] = \text{LEASE\_BEGIN}]$

$$\begin{aligned}
& \wedge \text{values}' = [\text{values} \text{ EXCEPT } ![mon] = \\
& \quad ((\text{last\_committed}[mon] + 1) :> \text{new\_value}'[mon]) @@ \text{values}[mon]] \\
& \wedge \text{Send\_set}(mon, \\
& \quad \{[type \quad \mapsto OP\_BEGIN, \\
& \quad \quad from \quad \mapsto mon, \\
& \quad \quad dest \quad \mapsto dest, \\
& \quad \quad last\_committed \mapsto last\_committed[mon], \\
& \quad \quad values \quad \mapsto \text{values}'[mon], \\
& \quad \quad pn \quad \mapsto \text{accepted\_pn}[mon]] : dest \in \{m \in \text{Monitors} \setminus \{mon\} : \text{quorum}[m]\} \\
& \quad \}) \\
& \wedge \text{pending\_pn}' = [\text{pending\_pn} \text{ EXCEPT } ![mon] = \text{accepted\_pn}[mon]] \\
& \wedge \text{pending\_v}' = [\text{pending\_v} \text{ EXCEPT } ![mon] = last\_committed[mon] + 1]
\end{aligned}$$

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, values and pending\_pn are assigned in order for the monitor to recover in case of a crash.

Variables changed: messages, message\_history, state, values, pending\_pn, pending\_v.

@type: (MONITOR, MESSAGE)  $\Rightarrow$  Bool;

handle\_begin(mon, msg)  $\triangleq$

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

$\wedge \text{IF } msg.pn < \text{accepted\_pn}[mon]$

THEN

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

$\wedge \text{UNCHANGED } \langle state, values, pending\_pn, pending\_v \rangle$

ELSE

$\wedge msg.pn = \text{accepted\_pn}[mon]$

$\wedge msg.last\_committed = last\_committed[mon]$

$\text{assign } values[mon][last\_committed[mon] + 1]$

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

$((last\_committed[mon] + 1) :> msg.values[last\_committed[mon] + 1]) @@ values[mon]]$

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

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

$\wedge \text{pending\_v}' = [\text{pending\_v} \text{ EXCEPT } ![mon] = last\_committed[mon] + 1]$

$\wedge \text{Reply}([type \quad \mapsto OP\_ACCEPT,$

$\quad from \quad \mapsto mon,$

$\quad dest \quad \mapsto msg.from,$

$\quad last\_committed \mapsto last\_committed[mon],$

$\quad pn \quad \mapsto \text{accepted\_pn}[mon]], msg)$

$\wedge \text{UNCHANGED } \langle epoch, quorum, quorum\_sz, isLeader, phase, monitor\_store,$   
 $\quad \text{accepted\_pn, first\_committed, last\_committed, uncommitted\_pn,$   
 $\quad \text{uncommitted\_v, uncommitted\_value} \rangle$

$\wedge \text{UNCHANGED } \langle collect\_vars, lease\_vars, commit\_vars \rangle$

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`, `message_history`, `accepted`

@type: (MONITOR, MESSAGE)  $\Rightarrow$  Bool;

```

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

```

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

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

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

@type: MONITOR  $\Rightarrow$  Bool;

```

post_accept(mon)  $\triangleq$ 
   $\wedge$  phase[mon] = PHASE_BEGIN
   $\wedge$   $\forall m \in$  Monitors : quorum[m]  $\Rightarrow$  accepted[mon][m] = TRUE
   $\wedge$  new_value[mon]  $\neq$  Nil
   $\wedge$   $\vee$  state[mon] = STATE_UPDATING_PREVIOUS
     $\vee$  state[mon] = STATE_UPDATING
   $\wedge$  last_committed' = [last_committed EXCEPT ![mon] = last_committed[mon] + 1]

   $\wedge$  IF first_committed[mon] = 0
    THEN first_committed' = [first_committed EXCEPT ![mon] = first_committed[mon] + 1]
    ELSE UNCHANGED first_committed

   $\wedge$  monitor_store' = [monitor_store EXCEPT ![mon] = values[mon][last_committed[mon] + 1]]
   $\wedge$  new_value' = [new_value EXCEPT ![mon] = Nil]
   $\wedge$  Send_set(mon,
    {[type            $\mapsto$  OP_COMMIT,
      from            $\mapsto$  mon,
      dest            $\mapsto$  dest,
      last_committed  $\mapsto$  last_committed'[mon],
      pn              $\mapsto$  accepted_pn[mon],
      values          $\mapsto$  values[mon]] : dest  $\in$  {m  $\in$  Monitors \ {mon} : quorum[m]}
    })

