

MODULE *paxos*

This is a specification of the paxos algorithm implemented in Ceph. The specification is based on the following source file: <https://github.com/ceph/ceph/blob/master/src/mon/Paxos.cc>

The main mechanism abstracted that may differ from the version implemented in Ceph are:

- The election logic. The leader is chosen randomly, and, for now, only one leader is chosen per epoch. When a new epoch begins, the messages from the previous epoch are discarded.
- Monitor quorum. The quorum is defined in the election phase, using all monitors that are up. Different epochs can have different quorums.
- The communication layer. The variable `messages` represents connections between monitors (e.g. `messages[mon1][mon2]` holds the messages sent from `mon1` to `mon2`). Within a connection the messages are sent and received in order.
- The transactions. Transactions are simplified to represent only a change of a value in the variable `monitor_store`.
- Failure model. A monitor can crash if the remaining number of monitors is sufficient to form a quorum. When a monitor crashes, new elections are triggered and the monitor is marked to not be part of a quorum until he recovers.
- Timeouts. A timeout can occur at any point in the algorithm and it will trigger new elections.

For a more detailed overview of the specification: <https://github.com/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  $x \triangleq \text{CHOOSE } x \in S : \text{TRUE}$   
 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.  
 @typeAlias: MONITOR = T;  
 @typeAlias: VALUE = T;  
 $SYMM \triangleq Permutations(Monitors) \cup Permutations(Value\_set)$

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. If is odd trigger an election.  
 VARIABLE @type: Int; epoch

Store messages waiting to be handled.

VARIABLE @type:  $MONITOR \rightarrow (MONITOR \rightarrow Seq(MESSAGE)); messages$

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:  $MONITOR \rightarrow Bool; quorum$

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 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       $\triangleq$   $\langle \text{epoch}, \text{messages}, \text{message\_history}, \text{quorum}, \text{quorum\_sz} \rangle$ 
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        $\triangleq$   $\langle \text{acked\_lease} \rangle$ 
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  $\triangleq$ 
   $\wedge \text{epoch} = 1$ 
   $\wedge \text{messages} = [\text{mon1} \in \text{Monitors} \mapsto [\text{mon2} \in \text{Monitors} \mapsto \langle \rangle]]$ 
   $\wedge \text{message\_history} = \{\}$ 
   $\wedge \text{quorum} = [\text{mon} \in \text{Monitors} \mapsto \text{TRUE}]$ 
   $\wedge \text{quorum\_sz} = \text{MonitorsLen}$ 
```

```
Init_state_vars  $\triangleq$ 
   $\wedge \text{isLeader} = [\text{mon} \in \text{Monitors} \mapsto \text{FALSE}]$ 
   $\wedge \text{state} = [\text{mon} \in \text{Monitors} \mapsto \text{STATE\_RECOVERING}]$ 
   $\wedge \text{phase} = [\text{mon} \in \text{Monitors} \mapsto \text{PHASE\_ELECTION}]$ 
```

```
Init_restart_vars  $\triangleq$ 
   $\wedge \text{pending\_pn} = [\text{mon} \in \text{Monitors} \mapsto 0]$ 
   $\wedge \text{pending\_v} = [\text{mon} \in \text{Monitors} \mapsto 0]$ 
   $\wedge \text{uncommitted\_pn} = [\text{mon} \in \text{Monitors} \mapsto 0]$ 
   $\wedge \text{uncommitted\_v} = [\text{mon} \in \text{Monitors} \mapsto 0]$ 
   $\wedge \text{uncommitted\_value} = [\text{mon} \in \text{Monitors} \mapsto \text{Nil}]$ 
```

```
Init_data_vars  $\triangleq$ 
   $\wedge \text{monitor\_store} = [\text{mon} \in \text{Monitors} \mapsto \text{Nil}]$ 
   $\wedge \text{values} = [\text{mon} \in \text{Monitors} \mapsto [\text{version} \in \{\} \mapsto \text{Nil}]]$ 
   $\wedge \text{accepted\_pn} = [\text{mon} \in \text{Monitors} \mapsto 0]$ 
   $\wedge \text{first\_committed} = [\text{mon} \in \text{Monitors} \mapsto 0]$ 
   $\wedge \text{last\_committed} = [\text{mon} \in \text{Monitors} \mapsto 0]$ 
```

