──────────────────────── MODULE *PrimSecRep* ────────────────────────

EXTENDS *Util*, *TLC*

────────────────────────────────────────────────────────────────────

CONSTANTS   Values:
            *Rep*,            replicas
            *Val*,

            \*\**HOLES*:
            Types:
            *Object*,         object *IDs* or addresses
            *DataType*,      data type
            *WriteType*,     write structure (the one that gets written to the data)
            *ErrorTypes*,    different errors returned by the secondaries

            Abstract methods ("holes") :

            *StoreUpdate*(\_, \_, \_),    (r,o,w) perform write $w$ ( $\in$ *WriteType*)
                                 on local data of object $o$ of replica $r$
            *StoreRead*(\_, \_),   *read(r, o)* peform local read at replica $r$
                                 of object $r$. Returns the value
            *PrepareWr*(\_, \_, \_),   *prepareWr(r, o, w)* prepares write $w$ to send out
                                 to secondaries. This preparation is done by the primary.
            *ReplyFinishedWrite*(\_, \_, \_),  (r, o, w) replica $r$ replies to a finished write on object $o$
            *GetWriteResponse*(\_, \_, \_, \_, \_),  (r,o,version,sender) returns the response that should be sent out
                                 In case of some error (*e.g.*, wrong vers), should return $[type \mapsto$ "*noResp*"],
                                 in which case nothing will be committed to disk or sent out

            *Reply*(\_, \_, \_),        (r, type, reply) reply to client
            *NoReply*,    no reply is sent yet, so just keep the aux variables unchanged
            *InitDataVal*,    initial data value for the disks
            *InitPrim*(\_),   *InitPrim(o) =* value of the initial primary of object $o$
            *InitSec*(\_),    initial secondaries for object $o$

            *NoRep*,
            *Client*,
            *NoWrite*

VARIABLES *master*,       reliable master state

            *data*,           data at each replica
             *cache*,         cache of the master state at each replica
             *stat*,           state of each replica + object locks of each replica
             *channel*,      in-channel of each replica
             *resps*         write acknowledgements for each replica from the others

*vars* $\triangleq$ ⟨*master*, *cache*, *data*, *channel*, *resps*, *stat*⟩

1

$NoRep \triangleq \text{CHOOSE } v : v \notin Rep$
$Client \triangleq \text{CHOOSE } v : v \notin Rep \land v \neq NoRep$
$NoWrite \triangleq \text{CHOOSE } v : v \notin WriteTypes$

---

$WrReqMsg \triangleq$ Write requests
$\quad [type : \{\text{"wrReq"}\}, sender : Rep \cup \{Client\}, object : Object,$
$\quad\quad w : WriteType, version : Nat]$
$WrRespMsg \triangleq$ Write responses
$\quad [type : \{\text{"wrResp"}\}, sender : Rep, object : Object, ack : (\{\text{"ok"}\} \cup ErrorTypes),$
$\quad\quad w : WriteType, version : Nat]$
$Messages \triangleq$ Message types
$\quad WrReqMsg \cup WrRespMsg$

$TypeInvariant \triangleq$
$\quad \land Print \ (\text{"1"}, \text{TRUE})$
$\quad \land master \in$
$\quad\quad\quad [objects : [Object \rightarrow [version : Nat, prim : Rep \cup \{NoRep\}, sec : \text{SUBSET } (Rep)]],$
$\quad\quad\quad health : [Rep \rightarrow \{\text{"dead"}, \text{"alive"}\}]]$

$\quad \land Print \ (\text{"2"}, \text{TRUE})$
$\quad \land data \in [Rep \rightarrow [Object \rightarrow DataType]]$
$\quad \land Print \ (\text{"3"}, \text{TRUE})$
$\quad \land cache \in [Rep \rightarrow [Object \rightarrow [version : Nat, prim : Rep \cup \{NoRep\}, sec : \text{SUBSET } (Rep)]]]$
$\quad \land Print \ (\text{"5"}, \text{TRUE})$
$\quad \land channel \in [Rep \rightarrow Seq(Messages)]$

$\quad \land Print \ (\text{"7"}, \text{TRUE})$
$\quad \land stat \in [Rep \rightarrow [phase : \{\text{"alive"}, \text{"dead"}\}, lock : [Object \rightarrow \{\text{"rdy"}, \text{"busy"}\}],$
$\quad\quad\quad\quad\quad\quad in\_progress : [Object \rightarrow [val : WriteType, version : Nat]]]]$