```

$\wedge state' = [state \text{ EXCEPT } ![mon] = STATE\_REFRESH]$   
 $\wedge phase' = [phase \text{ EXCEPT } ![mon] = PHASE\_COMMIT]$   
 $\wedge \text{UNCHANGED } \langle isLeader, values, accepted\_pn, pending\_proposal, accepted \rangle$   
 $\wedge \text{UNCHANGED } \langle epoch, quorum, quorum\_sz, restart\_vars, collect\_vars, lease\_vars \rangle$

Predicate that is called after *post\_accept*. The leader finishes the commit phase by updating his state to *STATE\_ACTIVE* and by extending the lease to his peers.

Variables changed: *state, phase, acked\_lease, messages, message\_history*.

@type: *MONITOR*  $\Rightarrow$  *Bool*;

$finish\_commit(mon) \triangleq$   
 $\wedge state[mon] = STATE\_REFRESH$   
 $\wedge phase[mon] = PHASE\_COMMIT$   
 $\wedge finish\_round(mon)$   
 $\wedge extend\_lease(mon)$   
 $\wedge \text{UNCHANGED } \langle epoch, quorum, quorum\_sz, isLeader \rangle$   
 $\wedge \text{UNCHANGED } \langle restart\_vars, data\_vars, collect\_vars, commit\_vars \rangle$

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

Variables changed: *messages, message\_history, values, first\_committed, last\_committed, monitor\_store, uncommitted\_v, uncommitted\_pn, uncommitted\_value*.

@type: (*MONITOR, MESSAGE*)  $\Rightarrow$  *Bool*;

$handle\_commit(mon, msg) \triangleq$   
 $\wedge isLeader[mon] = \text{FALSE}$   
 $\wedge store\_state(mon, msg)$   
 $\wedge check\_and\_correct\_uncommitted(mon)$   
 $\wedge Discard(msg)$   
 $\wedge \text{UNCHANGED } \langle epoch, quorum, quorum\_sz, accepted\_pn, pending\_pn, pending\_v \rangle$   
 $\wedge \text{UNCHANGED } \langle state\_vars, collect\_vars, lease\_vars, commit\_vars \rangle$

### Client Request

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

Variables changed: *pending\_proposal*.

@type: (*MONITOR, VALUE*)  $\Rightarrow$  *Bool*;

$client\_request(mon, v) \triangleq$   
 $\wedge isLeader[mon] = \text{TRUE}$   
 $\wedge state[mon] = STATE\_ACTIVE$   
 $\wedge pending\_proposal[mon] = Nil$   
 $\wedge pending\_proposal' = [pending\_proposal \text{ EXCEPT } ![mon] = v]$   
 $\wedge \text{UNCHANGED } \langle new\_value, accepted \rangle$   
 $\wedge \text{UNCHANGED } \langle global\_vars, state\_vars, restart\_vars, data\_vars, collect\_vars, lease\_vars \rangle$

Start a commit phase with the value on pending proposal.