```
Init_collect_vars  $\triangleq$ 
   $\wedge \text{num\_last} = [\text{mon} \in \text{Monitors} \mapsto 0]$ 
   $\wedge \text{peer\_first\_committed} = [\text{mon1} \in \text{Monitors} \mapsto [\text{mon2} \in \text{Monitors} \mapsto -1]]$ 
   $\wedge \text{peer\_last\_committed} = [\text{mon1} \in \text{Monitors} \mapsto [\text{mon2} \in \text{Monitors} \mapsto -1]]$ 
```

```
Init_lease_vars  $\triangleq$ 
   $\wedge \text{acked\_lease} = [\text{mon1} \in \text{Monitors} \mapsto [\text{mon2} \in \text{Monitors} \mapsto \text{FALSE}]]$ 
```

$Init\_commit\_vars \triangleq$   
 $\wedge pending\_proposal = [mon \in Monitors \mapsto Nil]$   
 $\wedge new\_value = [mon \in Monitors \mapsto Nil]$   
 $\wedge accepted = [mon1 \in Monitors \mapsto [mon2 \in Monitors \mapsto FALSE]]$   
 $Init \triangleq$   
 $\wedge Init\_global\_vars$   
 $\wedge Init\_state\_vars$   
 $\wedge Init\_restart\_vars$   
 $\wedge Init\_data\_vars$   
 $\wedge Init\_collect\_vars$   
 $\wedge Init\_lease\_vars$   
 $\wedge Init\_commit\_vars$   
 $\wedge step\_name = \text{"init"} \wedge number\_crashes = 0$

### Message manipulation

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

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

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

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

$@type: Set(MESSAGE) \Rightarrow Seq(MESSAGE);$

$SingleMessageSetToSeq(S) \triangleq$

$IF \exists elem \in S : TRUE THEN LET elem \triangleq CHOOSE x \in S : TRUE$   
 $IN \langle elem \rangle$   
 $ELSE \langle \rangle$

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

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

$WithMessage(m, msgs) \triangleq$

$[msgs EXCEPT ![m.from] =$   
 $[msgs[m.from] EXCEPT ![m.dest] = Append(msgs[m.from][m.dest], m)]$

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

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

$WithoutMessage(m, msgs) \triangleq$

$[msgs EXCEPT ![m.from] =$   
 $[msgs[m.from] EXCEPT ![m.dest] = Tail(msgs[m.from][m.dest])]$

Adds the message *m* to the network.

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

$@type: MESSAGE \Rightarrow Bool;$

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

Adds a set of messages to the network.

Variables changed:  $messages, message\_history$ .

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

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

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)  $\Rightarrow Bool$ ;

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

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;

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

Clear the variable *peer\_first\_committed*.

Variables changed: *peer\_first\_committed*.

@type: MONITOR  $\Rightarrow$  Bool;

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

$peer\_first\_committed' = [peer\_first\_committed \text{ EXCEPT } ![mon] =$   
 $[m \in Monitors \mapsto -1]]$

Clear the variable *peer\_last\_committed*.

Variables changed: *peer\_last\_committed*.

@type: MONITOR  $\Rightarrow$  Bool;

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

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

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

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

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

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

Choose peer values from *mon* last committed + 1 to peer last committed.

$\wedge \text{ LET } logs \triangleq (\text{DOMAIN } msg.values) \cap (last\_committed[mon] + 1 .. msg.last\_committed)$

IN  $\wedge values' = [values \text{ EXCEPT } ![mon] =$   
 $[i \in \text{DOMAIN } values[mon] \cup logs \mapsto$   
 IF  $i \in logs$   
 THEN  $msg.values[i]$   
 ELSE  $values[mon][i]]]$

Update last committed and first committed.

$\wedge last\_committed' = [last\_committed \text{ EXCEPT } ![mon] = \text{Max}(logs \cup \{last\_committed[mon]\})]$

$\wedge \text{ IF } logs \neq \{\} \wedge first\_committed[mon] = 0$

THEN  $first\_committed' =$   
 $[first\_committed \text{ EXCEPT } ![mon] = \text{Min}(logs)]$

ELSE  $first\_committed' =$   
 $[first\_committed \text{ EXCEPT } ![mon] = \text{Min}(logs \cup \{first\_committed[mon]\})]$

Update monitor store.

$\wedge \text{ IF } last\_committed'[mon] = 0$

THEN UNCHANGED *monitor\_store*

ELSE  $monitor\_store' = [monitor\_store \text{ EXCEPT } ![mon] = values'[mon][last\_committed'[mon]]]$



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

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

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

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

Trigger new election by incrementing epoch.

Variables changed: *epoch*.

@type: *Bool*;

*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*] = 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*;

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

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;

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

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;

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

### 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
     $\vee$  state'[mon] = STATE_UPDATING_PREVIOUS
   $\wedge$  quorum_sz = 1  $\vee$  num_last[mon] > MonitorsLen  $\div$  2
   $\wedge$  new_value[mon] = Nil
   $\wedge$  accepted' = [accepted EXCEPT ![mon] =
    [m  $\in$  Monitors  $\mapsto$  IF m = mon THEN TRUE ELSE FALSE]]
   $\wedge$  new_value' = [new_value EXCEPT ![mon] = v]
   $\wedge$  phase' = [phase EXCEPT ![mon] = PHASE_BEGIN]
   $\wedge$  values' = [values EXCEPT ![mon] =
    ((last_committed[mon] + 1) :> new_value'[mon]) @@ values[mon]]
   $\wedge$  Send_set(mon,
    {[type  $\mapsto$  OP_BEGIN,
      from  $\mapsto$  mon,
      dest  $\mapsto$  dest,
      last_committed  $\mapsto$  last_committed[mon],
      values  $\mapsto$  values'[mon],
      pn  $\mapsto$  accepted_pn[mon]] : dest  $\in$  {m  $\in$  Monitors \ {mon} : quorum[m]}
    })
   $\wedge$  pending_pn' = [pending_pn EXCEPT ![mon] = accepted_pn[mon]]
   $\wedge$  pending_v' = [pending_v EXCEPT ![mon] = last_committed[mon] + 1]

