

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

**Constants**

Set of *Monitors*.

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

Sequence of monitors.

@type: *Seq(MONITOR)*;

*MonitorsSeq*  $\triangleq \text{TLCEval}(\text{SetToSeq}(\text{Monitors}))$

Number of monitors.

@type: *Int*;

*MonitorsLen*  $\triangleq \text{TLCEval}(\text{Len}(\text{MonitorsSeq}))$

Rank predicate, used as a monitor identifier.

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

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

Set of possible values.

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

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

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

Reserved value.  
 CONSTANTS @type: VALUE; Nil

Paxos states.  
 CONSTANTS @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.  
 CONSTANTS @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.  
 CONSTANTS @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. When epoch is odd trigger an election.  
 VARIABLE @type: Int; epoch

Store messages waiting to be handled.  
 VARIABLE @type: MONITOR → (MONITOR → Seq(MESSAGE)); messages

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

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

Size of the current quorum.  
 VARIABLE @type: Int; quorum\_sz

## State variables

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

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 specification, 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*

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 initialization

@typeAlias: *VALUE\_VERSION* = *Int*;

@typeAlias: *PN* = *Int*;

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

*state\_vars*  $\triangleq$   $\langle \text{isLeader}, \text{state}, \text{phase} \rangle$

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

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

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

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

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

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

*Init\\_global\\_vars*  $\triangleq$

$\wedge \text{epoch} = 1$

$\wedge \text{messages} = [\text{mon1} \in \text{Monitors} \mapsto [\text{mon2} \in \text{Monitors} \mapsto \langle \rangle]]$

$$\begin{aligned}
& \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"}
\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  $\exists \text{ elem} \in S : \text{TRUE}$  THEN LET *elem*  $\triangleq$  CHOOSE  $x \in S : \text{TRUE}$   
 IN  $\langle \text{elem} \rangle$   
 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*)  $\triangleq$   
 ∧ *messages*' = WithMessage(*m*, *messages*)  
 ∧ *message\_history*' = *message\_history* ∪ {*m*}  
 ∧ 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*)  $\triangleq$   
 ∧ *messages*' = [*messages* EXCEPT ![*from*] =  
 [*mon* ∈ *Monitors* ↦  
*messages*[*from*][*mon*] ∘ SingleMessageSetToSeq({*m* ∈ *m\_set* : *m.dest* = *mon*})]  
 ∧ *message\_history*' = *message\_history* ∪ *m\_set*  
 ∧ UNCHANGED *message\_history*

Removes the request from network and adds the response.

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

@type: (MESSAGE, MESSAGE) ⇒ 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:  $oldpn = 302, epoch = 4, newpn = 404$ .

@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:  $MONITOR \Rightarrow Bool$ ;

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

Clear the variable  $peer\_last\_committed$ .

Variables changed:  $peer\_last\_committed$ .

$\text{@type: } \text{MONITOR} \Rightarrow \text{Bool};$   
 $\text{clear\_peer\_last\_committed}(\text{mon}) \triangleq$   
 $\text{peer\_last\_committed}' = [\text{peer\_last\_committed} \text{ EXCEPT } ![mon] = [m \in \text{Monitors} \mapsto -1]]$

Store peer values and update *first\_committed*, *last\_committed* and *monitor\_store* accordingly.  
 Variables changed: *values*, *first\_committed*, *last\_committed*, *monitor\_store*.  
 $\text{@type: (MONITOR, MESSAGE)} \Rightarrow \text{Bool};$   
 $\text{store\_state}(\text{mon}, \text{msg}) \triangleq$