Variables changed: *state, pending\_proposal, accepted, new\_value, phase, messages, message\_history, values,*

$pending\_pn, pending\_v.$   
 $@type: MONITOR \Rightarrow Bool;$   
 $propose\_pending(mon) \triangleq$   
 $\wedge phase[mon] = PHASE\_LEASE \vee phase[mon] = PHASE\_ELECTION$   
 $\wedge state[mon] = STATE\_ACTIVE$   
 $\wedge pending\_proposal[mon] \neq Nil$   
 $\wedge pending\_proposal' = [pending\_proposal \text{ EXCEPT } ![mon] = Nil]$   
 $\wedge state' = [state \text{ EXCEPT } ![mon] = STATE\_UPDATING]$   
 $\wedge begin(mon, pending\_proposal[mon])$   
 $\wedge UNCHANGED \langle isLeader, monitor\_store, accepted\_pn, first\_committed, last\_committed,$   
 $epoch, quorum, quorum\_sz, uncommitted\_v, uncommitted\_pn, uncommitted\_value \rangle$   
 $\wedge UNCHANGED \langle collect\_vars, lease\_vars \rangle$

### Collect phase predicates

Start collect phase. This first part of the collect phase is divided in two parts (collect and *send\_collect*) in order to simplify variable changes (when collect is triggered from *handle\_last*).

Variables changed: *accepted\_pn*, *phase*.

$@type: (MONITOR, Int) \Rightarrow Bool;$

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

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

Send collect messages to the peers.

Variables changed: *peer\_first\_committed*, *peer\_last\_committed*, *uncommitted\_pn*, *uncommitted\_v*, *uncommitted\_value*, *num\_last*, *messages*, *message\_history*, *phase*.

$@type: MONITOR \Rightarrow Bool;$

$send\_collect(mon) \triangleq$   
 $\wedge state[mon] = STATE\_RECOVERING$   
 $\wedge isLeader[mon] = TRUE$   
 $\wedge phase[mon] = PHASE\_SEND\_COLLECT$   
 $\wedge clear\_peer\_first\_committed(mon)$   
 $\wedge clear\_peer\_last\_committed(mon)$   
 $\wedge IF \ last\_committed[mon] + 1 \in DOMAIN \ values[mon]$   
 $THEN \wedge uncommitted\_v' =$   
 $\quad [uncommitted\_v \text{ EXCEPT } ![mon] = last\_committed[mon] + 1]$   
 $\wedge uncommitted\_value' =$   
 $\quad [uncommitted\_value \text{ EXCEPT } ![mon] = values[mon][last\_committed[mon] + 1]]$   
 $\wedge uncommitted\_pn' = [uncommitted\_pn \text{ EXCEPT } ![mon] = pending\_pn[mon]]$

```

       $\wedge$  UNCHANGED  $\langle \text{pending\_pn}, \text{pending\_v} \rangle$ 
    ELSE UNCHANGED  $\langle \text{restart\_vars} \rangle$ 

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

```

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

Variables changed: *messages*, *message\_history*, *epoch*, *state*, *accepted\_pn*.

@type: (MONITOR, MESSAGE)  $\Rightarrow$  Bool;

```

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

```

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

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

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

Variables changed: *messages*, *message\_history*, *epoch*, *phase*, *uncommitted\_pn*, *uncommitted\_v*, *uncommitted\_value*, *monitor\_store*, *accepted\_pn*, *first\_committed*, *last\_committed*, *num\_last*, *peer\_first\_committed*, *peer\_last\_committed*.

@type: (MONITOR, MESSAGE)  $\Rightarrow$  Bool;

$handle\_last(mon, msg) \triangleq$

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

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

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

$\wedge \text{IF } msg.first\_committed > last\_committed[mon] + 1$

THEN

$\wedge bootstrap$

$\wedge Discard(msg)$

$\wedge \text{UNCHANGED } \langle num\_last, accepted\_pn, values, phase, monitor\_store \rangle$

$\wedge \text{UNCHANGED } \langle first\_committed, last\_committed, uncommitted\_pn, uncommitted\_v, uncommitted\_value \rangle$