```

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] = FALSE
   $\wedge$  IF msg.pn < accepted_pn[mon]
    THEN
       $\wedge$  Discard(msg)
       $\wedge$  UNCHANGED  $\langle$ state, values, pending_pn, pending_v $\rangle$ 
    ELSE
       $\wedge$  msg.pn = accepted_pn[mon]
       $\wedge$  msg.last_committed = last_committed[mon]

      assign values[mon][last_committed[mon] + 1]
       $\wedge$  values' = [values EXCEPT ![mon] =

```

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

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;

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

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;

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

```

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

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

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

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;

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

### 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;

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

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' =
      [uncommitted_v EXCEPT ![mon] = last_committed[mon] + 1]
       $\wedge$  uncommitted_value' =
        [uncommitted_value EXCEPT ![mon] = values[mon][last_committed[mon] + 1]]
       $\wedge$  uncommitted_pn' = [uncommitted_pn EXCEPT ![mon] = pending_pn[mon]]
       $\wedge$  UNCHANGED  $\langle$  pending_pn, pending_v  $\rangle$ 
    ELSE UNCHANGED  $\langle$  restart_vars  $\rangle$ 

   $\wedge$  num_last' = [num_last EXCEPT ![mon] = 1]
   $\wedge$  Send_set(mon,
    {[type  $\mapsto$  OP_COLLECT,
      from  $\mapsto$  mon,
      dest  $\mapsto$  dest,
      first_committed  $\mapsto$  first_committed[mon],
      last_committed  $\mapsto$  last_committed[mon],
      pn  $\mapsto$  accepted_pn[mon]] : dest  $\in$  {m  $\in$  Monitors  $\setminus$  {mon} : quorum[m]}
    })
   $\wedge$  phase' = [phase EXCEPT ![mon] = PHASE_COLLECT]
   $\wedge$  UNCHANGED  $\langle$  isLeader, state  $\rangle$ 
   $\wedge$  UNCHANGED  $\langle$  epoch, quorum, quorum_sz, data_vars, lease_vars, commit_vars  $\rangle$ 

```

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

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