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

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

Update monitor store.  
 $\wedge \text{IF } \text{last\_committed}'[mon] = 0$   
 THEN UNCHANGED *monitor\_store*  
 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*.  
 $\text{@type: } \text{MONITOR} \Rightarrow \text{Bool};$   
 $\text{check\_and\_correct\_uncommitted}(\text{mon}) \triangleq$   
 IF  $\text{uncommitted\_v}[mon] \leq \text{last\_committed}'[mon]$   
 THEN  $\wedge \text{uncommitted\_v}' = [\text{uncommitted\_v} \text{ EXCEPT } ![mon] = 0]$   
 $\wedge \text{uncommitted\_pn}' = [\text{uncommitted\_pn} \text{ EXCEPT } ![mon] = 0]$   
 $\wedge \text{uncommitted\_value}' = [\text{uncommitted\_value} \text{ EXCEPT } ![mon] = \text{Nil}]$   
 ELSE UNCHANGED  $\langle \text{uncommitted\_pn}, \text{uncommitted\_v}, \text{uncommitted\_value} \rangle$

Trigger new election by incrementing epoch.  
 Variables changed: *epoch*.  
 $\text{@type: } \text{Bool};$   
 $\text{bootstrap} \triangleq$   
 $\wedge \text{epoch}' = \text{epoch} + 1$

#### Lease phase predicates



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

Variables changed: *state*.

```
@type: MONITOR  $\Rightarrow$  Bool;
finish_round(mon)  $\triangleq$ 
   $\wedge$  isLeader[mon] = TRUE
   $\wedge$  state' = [state EXCEPT ![mon] = 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] = TRUE
   $\wedge$  acked_lease' = [acked_lease EXCEPT ![mon] =
    [m  $\in$  Monitors  $\mapsto$  IF m = mon THEN TRUE ELSE FALSE]]
   $\wedge$  Send_set(mon,
    {[type  $\mapsto$  OP_LEASE,
      from  $\mapsto$  mon,
      dest  $\mapsto$  dest,
      last_committed  $\mapsto$  last_committed[mon]] : dest  $\in$  {m  $\in$  Monitors \ {mon} : quorum[m]}}
    })
   $\wedge$  phase' = [phase EXCEPT ![mon] = PHASE_LEASE]
```

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

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

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

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

```
@type: (MONITOR, MESSAGE)  $\Rightarrow$  Bool;
handle_lease(mon, msg)  $\triangleq$ 
   $\wedge$  discard if not peon or peon is behind
  IF  $\vee$  isLeader[mon] = TRUE
     $\vee$  last_committed[mon]  $\neq$  msg.last_committed
  THEN  $\wedge$  Discard(msg)
     $\wedge$  UNCHANGED state
  ELSE  $\wedge$  state' = [state EXCEPT ![mon] = STATE_ACTIVE]
     $\wedge$  Reply([type  $\mapsto$  OP_LEASE_ACK,
      from  $\mapsto$  mon,
      dest  $\mapsto$  msg.from,
      first_committed  $\mapsto$  first_committed[mon],
      last_committed  $\mapsto$  last_committed[mon]], msg)
     $\wedge$  Discard(msg)
   $\wedge$  UNCHANGED  $\langle$ epoch, quorum, quorum_sz, isLeader, phase $\rangle$ 
   $\wedge$  UNCHANGED  $\langle$ restart_vars, data_vars, collect_vars, lease_vars, commit_vars $\rangle$ 
```

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

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

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

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

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

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

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;

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

### 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;

$begin(mon, v) \triangleq$   
 $\wedge isLeader[mon] = TRUE$   
 $\wedge \vee state'[mon] = STATE\_UPDATING$   
 $\quad \vee state'[mon] = STATE\_UPDATING\_PREVIOUS$   
 $\wedge quorum\_sz = 1 \vee num\_last[mon] > MonitorsLen \div 2$   
 $\wedge new\_value[mon] = Nil$   
 $\wedge accepted' = [accepted \text{ EXCEPT } ![mon] =$   
 $\quad [m \in Monitors \mapsto \text{IF } m = mon \text{ THEN TRUE ELSE FALSE}]]$   
 $\wedge new\_value' = [new\_value \text{ EXCEPT } ![mon] = v]$   
 $\wedge phase' = [phase \text{ EXCEPT } ![mon] = PHASE\_BEGIN]$   
 $\wedge values' = [values \text{ EXCEPT } ![mon] =$   
 $\quad ((last\_committed[mon] + 1) :> new\_value'[mon]) @@ values[mon]]$   
 $\wedge Send\_set(mon,$   
 $\quad \{[type \quad \mapsto OP\_BEGIN,$   
 $\quad \quad from \quad \mapsto mon,$