ELSE

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

$\wedge \text{IF } \exists peer \in Monitors :$

$\wedge peer \neq mon$

$\wedge peer\_last\_committed'[mon][peer] \neq -1$

$\wedge peer\_last\_committed'[mon][peer] + 1 < first\_committed[mon]$

$\wedge first\_committed[mon] > 1$

THEN

$\wedge bootstrap$

$\wedge check\_and\_correct\_uncommitted(mon)$

$\wedge Discard(msg)$

$\wedge \text{UNCHANGED } \langle phase, accepted\_pn, num\_last \rangle$

ELSE

$\wedge \text{LET } monitors\_behind \triangleq \{peer \in Monitors :$

$\wedge peer \neq mon$

$\wedge peer\_last\_committed'[mon][peer] \neq -1$

$\wedge peer\_last\_committed'[mon][peer] < last\_committed[mon]$

$\wedge quorum[peer]\}$

IN  $Reply\_set(mon,$

$\{[type \mapsto OP\_COMMIT,$

$from \mapsto mon,$

$dest \mapsto dest,$

$last\_committed \mapsto last\_committed'[mon],$

$pn \mapsto accepted\_pn[mon],$

$values \mapsto values[mon]] : dest \in monitors\_behind$

$\}, msg)$

$\wedge \vee \wedge msg.pn > accepted\_pn[mon]$



$$\begin{aligned}
& \wedge \text{collect}(\text{mon}, \text{msg.pn}) \\
& \wedge \text{check\_and\_correct\_uncommitted}(\text{mon}) \\
& \wedge \text{UNCHANGED } \text{num\_last} \\
\vee & \wedge \text{msg.pn} = \text{accepted\_pn}[\text{mon}] \\
& \wedge \text{num\_last}' = [\text{num\_last} \text{ EXCEPT } ![\text{mon}] = \text{num\_last}[\text{mon}] + 1] \\
& \wedge \text{IF } \wedge \text{msg.last\_committed} + 1 \in \text{DOMAIN } \text{msg.values} \\
& \quad \wedge \text{msg.last\_committed} \geq \text{last\_committed}'[\text{mon}] \\
& \quad \wedge \text{msg.last\_committed} + 1 \geq \text{uncommitted\_v}[\text{mon}] \\
& \quad \wedge \text{msg.uncommitted\_pn} \geq \text{uncommitted\_pn}[\text{mon}] \\
& \quad \text{THEN } \wedge \text{uncommitted\_v}' = \\
& \quad \quad [\text{uncommitted\_v} \text{ EXCEPT } ![\text{mon}] = \text{msg.last\_committed} + 1] \\
& \quad \wedge \text{uncommitted\_pn}' = \\
& \quad \quad [\text{uncommitted\_pn} \text{ EXCEPT } ![\text{mon}] = \text{msg.uncommitted\_pn}] \\
& \quad \wedge \text{uncommitted\_value}' = \\
& \quad \quad [\text{uncommitted\_value} \text{ EXCEPT } ![\text{mon}] = \text{msg.values}[\text{msg.last\_committed} + 1]] \\
& \quad \text{ELSE } \text{check\_and\_correct\_uncommitted}(\text{mon}) \\
& \wedge \text{UNCHANGED } \langle \text{phase}, \text{accepted\_pn} \rangle \\
\vee & \wedge \text{msg.pn} < \text{accepted\_pn}[\text{mon}] \\
& \wedge \text{check\_and\_correct\_uncommitted}(\text{mon}) \\
& \wedge \text{UNCHANGED } \langle \text{phase}, \text{accepted\_pn}, \text{num\_last} \rangle \\
& \wedge \text{UNCHANGED } \text{epoch} \\
& \wedge \text{UNCHANGED } \text{epoch} \\
& \wedge \text{UNCHANGED } \langle \text{quorum}, \text{quorum\_sz}, \text{isLeader}, \text{state}, \text{pending\_pn}, \text{pending\_v} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{lease\_vars}, \text{commit\_vars} \rangle
\end{aligned}$$

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

Variables changed: *peer\_first\_committed*, *peer\_last\_committed*, *state*, *accepted*, *new\_value*, *phase*, *messages*, *message\_history*, *values*, *pending\_pn*, *pending\_v*, *acked\_lease*.