$\quad \land Print \ (\text{"8"}, \text{TRUE})$
$\quad \land resps \in [Rep \rightarrow [Object \rightarrow [Rep \rightarrow \{\text{"ok"}, \text{"waiting"}, \text{"timeout"}, \text{"n/a"}\} \cup ErrorTypes]]]$
$\quad\quad\quad\quad$ separate set of responses per object
$\quad \land Print \ (\text{"9"}, \text{TRUE})$
$Init \triangleq$
$\quad \land master =$
$\quad\quad [objects \mapsto [o \in Object \mapsto$
$\quad\quad\quad\quad [version \mapsto 0, prim \mapsto InitPrim(o), sec \mapsto InitSec(o)]],$
$\quad\quad health \mapsto [r \in Rep \mapsto \text{"alive"}]]$ but all replicas are alive

$\quad \land cache = [r \in Rep \mapsto [o \in Object \mapsto$
$\quad\quad\quad\quad [version \mapsto 0, prim \mapsto InitPrim(o), sec \mapsto InitSec(o)]]]$

$\quad \land data = [r \in Rep \mapsto [o \in Object \mapsto InitDataVal]]$
$\quad \land channel = [r \in Rep \mapsto \langle\rangle]$ no message yet
$\quad \land resps = [r \in Rep \mapsto [o \in Object \mapsto [r1 \in Rep \mapsto \text{"waiting"}]]]$

2

$$\land\ stat = [r \in Rep \mapsto [phase \mapsto \text{"alive"},\ lock \mapsto [o \in Object \mapsto \text{"rdy"}],$$
$$in\_progress \mapsto [o \in Object \mapsto [val \mapsto NoWrite,\ version \mapsto 0]]]]]$$

all reps alive and ready to accept updates

---

In the most general case, the master does not participate actively. It's only for mitigation and for a reliable source of accurate information. In master-monitored replicated systems (like *GFS*), the master is active and monitors replicas, releases primary leases, removes replicas, adds replicas, etc. In self-monitoring replica sys, (like Blue), master is not active in the spec.

$MasterActions \triangleq \text{FALSE}$

---

TIME AND CHANNELS: (sources of timeouts and losses)

The channel to replica $r$ loses one of the messages that should have arrived there.

$\_TransLoss(r) \triangleq$
    $\land\ channel[r] \neq \langle\rangle$  the channel is not empty

    $\land\ channel' = [channel \text{ EXCEPT } ![r] = Tail(@)]$
    $\land\ \text{UNCHANGED } \langle master,\ cache,\ data,\ stat,\ resps\rangle$
$TransLoss \triangleq \exists\, r \in Rep : \_TransLoss(r)$

(Ex-)Primary $r$ times out waiting for response from replica $s$.

$\_Timeout(r,\ o,\ s) \triangleq$
  LET $updateMsgIsLost(chan,\ update) \triangleq$
        $\forall\, i \in 1\,..\,Len(chan) : chan[i].w \neq update.val$
  IN
    $\land\ stat[r].phase = \text{"alive"}$  $r$ is alive
    $\land\ stat[r].lock[o] = \text{"busy"}$  $r$ is indeed waiting for responses to a write request
    $\land\ resps[r][o][s] = \text{"waiting"}$  $r$ is still waiting for response from $s$

    $\land\ updateMsgIsLost(channel[s],\ stat[r].in\_progress[o]) \setminus *$ REDUCE SPACE: the update message sent to $s$ was lost – *TODO*

    $\land\ resps' = [resps \text{ EXCEPT } ![r][o][s] = \text{"timeout"}]$  response of $s$ timed out
    $\land\ \text{UNCHANGED } \langle master,\ cache,\ data,\ stat,\ channel\rangle$
$Timeout \triangleq$
  $\exists\, r \in Rep,\ o \in Object,\ s \in Rep :$
    $\land\ r \neq s$
    $\land\ \_Timeout(r,\ o,\ s)$

$TimeActions \triangleq$
  $\lor\ TransLoss \qquad \lor\ Timeout$

---

REPLICA actions.

Replica $r$ suddenly dies.