$$\begin{aligned}
& dest \quad \mapsto dest, \\
& last\_committed \mapsto last\_committed[mon], \\
& values \quad \mapsto values'[mon], \\
& pn \quad \mapsto accepted\_pn[mon] : dest \in \{m \in Monitors \setminus \{mon\} : quorum[m]\} \\
& }) \\
& \wedge pending\_pn' = [pending\_pn \text{ EXCEPT } ![mon] = accepted\_pn[mon]] \\
& \wedge pending\_v' = [pending\_v \text{ EXCEPT } ![mon] = last\_committed[mon] + 1]
\end{aligned}$$

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;

$$\begin{aligned}
& handle\_begin(mon, msg) \triangleq \\
& \quad \wedge isLeader[mon] = FALSE \\
& \quad \wedge IF \ msg.pn < accepted\_pn[mon] \\
& \quad \quad THEN \\
& \quad \quad \quad \wedge Discard(msg) \\
& \quad \quad \quad \wedge UNCHANGED \langle state, values, pending\_pn, pending\_v \rangle \\
& \quad \quad ELSE \\
& \quad \quad \quad \wedge msg.pn = accepted\_pn[mon] \\
& \quad \quad \quad \wedge msg.last\_committed = last\_committed[mon] \\
& \quad \quad \quad assign values[mon][last\_committed[mon] + 1] \\
& \quad \quad \quad \wedge values' = [values \text{ EXCEPT } ![mon] = \\
& \quad \quad \quad \quad ((last\_committed[mon] + 1) :> msg.values[last\_committed[mon] + 1]) @@ values[mon]] \\
& \quad \quad \quad \wedge state' = [state \text{ EXCEPT } ![mon] = STATE\_UPDATING] \\
& \quad \quad \quad \wedge pending\_pn' = [pending\_pn \text{ EXCEPT } ![mon] = accepted\_pn[mon]] \\
& \quad \quad \quad \wedge pending\_v' = [pending\_v \text{ EXCEPT } ![mon] = last\_committed[mon] + 1] \\
& \quad \quad \quad \wedge Reply([type \quad \mapsto OP\_ACCEPT, \\
& \quad \quad \quad \quad from \quad \mapsto mon, \\
& \quad \quad \quad \quad dest \quad \mapsto msg.from, \\
& \quad \quad \quad \quad last\_committed \mapsto last\_committed[mon], \\
& \quad \quad \quad \quad pn \quad \mapsto accepted\_pn[mon]], msg) \\
& \quad \quad \quad \wedge UNCHANGED \langle epoch, quorum, quorum\_sz, isLeader, phase, monitor\_store, \\
& \quad \quad \quad \quad \quad accepted\_pn, first\_committed, last\_committed, uncommitted\_pn, \\
& \quad \quad \quad \quad \quad uncommitted\_v, uncommitted\_value \rangle \\
& \quad \quad \quad \wedge UNCHANGED \langle collect\_vars, lease\_vars, commit\_vars \rangle
\end{aligned}$$

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

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

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

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

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

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*;