@type: *MONITOR*  $\Rightarrow$  *Bool*;

*post\_last(mon)*  $\triangleq$

$$\begin{aligned}
& \wedge \text{isLeader}[\text{mon}] = \text{TRUE} \\
& \wedge \text{num\_last}[\text{mon}] = \text{quorum\_sz} \\
& \wedge \text{phase}[\text{mon}] = \text{PHASE\_COLLECT}
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{clear\_peer\_first\_committed}(\text{mon}) \\
& \wedge \text{clear\_peer\_last\_committed}(\text{mon})
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{IF } \wedge \text{uncommitted\_v}[\text{mon}] = \text{last\_committed}[\text{mon}] + 1 \\
& \quad \wedge \text{uncommitted\_value}[\text{mon}] \neq \text{Nil} \\
& \quad \text{THEN } \wedge \text{state}' = [\text{state} \text{ EXCEPT } ![\text{mon}] = \text{STATE\_UPDATING\_PREVIOUS}] \\
& \quad \wedge \text{begin}(\text{mon}, \text{uncommitted\_value}[\text{mon}]) \\
& \quad \wedge \text{UNCHANGED } \langle \text{acked\_lease}, \text{uncommitted\_v}, \text{uncommitted\_pn}, \text{uncommitted\_value} \rangle
\end{aligned}$$

ELSE  $\wedge$  *finish\_round*(*mon*)  
 $\wedge$  *extend\_lease*(*mon*)  
 $\wedge$  UNCHANGED  $\langle$  *accepted*, *new\_value*, *values*, *restart\_vars*  $\rangle$   
 $\wedge$  UNCHANGED  $\langle$  *isLeader*, *monitor\_store*, *accepted\_pn*, *first\_committed*, *last\_committed*  $\rangle$   
 $\wedge$  UNCHANGED  $\langle$  *epoch*, *quorum*, *quorum\_sz*, *num\_last*, *pending\_proposal*  $\rangle$

#### Leader election

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

Variables changed: *isLeader*, *state*, *phase*, *new\_value*, *pending\_proposal*, *epoch*.

@type: *Bool*;

*leader\_election*  $\triangleq$

$\wedge \exists \text{mon} \in \text{Monitors} :$   
 $\wedge \text{quorum}[\text{mon}]$   
 $\wedge \text{isLeader}' = [m \in \text{Monitors} \mapsto \text{IF } m = \text{mon} \text{ THEN TRUE ELSE FALSE}]$   
 $\wedge \text{state}' = [m \in \text{Monitors} \mapsto$   
 $\quad \text{IF } \text{quorum\_sz} = 1 \text{ THEN } \text{STATE\_ACTIVE} \text{ ELSE } \text{STATE\_RECOVERING}]$   
 $\wedge \text{phase}' = [m \in \text{Monitors} \mapsto \text{PHASE\_ELECTION}]$   
 $\wedge \text{new\_value}' = [m \in \text{Monitors} \mapsto \text{Nil}]$   
 $\wedge \text{pending\_proposal}' = [m \in \text{Monitors} \mapsto \text{Nil}]$   
 $\wedge \text{epoch}' = \text{epoch} + 1$   
 $\wedge \text{messages}' = [\text{mon1} \in \text{Monitors} \mapsto [\text{mon2} \in \text{Monitors} \mapsto \langle \rangle]]$   
 $\wedge$  UNCHANGED  $\langle$  *quorum*, *quorum\_sz*, *accepted*, *message\_history*  $\rangle$   
 $\wedge$  UNCHANGED  $\langle$  *data\_vars*, *restart\_vars*, *collect\_vars*, *lease\_vars*  $\rangle$

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

Variables changed: *accepted\_pn*, *phase*.