```

@type: (MONITOR, MESSAGE)  $\Rightarrow$  Bool;
handle_collect(mon, msg)  $\triangleq$ 
   $\wedge$  isLeader[mon] = FALSE
   $\wedge$  state' = [state EXCEPT ![mon] = STATE_RECOVERING]
   $\wedge$   $\vee$   $\wedge$  msg.first_committed > last_committed[mon] + 1
     $\wedge$  bootstrap
     $\wedge$  Discard(msg)
     $\wedge$  UNCHANGED  $\langle$  accepted_pn  $\rangle$ 
   $\vee$   $\wedge$  msg.first_committed  $\leq$  last_committed[mon] + 1
     $\wedge$  IF msg.pn > accepted_pn[mon]
      THEN accepted_pn' = [accepted_pn EXCEPT ![mon] = msg.pn]
      ELSE UNCHANGED accepted_pn
     $\wedge$  Reply([type  $\mapsto$  OP_LAST,

```

$from \mapsto mon,$   
 $dest \mapsto msg.from,$   
 $first\_committed \mapsto first\_committed[mon],$   
 $last\_committed \mapsto last\_committed[mon],$   
 $values \mapsto values[mon],$   
 $uncommitted\_pn \mapsto pending\_pn[mon],$   
 $pn \mapsto accepted\_pn'[mon]], msg)$   
 $\wedge \text{UNCHANGED } epoch$   
 $\wedge \text{UNCHANGED } \langle isLeader, phase, values, first\_committed, last\_committed, monitor\_store \rangle$   
 $\wedge \text{UNCHANGED } \langle quorum, quorum\_sz, restart\_vars, collect\_vars, lease\_vars, commit\_vars \rangle$

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

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

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

Variables changed: messages, *message\_history*, epoch, phase, *uncommitted\_pn*, *uncommitted\_v*, *uncommitted\_value*, *monitor\_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

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



$$\begin{aligned}
& \wedge \text{peer} \neq \text{mon} \\
& \wedge \text{peer\_last\_committed}'[\text{mon}][\text{peer}] \neq -1 \\
& \wedge \text{peer\_last\_committed}'[\text{mon}][\text{peer}] < \text{last\_committed}[\text{mon}] \\
& \wedge \text{quorum}[\text{peer}] \} \\
\text{IN } & \text{Reply\_set}(\text{mon}, \\
& \{ [\text{type} \quad \mapsto \text{OP\_COMMIT}, \\
& \quad \text{from} \quad \mapsto \text{mon}, \\
& \quad \text{dest} \quad \mapsto \text{dest}, \\
& \quad \text{last\_committed} \mapsto \text{last\_committed}'[\text{mon}], \\
& \quad \text{pn} \quad \mapsto \text{accepted\_pn}[\text{mon}], \\
& \quad \text{values} \quad \mapsto \text{values}[\text{mon}]] : \text{dest} \in \text{monitors\_behind} \\
& \}, \text{msg}) \\
& \wedge \vee \wedge \text{msg.pn} > \text{accepted\_pn}[\text{mon}] \\
& \quad \wedge \text{collect}(\text{mon}, \text{msg.pn}) \\
& \quad \wedge \text{check\_and\_correct\_uncommitted}(\text{mon}) \\
& \quad \wedge \text{UNCHANGED num\_last} \\
& \vee \wedge \text{msg.pn} = \text{accepted\_pn}[\text{mon}] \\
& \quad \wedge \text{num\_last}' = [\text{num\_last} \text{ EXCEPT } ![\text{mon}] = \text{num\_last}[\text{mon}] + 1] \\
& \quad \wedge \text{IF } \wedge \text{msg.last\_committed} + 1 \in \text{DOMAIN msg.values} \\
& \quad \quad \wedge \text{msg.last\_committed} \geq \text{last\_committed}'[\text{mon}] \\
& \quad \quad \wedge \text{msg.last\_committed} + 1 \geq \text{uncommitted\_v}[\text{mon}] \\
& \quad \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 \quad \wedge \text{uncommitted\_pn}' = \\
& \quad \quad [\text{uncommitted\_pn} \text{ EXCEPT } ![\text{mon}] = \text{msg.uncommitted\_pn}] \\
& \quad \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}) \\
& \quad \wedge \text{UNCHANGED } \langle \text{phase}, \text{accepted\_pn} \rangle \\
& \vee \wedge \text{msg.pn} < \text{accepted\_pn}[\text{mon}] \\
& \quad \wedge \text{check\_and\_correct\_uncommitted}(\text{mon}) \\
& \quad \wedge \text{UNCHANGED } \langle \text{phase}, \text{accepted\_pn}, \text{num\_last} \rangle \\
& \wedge \text{UNCHANGED epoch} \\
& \wedge \text{UNCHANGED 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)  $\triangleq$ 
   $\wedge$  isLeader[mon] = TRUE
   $\wedge$  num_last[mon] = quorum_sz
   $\wedge$  phase[mon] = 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 EXCEPT ![mon] = 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)
       $\wedge$  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  $\triangleq$ 
   $\wedge \exists$  mon  $\in$  Monitors :
     $\wedge$  quorum[mon]
     $\wedge$  isLeader' = [m  $\in$  Monitors  $\mapsto$  IF m = mon THEN TRUE ELSE FALSE]
     $\wedge$  state' = [m  $\in$  Monitors  $\mapsto$ 
      IF quorum_sz = 1 THEN STATE_ACTIVE ELSE STATE_RECOVERING]
     $\wedge$  phase' = [m  $\in$  Monitors  $\mapsto$  PHASE_ELECTION]
     $\wedge$  new_value' = [m  $\in$  Monitors  $\mapsto$  Nil]
     $\wedge$  pending_proposal' = [m  $\in$  Monitors  $\mapsto$  Nil]
     $\wedge$  epoch' = epoch + 1
     $\wedge$  messages' = [mon1  $\in$  Monitors  $\mapsto$  [mon2  $\in$  Monitors  $\mapsto \langle \rangle$ ]]
     $\wedge$  UNCHANGED  $\langle$ quorum, quorum_sz, accepted, message_history $\rangle$ 
     $\wedge$  UNCHANGED  $\langle$ data_vars, restart_vars, collect_vars, lease_vars $\rangle$ 