$post\_accept(mon) \triangleq$

$\wedge phase[mon] = \text{PHASE\_BEGIN}$   
 $\wedge \forall m \in \text{Monitors} : quorum[m] \Rightarrow accepted[mon][m] = \text{TRUE}$   
 $\wedge new\_value[mon] \neq \text{Nil}$   
 $\wedge \vee state[mon] = \text{STATE\_UPDATING\_PREVIOUS}$   
 $\vee state[mon] = \text{STATE\_UPDATING}$   
 $\wedge last\_committed' = [last\_committed \text{ EXCEPT } ![mon] = last\_committed[mon] + 1]$   
 $\wedge \text{IF } first\_committed[mon] = 0$   
 $\text{THEN } first\_committed' = [first\_committed \text{ EXCEPT } ![mon] = first\_committed[mon] + 1]$   
 $\text{ELSE UNCHANGED } first\_committed$   
 $\wedge monitor\_store' = [monitor\_store \text{ EXCEPT } ![mon] = values[mon][last\_committed[mon] + 1]]$   
 $\wedge new\_value' = [new\_value \text{ EXCEPT } ![mon] = \text{Nil}]$   
 $\wedge Send\_set(mon,$   
 $\quad \{[type \quad \mapsto OP\_COMMIT,$   
 $\quad \quad from \quad \mapsto mon,$   
 $\quad \quad dest \quad \mapsto dest,$   
 $\quad \quad last\_committed \mapsto last\_committed'[mon],$   
 $\quad \quad pn \quad \mapsto accepted\_pn[mon],$   
 $\quad \quad values \quad \mapsto values[mon]] : dest \in \{m \in \text{Monitors} \setminus \{mon\} : quorum[m]\}$   
 $\quad \})$   
 $\wedge state' = [state \text{ EXCEPT } ![mon] = \text{STATE\_REFRESH}]$   
 $\wedge phase' = [phase \text{ EXCEPT } ![mon] = \text{PHASE\_COMMIT}]$   
 $\wedge \text{UNCHANGED } \langle isLeader, values, accepted\_pn, pending\_proposal, accepted \rangle$   
 $\wedge \text{UNCHANGED } \langle epoch, quorum, quorum\_sz, restart\_vars, collect\_vars, lease\_vars \rangle$

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

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

@type: *MONITOR*  $\Rightarrow$  *Bool*;  
 $finish\_commit(mon) \triangleq$   
 $\wedge state[mon] = STATE\_REFRESH$   
 $\wedge phase[mon] = PHASE\_COMMIT$   
 $\wedge finish\_round(mon)$   
 $\wedge extend\_lease(mon)$   
 $\wedge UNCHANGED \langle epoch, quorum, quorum\_sz, isLeader \rangle$   
 $\wedge UNCHANGED \langle restart\_vars, data\_vars, collect\_vars, commit\_vars \rangle$

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

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

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*;  
 $propose\_pending(mon) \triangleq$   
 $\wedge phase[mon] = PHASE\_LEASE \vee phase[mon] = PHASE\_ELECTION$   
 $\wedge state[mon] = STATE\_ACTIVE$

$\wedge \text{pending\_proposal}[mon] \neq Nil$   
 $\wedge \text{pending\_proposal}' = [\text{pending\_proposal} \text{ EXCEPT } ![mon] = Nil]$   
 $\wedge \text{state}' = [\text{state} \text{ EXCEPT } ![mon] = STATE\_UPDATING]$   
 $\wedge \text{begin}(mon, \text{pending\_proposal}[mon])$   
 $\wedge \text{UNCHANGED } \langle isLeader, monitor\_store, accepted\_pn, first\_committed, last\_committed,$   
 $\quad epoch, quorum, quorum\_sz, uncommitted\_v, uncommitted\_pn, uncommitted\_value \rangle$   
 $\wedge \text{UNCHANGED } \langle collect\_vars, lease\_vars \rangle$

### 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;

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

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;