@type: *MONITOR*  $\Rightarrow$  *Bool*;

*election\_recover*(*mon*)  $\triangleq$

$\wedge \text{quorum\_sz} > 1$   
 $\wedge \text{phase}[\text{mon}] = \text{PHASE\_ELECTION}$   
 $\wedge \text{collect}(\text{mon}, 0)$   
 $\wedge$  UNCHANGED  $\langle$  *isLeader*, *state*, *values*, *first\_committed*, *last\_committed*, *monitor\_store*  $\rangle$   
 $\wedge$  UNCHANGED  $\langle$  *global\_vars*, *restart\_vars*, *collect\_vars*, *lease\_vars*, *commit\_vars*  $\rangle$

#### Timeouts and restart

Remove monitor from quorum, if there are enough monitors in the quorum.

@type: *MONITOR*  $\Rightarrow$  *Bool*;

*crash\_mon*(*mon*)  $\triangleq$

$\wedge \text{quorum\_sz} > (\text{MonitorsLen} \div 2) + 1$   
 $\wedge \text{quorum}[\text{mon}] = \text{TRUE}$   
 $\wedge \text{quorum}' = [\text{quorum} \text{ EXCEPT } ![\text{mon}] = \text{FALSE}]$   
 $\wedge \text{quorum\_sz}' = \text{quorum\_sz} - 1$

$\wedge$  *bootstrap*  
 $\wedge$  *number\_crashes'* = *number\_crashes* + 1  
 $\wedge$  UNCHANGED  $\langle$  *messages*, *message\_history*  $\rangle$   
 $\wedge$  UNCHANGED  $\langle$  *state\_vars*, *restart\_vars*, *data\_vars*, *collect\_vars*, *lease\_vars*, *commit\_vars*  $\rangle$

Add monitor to the quorum.

@type: *MONITOR*  $\Rightarrow$  *Bool*;

*restore\_mon(mon)*  $\triangleq$

$\wedge$  *quorum*[*mon*] = FALSE  
 $\wedge$  *quorum'* = [*quorum* EXCEPT ![*mon*] = TRUE]  
 $\wedge$  *quorum\_sz'* = *quorum\_sz* + 1  
 $\wedge$  *bootstrap*  
 $\wedge$  UNCHANGED  $\langle$  *messages*, *message\_history*  $\rangle$   
 $\wedge$  UNCHANGED  $\langle$  *state\_vars*, *restart\_vars*, *data\_vars*, *collect\_vars*, *lease\_vars*, *commit\_vars*  $\rangle$

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

Variables changed: epoch.

@type: *MONITOR*  $\Rightarrow$  *Bool*;

*Timeout(mon)*  $\triangleq$

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

#### Dispatchers and next statement

Handle a message.

@type: *MESSAGE*  $\Rightarrow$  *Bool*;

*Receive(msg)*  $\triangleq$

$\wedge \vee \wedge$  *msg.type* = *OP\_COLLECT*  
 $\wedge$  *handle\_collect(msg.dest, msg)*  
 $\wedge$  *step\_name'* = "receive collect"  
  
 $\vee \wedge$  *msg.type* = *OP\_LAST*  
 $\wedge$  *handle\_last(msg.dest, msg)*  
 $\wedge$  *step\_name'* = "receive last"  
  
 $\vee \wedge$  *msg.type* = *OP\_LEASE*  
 $\wedge$  *handle\_lease(msg.dest, msg)*  
 $\wedge$  *step\_name'* = "receive lease"  
  
 $\vee \wedge$  *msg.type* = *OP\_LEASE\_ACK*  
 $\wedge$  *handle\_lease\_ack(msg.dest, msg)*  
 $\wedge$  *step\_name'* = "receive lease\_ack"  
  