```

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

Variables changed: accepted\_pn, phase.

```

@type: MONITOR  $\Rightarrow$  Bool;
election_recover(mon)  $\triangleq$ 
   $\wedge$  quorum_sz > 1
   $\wedge$  phase[mon] = 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*;

$\text{crash\_mon}(\text{mon}) \triangleq$

$\wedge \text{quorum\_sz} > (\text{MonitorsLen} \div 2) + 1$   
 $\wedge \text{quorum}[\text{mon}] = \text{TRUE}$   
 $\wedge \text{quorum}' = [\text{quorum} \text{ EXCEPT } ![\text{mon}] = \text{FALSE}]$   
 $\wedge \text{quorum\_sz}' = \text{quorum\_sz} - 1$   
 $\wedge \text{bootstrap}$   
 $\wedge \text{number\_crashes}' = \text{number\_crashes} + 1$   
 $\wedge \text{UNCHANGED } \langle \text{messages}, \text{message\_history} \rangle$   
 $\wedge \text{UNCHANGED } \langle \text{state\_vars}, \text{restart\_vars}, \text{data\_vars}, \text{collect\_vars}, \text{lease\_vars}, \text{commit\_vars} \rangle$

Add monitor to the quorum.

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

$\text{restore\_mon}(\text{mon}) \triangleq$

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

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

Variables changed: epoch.

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

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

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

#### Dispatchers and next statement

Handle a message.

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

$\text{Receive}(\text{msg}) \triangleq$

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

$$\begin{aligned}
& \wedge \text{handle\_last}(\text{msg.dest}, \text{msg}) \\
& \wedge \text{step\_name}' = \text{"receive last"} \\
\vee & \wedge \text{msg.type} = \text{OP\_LEASE} \\
& \wedge \text{handle\_lease}(\text{msg.dest}, \text{msg}) \\
& \wedge \text{step\_name}' = \text{"receive lease"} \\
\vee & \wedge \text{msg.type} = \text{OP\_LEASE\_ACK} \\
& \wedge \text{handle\_lease\_ack}(\text{msg.dest}, \text{msg}) \\
& \wedge \text{step\_name}' = \text{"receive lease\_ack"} \\
\vee & \wedge \text{msg.type} = \text{OP\_BEGIN} \\
& \wedge \text{handle\_begin}(\text{msg.dest}, \text{msg}) \\
& \wedge \text{step\_name}' = \text{"receive begin"} \\
\vee & \wedge \text{msg.type} = \text{OP\_ACCEPT} \\
& \wedge \text{handle\_accept}(\text{msg.dest}, \text{msg}) \\
& \wedge \text{step\_name}' = \text{"receive accept"} \\
\vee & \wedge \text{msg.type} = \text{OP\_COMMIT} \\
& \wedge \text{handle\_commit}(\text{msg.dest}, \text{msg}) \\
& \wedge \text{step\_name}' = \text{"receive commit"}
\end{aligned}$$

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;

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

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

#### Safety invariants

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

@type: Bool;

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

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.

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

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

```
Inv_find_state(  
  ∧ 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  
  ∧ epoch = 2  
)
```

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

```
\ * Modification History  
\ * Last modified Thu Apr 15 13:49:52 WEST 2021 by afonsonf  
\ * Created Mon Jan 11 16:15:26 WET 2021 by afonsonf
```