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

```

    {[type      ↦ OP_COLLECT,
     from      ↦ mon,
     dest      ↦ dest,
     first_committed ↦ first_committed[mon],
     last_committed ↦ last_committed[mon],
     pn        ↦ accepted_pn[mon]] : dest ∈ {m ∈ Monitors \ {mon} : quorum[m]}
    })
  ∧ phase' = [phase EXCEPT ![mon] = PHASE_COLLECT]
  ∧ UNCHANGED ⟨isLeader, state⟩
  ∧ UNCHANGED ⟨epoch, quorum, quorum_sz, data_vars, lease_vars, commit_vars⟩

```

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) ⇒ Bool;

```

handle_collect(mon, msg)  $\triangleq$ 
  ∧ isLeader[mon] = FALSE
  ∧ state' = [state EXCEPT ![mon] = STATE_RECOVERING]
  ∧ ∨ ∧ msg.first_committed > last_committed[mon] + 1
    ∧ bootstrap
    ∧ Discard(msg)
    ∧ UNCHANGED ⟨accepted_pn⟩
  ∨ ∧ msg.first_committed ≤ last_committed[mon] + 1
    ∧ IF msg.pn > accepted_pn[mon]
      THEN accepted_pn' = [accepted_pn EXCEPT ![mon] = msg.pn]
      ELSE UNCHANGED accepted_pn
    ∧ Reply([type      ↦ OP_LAST,
             from      ↦ mon,
             dest      ↦ msg.from,
             first_committed ↦ first_committed[mon],
             last_committed ↦ last_committed[mon],
             values      ↦ values[mon],
             uncommitted_pn ↦ pending_pn[mon],
             pn          ↦ accepted_pn'[mon]], msg)
    ∧ UNCHANGED epoch
  ∧ UNCHANGED ⟨isLeader, phase, values, first_committed, last_committed, monitor_store⟩
  ∧ UNCHANGED ⟨quorum, quorum_sz, restart_vars, collect_vars, lease_vars, commit_vars⟩

```

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, values, accepted\_pn, first\_committed, last\_committed, num\_last, peer\_first\_committed, peer\_last\_committed.

@type: (MONITOR, MESSAGE) ⇒ Bool;