 $\vee \wedge$  *msg.type* = *OP\_BEGIN*  
 $\wedge$  *handle\_begin(msg.dest, msg)*

$$\begin{aligned}
& \wedge \text{step\_name}' = \text{"receive begin"} \\
& \vee \wedge \text{msg.type} = \text{OP\_ACCEPT} \\
& \quad \wedge \text{handle\_accept}(\text{msg.dest}, \text{msg}) \\
& \quad \wedge \text{step\_name}' = \text{"receive accept"} \\
& \vee \wedge \text{msg.type} = \text{OP\_COMMIT} \\
& \quad \wedge \text{handle\_commit}(\text{msg.dest}, \text{msg}) \\
& \quad \wedge \text{step\_name}' = \text{"receive commit"}
\end{aligned}$$

Limit some variables to reduce search space.

@type: Bool;

$$\begin{aligned}
\text{reduce\_search\_space} & \triangleq \\
& \wedge \text{epoch} \neq 8 \\
& \wedge \forall \text{mon} \in \text{Monitors} : \text{last\_committed}[\text{mon}] < 2 \\
& \quad \vee \forall \text{mon2} \in \text{Monitors} : \text{new\_value}[\text{mon2}] = \text{Nil} \\
& \wedge \forall \text{mon} \in \text{Monitors} : \text{accepted\_pn}[\text{mon}] < 300 \\
& \wedge \text{number\_crashes} \neq 4
\end{aligned}$$

State transitions.

@type: Bool;

$$\begin{aligned}
\text{Next} & \triangleq \\
& \wedge \text{reduce\_search\_space} \\
& \wedge \text{IF } \text{epoch} \% 2 = 1 \text{ THEN} \\
& \quad \wedge \text{leader\_election} \\
& \quad \wedge \text{step\_name}' = \text{"election"} \\
& \quad \wedge \text{UNCHANGED } \text{number\_crashes} \\
& \text{ELSE} \\
& \quad \vee \wedge \exists \text{mon} \in \text{Monitors} : \text{election\_recover}(\text{mon}) \\
& \quad \quad \wedge \text{step\_name}' = \text{"election\_recover"} \\
& \quad \quad \wedge \text{UNCHANGED } \text{number\_crashes} \\
& \quad \vee \wedge \exists \text{mon} \in \text{Monitors} : \text{send\_collect}(\text{mon}) \\
& \quad \quad \wedge \text{step\_name}' = \text{"send\_collect"} \\
& \quad \quad \wedge \text{UNCHANGED } \text{number\_crashes} \\
& \quad \vee \wedge \exists \text{mon} \in \text{Monitors} : \text{post\_last}(\text{mon}) \\
& \quad \quad \wedge \text{step\_name}' = \text{"post\_last"} \\
& \quad \quad \wedge \text{UNCHANGED } \text{number\_crashes} \\
& \quad \vee \wedge \exists \text{mon} \in \text{Monitors} : \text{post\_lease\_ack}(\text{mon}) \\
& \quad \quad \wedge \text{step\_name}' = \text{"post\_lease\_ack"} \\
& \quad \quad \wedge \text{UNCHANGED } \text{number\_crashes} \\
& \quad \vee \wedge \exists \text{mon} \in \text{Monitors} : \text{post\_accept}(\text{mon}) \\
& \quad \quad \wedge \text{step\_name}' = \text{"post\_accept"} \\
& \quad \quad \wedge \text{UNCHANGED } \text{number\_crashes}
\end{aligned}$$