$$\begin{aligned}
& \text{handle\_last}(\text{mon}, \text{msg}) \triangleq \\
& \quad \wedge \text{isLeader}[\text{mon}] = \text{TRUE} \\
& \quad \wedge \text{peer\_first\_committed}' = [\text{peer\_first\_committed} \text{ EXCEPT } ![\text{mon}] = \\
& \quad \quad [\text{peer\_first\_committed}[\text{mon}] \text{ EXCEPT } ![\text{msg.from}] = \text{msg.first\_committed}]] \\
& \quad \wedge \text{peer\_last\_committed}' = [\text{peer\_last\_committed} \text{ EXCEPT } ![\text{mon}] = \\
& \quad \quad [\text{peer\_last\_committed}[\text{mon}] \text{ EXCEPT } ![\text{msg.from}] = \text{msg.last\_committed}]] \\
& \quad \wedge \text{IF } \text{msg.first\_committed} > \text{last\_committed}[\text{mon}] + 1 \\
& \quad \quad \text{THEN} \\
& \quad \quad \quad \wedge \text{bootstrap} \\
& \quad \quad \quad \wedge \text{Discard}(\text{msg}) \\
& \quad \quad \quad \wedge \text{UNCHANGED } \langle \text{first\_committed}, \text{last\_committed}, \text{num\_last}, \text{accepted\_pn}, \text{values}, \text{phase} \rangle \\
& \quad \quad \quad \wedge \text{UNCHANGED } \langle \text{monitor\_store}, \text{uncommitted\_pn}, \text{uncommitted\_v}, \text{uncommitted\_value} \rangle \\
& \quad \quad \text{ELSE} \\
& \quad \quad \quad \wedge \text{store\_state}(\text{mon}, \text{msg}) \\
& \quad \quad \quad \wedge \text{IF } \exists \text{peer} \in \text{Monitors} : \\
& \quad \quad \quad \quad \wedge \text{peer} \neq \text{mon} \\
& \quad \quad \quad \quad \wedge \text{peer\_last\_committed}'[\text{mon}][\text{peer}] \neq -1 \\
& \quad \quad \quad \quad \wedge \text{peer\_last\_committed}'[\text{mon}][\text{peer}] + 1 < \text{first\_committed}[\text{mon}] \\
& \quad \quad \quad \quad \wedge \text{first\_committed}[\text{mon}] > 1 \\
& \quad \quad \quad \text{THEN} \\
& \quad \quad \quad \quad \wedge \text{bootstrap} \\
& \quad \quad \quad \quad \wedge \text{check\_and\_correct\_uncommitted}(\text{mon}) \\
& \quad \quad \quad \quad \wedge \text{Discard}(\text{msg}) \\
& \quad \quad \quad \quad \wedge \text{UNCHANGED } \langle \text{phase}, \text{accepted\_pn}, \text{num\_last} \rangle \\
& \quad \quad \quad \text{ELSE} \\
& \quad \quad \quad \quad \text{Send commit messages with my values to peers that are behind.} \\
& \quad \quad \quad \text{LET } \text{monitors\_behind} \triangleq \{ \text{peer} \in \text{Monitors} : \\
& \quad \quad \quad \quad \wedge \text{peer} \neq \text{mon} \\
& \quad \quad \quad \quad \wedge \text{peer\_last\_committed}'[\text{mon}][\text{peer}] \neq -1 \\
& \quad \quad \quad \quad \wedge \text{peer\_last\_committed}'[\text{mon}][\text{peer}] < \text{last\_committed}[\text{mon}] \\
& \quad \quad \quad \quad \wedge \text{quorum}[\text{peer}] \} \\
& \quad \quad \quad \text{IN } \text{Reply\_set}(\text{mon}, \\
& \quad \quad \quad \quad \{ [\text{type} \quad \quad \quad \mapsto \text{OP\_COMMIT}, \\
& \quad \quad \quad \quad \text{from} \quad \quad \quad \mapsto \text{mon}, \\
& \quad \quad \quad \quad \text{dest} \quad \quad \quad \mapsto \text{dest}, \\
& \quad \quad \quad \quad \text{last\_committed} \mapsto \text{last\_committed}'[\text{mon}], \\
& \quad \quad \quad \quad \text{pn} \quad \quad \quad \mapsto \text{accepted\_pn}[\text{mon}], \\
& \quad \quad \quad \quad \text{values} \quad \quad \mapsto \text{values}[\text{mon}]] : \text{dest} \in \text{monitors\_behind} \\
& \quad \quad \quad \quad \}, \text{msg}) \\
& \quad \quad \wedge \vee \text{Start new collect phase with a bigger proposal number.} \\
& \quad \quad \quad \wedge \text{msg.pn} > \text{accepted\_pn}[\text{mon}] \\
& \quad \quad \quad \wedge \text{collect}(\text{mon}, \text{msg.pn}) \\
& \quad \quad \quad \wedge \text{check\_and\_correct\_uncommitted}(\text{mon})
\end{aligned}$$



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

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

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

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

$post\_last(mon) \triangleq$

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

$\wedge num\_last[mon] = quorum\_sz$

$\wedge phase[mon] = \text{PHASE\_COLLECT}$

$\wedge clear\_peer\_first\_committed(mon)$

$\wedge clear\_peer\_last\_committed(mon)$

$\wedge$  IF  $\wedge uncommitted\_v[mon] = last\_committed[mon] + 1$

$\wedge uncommitted\_value[mon] \neq Nil$

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

$\wedge begin(mon, uncommitted\_value[mon])$

$\wedge$  UNCHANGED  $\langle acked\_lease, uncommitted\_v, uncommitted\_pn, uncommitted\_value \rangle$

ELSE  $\wedge finish\_round(mon)$

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

### 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*  $\triangleq$

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

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

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

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

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

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

### Timeouts and restart

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

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

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

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

$\wedge$  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*;

$restore\_mon(mon) \triangleq$

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

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

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

$\wedge bootstrap$

$\wedge$  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*;

$Receive(msg) \triangleq$

$\wedge \vee \wedge msg.type = OP\_COLLECT$   
 $\wedge handle\_collect(msg.dest, msg)$   
 $\wedge step\_name' = \text{"receive collect"}$

$\vee \wedge msg.type = OP\_LAST$   
 $\wedge handle\_last(msg.dest, msg)$   
 $\wedge step\_name' = \text{"receive last"}$

$\vee \wedge msg.type = OP\_LEASE$   
 $\wedge handle\_lease(msg.dest, msg)$   
 $\wedge step\_name' = \text{"receive lease"}$

$\vee \wedge msg.type = OP\_LEASE\_ACK$   
 $\wedge handle\_lease\_ack(msg.dest, msg)$   
 $\wedge step\_name' = \text{"receive lease\_ack"}$

$\vee \wedge msg.type = OP\_BEGIN$   
 $\wedge handle\_begin(msg.dest, msg)$   
 $\wedge step\_name' = \text{"receive begin"}$

$\vee \wedge msg.type = OP\_ACCEPT$

$\wedge \text{handle\_accept}(\text{msg.dest}, \text{msg})$   
 $\wedge \text{step\_name}' = \text{"receive accept"}$

$\vee \wedge \text{msg.type} = \text{OP\_COMMIT}$   
 $\wedge \text{handle\_commit}(\text{msg.dest}, \text{msg})$   
 $\wedge \text{step\_name}' = \text{"receive commit"}$

Limit some variables to reduce search space.

@type: Bool;

$\text{reduce\_search\_space} \triangleq$   
 $\wedge \text{epoch} \neq 14$   
 $\wedge \forall \text{mon} \in \text{Monitors} : \text{last\_committed}[\text{mon}] < 2$   
 $\wedge \forall \text{mon} \in \text{Monitors} : \text{accepted\_pn}[\text{mon}] < 400$

State transitions.

@type: Bool;

$\text{Next} \triangleq$   
 $\wedge \text{reduce\_search\_space}$   
 $\wedge \text{IF } \text{epoch} \% 2 = 1 \text{ THEN}$   
 $\quad \wedge \text{leader\_election}$   
 $\quad \wedge \text{step\_name}' = \text{"election"}$   
 $\text{ELSE}$   
 $\quad \vee \wedge \exists \text{mon} \in \text{Monitors} : \text{election\_recover}(\text{mon})$   
 $\quad \quad \wedge \text{step\_name}' = \text{"election\_recover"}$   
 $\quad \vee \wedge \exists \text{mon} \in \text{Monitors} : \text{send\_collect}(\text{mon})$   
 $\quad \quad \wedge \text{step\_name}' = \text{"send\_collect"}$   
 $\quad \vee \wedge \exists \text{mon} \in \text{Monitors} : \text{post\_last}(\text{mon})$   
 $\quad \quad \wedge \text{step\_name}' = \text{"post\_last"}$   
 $\quad \vee \wedge \exists \text{mon} \in \text{Monitors} : \text{post\_lease\_ack}(\text{mon})$   
 $\quad \quad \wedge \text{step\_name}' = \text{"post\_lease\_ack"}$   
 $\quad \vee \wedge \exists \text{mon} \in \text{Monitors} : \text{post\_accept}(\text{mon})$   
 $\quad \quad \wedge \text{step\_name}' = \text{"post\_accept"}$   
 $\quad \vee \wedge \exists \text{mon} \in \text{Monitors} : \text{finish\_commit}(\text{mon})$   
 $\quad \quad \wedge \text{step\_name}' = \text{"finish\_commit"}$   
 $\quad \vee \wedge \exists \text{mon} \in \text{Monitors} : \exists v \in \text{Value\_set} : \text{client\_request}(\text{mon}, v)$   
 $\quad \quad \wedge \text{step\_name}' = \text{"client\_request"}$   
 $\quad \vee \wedge \exists \text{mon} \in \text{Monitors} : \text{propose\_pending}(\text{mon})$   
 $\quad \quad \wedge \text{step\_name}' = \text{"propose\_pending"}$   
 $\quad \vee \wedge \exists \text{mon1}, \text{mon2} \in \text{Monitors} :$   
 $\quad \quad \quad \wedge \text{mon1} \neq \text{mon2}$