$$\begin{aligned}
& \vee \wedge \exists mon \in Monitors : finish\_commit(mon) \\
& \quad \wedge step\_name' = \text{"finish\_commit"} \\
& \quad \wedge UNCHANGED\ number\_crashes \\
& \vee \wedge \exists mon \in Monitors : \exists v \in Value\_set : client\_request(mon, v) \\
& \quad \wedge step\_name' = \text{"client\_request"} \\
& \quad \wedge UNCHANGED\ number\_crashes \\
& \vee \wedge \exists mon \in Monitors : propose\_pending(mon) \\
& \quad \wedge step\_name' = \text{"propose\_pending"} \\
& \quad \wedge UNCHANGED\ number\_crashes \\
& \vee \wedge \exists mon1, mon2 \in Monitors : \\
& \quad \wedge mon1 \neq mon2 \\
& \quad \wedge Len(messages[mon1][mon2]) > 0 \\
& \quad \wedge Receive(messages[mon1][mon2][1]) \\
& \quad \wedge UNCHANGED\ number\_crashes \\
& \vee \wedge \exists mon \in Monitors : crash\_mon(mon) \\
& \quad \wedge step\_name' = \text{"crash\_mon"} \\
& \quad \wedge UNCHANGED\ number\_crashes \\
& \vee \wedge \exists mon \in Monitors : restore\_mon(mon) \\
& \quad \wedge step\_name' = \text{"restore\_mon"} \\
& \quad \wedge UNCHANGED\ number\_crashes \\
& \vee \wedge \exists mon \in Monitors : Timeout(mon) \\
& \quad \wedge step\_name' = \text{"timeout\_and\_restart"} \\
& \quad \wedge UNCHANGED\ number\_crashes
\end{aligned}$$

#### Safety invariants

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

@type: Bool;

$$\begin{aligned}
same\_monitor\_store &\triangleq \forall mon1, mon2 \in Monitors : \\
& \quad state[mon1] = STATE\_ACTIVE \wedge state[mon2] = STATE\_ACTIVE \\
& \quad \Rightarrow monitor\_store[mon1] = monitor\_store[mon2]
\end{aligned}$$

Invariant.

@type: Bool;

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

#### Test/Debug invariants

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

$$Inv\_find\_state(x) \triangleq \neg x$$

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

$TLCGet("level")$  not supported by snowcat typechecker.

$Inv\_diam(size) \triangleq TLCGet("level") \neq size - 1$

Invariants to test in model check

$DEBUG\_Inv \triangleq \wedge TRUE$   
 $\wedge Inv\_diam(20)$

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.

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

Find a state where a monitor crashed during the collect phase and fails to send a  $OP\_LAST$  message.

$Inv\_find\_state($   
   $\wedge step\_name = "crash\ mon"$   
  
   $\backslash * The\ system\ is\ in\ collect\ phase\ and\ no\ OP\_LAST\ message\ has\ been\ received.$   
   $\backslash * isLeader[mon] = TRUE\ assures\ that\ the\ leader\ was\ not\ the\ one\ that\ crashed.$   
   $\wedge \exists mon \in Monitors :$   
     $\wedge isLeader[mon] = TRUE$   
     $\wedge phase[mon] = PHASE\_COLLECT$   
     $\wedge num\_last[mon] = 1$   
  
   $\backslash * All\ the\ collect\ requests\ have\ been\ handled\ by\ the\ peers.$   
   $\wedge \forall mon1, mon2 \in Monitors :$   
     $\forall i \in 1 \dots Len(messages[mon1][mon2]) : messages[mon1][mon2][i].type \neq OP\_COLLECT$   
  
   $\wedge epoch = 2$   
 $)$

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

$Inv\_find\_state($   
   $\wedge step\_name = "crash\ mon"$   
   $\wedge \exists mon1, mon2 \in Monitors :$   
     $\exists i \in 1 \dots Len(messages[mon1][mon2]) : messages[mon1][mon2][i].type = OP\_ACCEPT$   
   $\wedge \forall mon \in Monitors :$   
     $isLeader[mon] = FALSE$   
   $\wedge epoch = 2$   
 $)$

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

$\backslash * Modification\ History$   
 $\backslash * Last\ modified\ Sun\ Apr\ 18\ 22:54:31\ WEST\ 2021\ by\ afonsonf$   
 $\backslash * Created\ Mon\ Jan\ 11\ 16:15:26\ WET\ 2021\ by\ afonsonf$