$$\begin{aligned}
& \wedge \text{Len}(\text{messages}[\text{mon1}][\text{mon2}]) > 0 \\
& \wedge \text{Receive}(\text{messages}[\text{mon1}][\text{mon2}][1]) \\
\vee \wedge \exists \text{mon} \in \text{Monitors} : & \text{crash\_mon}(\text{mon}) \\
& \wedge \text{step\_name}' = \text{"crash\_mon"} \\
\vee \wedge \exists \text{mon} \in \text{Monitors} : & \text{restore\_mon}(\text{mon}) \\
& \wedge \text{step\_name}' = \text{"restore\_mon"} \\
\vee \wedge \exists \text{mon} \in \text{Monitors} : & \text{Timeout}(\text{mon}) \\
& \wedge \text{step\_name}' = \text{"timeout\_and\_restart"}
\end{aligned}$$

### Safety invariants

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

@type: Bool;

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

In each version, every monitor chooses the same value.

@type: Bool;

$$\begin{aligned}
\text{same\_monitor\_values} & \triangleq \\
& \forall \text{version} \in 1 \dots \text{Max}(\{\text{last\_committed}[\text{mon}] : \text{mon} \in \text{Monitors}\}) : \\
& \quad \exists \text{val} \in \text{Value\_set} : \\
& \quad \forall \text{mon} \in \text{Monitors} : \\
& \quad \quad \text{last\_committed}[\text{mon}] < \text{version} \vee \text{values}[\text{mon}][\text{version}] = \text{val}
\end{aligned}$$

Invariant.

@type: Bool;

$$\begin{aligned}
\text{Inv} & \triangleq \wedge \text{same\_monitor\_store} \\
& \wedge \text{same\_monitor\_values}
\end{aligned}$$

### Test/Debug invariants

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

$$\text{Inv\_find\_state}(x) \triangleq \neg x$$

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

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

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

Invariants to test in model check

$$\begin{aligned}
\text{DEBUG\_Inv} & \triangleq \wedge \text{TRUE} \\
& \wedge \text{Inv\_diam}(20)
\end{aligned}$$

Examples:

Find a behavior with a diameter of size 60.

*Inv\_diam*(60)

Find a behavior where two different monitors assume the role of a leader (in different epochs).  
Note: In the message manipulation predicates, the variable *message\_history* needs to be updated.

```
Inv_find_state(
  ∃ msg1, msg2 ∈ message_history :
    ∧ msg1.type = OP_COLLECT ∧ msg2.type = OP_COLLECT
    ∧ msg1.from ≠ msg2.from
)
```

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

```
Inv_find_state(
  ∧ step_name = "crash_mon"

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

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

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

```
Inv_find_state(
  ∧ step_name = "crash_mon"
  ∧ ∃ mon1, mon2 ∈ Monitors :
    ∃ i ∈ 1 .. Len(messages[mon1][mon2]) : messages[mon1][mon2][i].type = OP_ACCEPT
  ∧ ∀ mon ∈ Monitors :
    isLeader[mon] = FALSE ∨ quorum[mon] = FALSE
)
```

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

```
\ * Modification History
\ * Last modified Thu Jun 10 18:34:59 WEST 2021 by afonsonf
\ * Created Mon Jan 11 16:15:26 WET 2021 by afonsonf
```