$\_ReplicaDeath(r) \triangleq$
    $\wedge\ stat[r].phase =$ "alive"   replica used to be alive

    $\wedge\ stat' = [stat$ EXCEPT $![r].phase =$ "dead"]   declares itself dead
    $\wedge$ UNCHANGED $\langle master,\ cache,\ data,\ channel,\ resps \rangle$
$ReplicaDeath \triangleq \exists\, r \in Rep : \_ReplicaDeath(r)$

Replica $r$ updates its cache w/ the accurate version from master. One could imagine multiple times when this could be triggered in the actual implem of, *e.g.*, *GFS* (*e.g.*, during *HeartBeat* protocol, after a client comes w/ a higher *version* $\neq$ , etc.). Since I don't know what happens in the real protocol exactly, I will assume it can happen anytime.

$\_ReadVersion(r,\ o) \triangleq$
    $\wedge\ stat[r].phase \neq$ "dead"   replica is not dead yet
    $\wedge\ cache[r][o].version \neq master.objects[o].version$   $r$'s cache is out-of-date

    $\wedge\ cache' = [cache$ EXCEPT $![r][o] = master.objects[o]]$   cache new version
    $\wedge$ UNCHANGED $\langle master,\ data,\ channel,\ resps,\ stat \rangle$
$ReadVersion \triangleq$
    $\wedge \exists\, r \in Rep,\ o \in Object : \_ReadVersion(r,\ o)$

Replica $r$ drops a write req $w$. The reasons might be various. ( *E.g.*, in *GFS*, one reason can be that prim cannot find the data associated w/ this request in its *LRU*).

$DropWrite(r,\ o,\ w) \triangleq$   changes $\langle data,\ channel,\ stat,\ resps \rangle$
    $\wedge$ UNCHANGED $\langle data,\ channel,\ stat,\ resps \rangle$

Primary continues w/ a client request. It pushes the data to its local store (always considered to succeed from this at this point, b/c anyway I provide for general request dropping). The primary then sends update messages to all secondaries. The write needs preparation before it's given out to secondaries. In *GFS* (*e.g.*), the preparation means deciding on an address to write at and setting this adr in the msg.

$PrimaryContinuesWrite(r,\ o,\ w,\ allreplicas,\ version) \triangleq$   changes $\langle data,\ channel,\ stat,\ resps \rangle$, sends *Reply*
    $\wedge$ LET
        $prepared\_wr \triangleq PrepareWr(r,\ o,\ w)$
        $wr\_req \triangleq [type \mapsto$ "wrReq"$,\ sender \mapsto r,\ object \mapsto o,$
                $version \mapsto version,\ w \mapsto prepared\_wr]$
                    the write needs to be prepared for the secondaries
        IN
            $\wedge$ IF $prepared\_wr \neq NoWrite$ THEN   No error while preparing

                $\wedge\ StoreUpdate(r,\ o,\ prepared\_wr)$   perform the update locally
                $\wedge\ channel' = [r1 \in Rep \mapsto$   send the update req to secondaries
                    IF $r1 \in allreplicas \setminus \{r\}$ THEN
                        $channel[r1] \circ \langle wr\_req \rangle$
                    ELSE $channel[r1]]$   not a secondary

                $\wedge\ stat' = [stat$ EXCEPT $![r].lock[o] =$ "busy",   the replica becomes busy
                                                and will not accept further updates on object $o$
                                                until this write finishes.

4

$$![r].in\_progress[o] = [val \mapsto prepared\_wr, \; version \mapsto version]]$$
maintain the write, for future reference
$$\wedge \; resps' = [resps \; \text{EXCEPT} \; ![r][o] = [r1 \in Rep \mapsto$$
$$\text{IF} \; r = r1 \; \text{THEN} \; \text{"ok"} \quad \text{this replica has already answered}$$
$$\text{ELSE} \; \text{IF} \; r1 \in allreplicas$$
$$\text{THEN} \; \text{"waiting"} \quad \text{wait for } ack \text{ from secondaries}$$
$$\text{ELSE} \; \text{"n/a"}]] \quad \text{don't wait for } ack \text{ from non-secondaries}$$
$\wedge \; NoReply$   don't return yet the response

$\wedge$ UNCHANGED $data$

ELSE   Preparing the message resulted in error, indicating
that this write shouldn't go ahead.
Skip it and announce the client of error
$\wedge$ UNCHANGED $\langle data, \; channel, \; stat, \; resps \rangle$
$\wedge \; Reply(r, \; \text{"NoWrite"}, \; prepared\_wr)$

Primary $r$ handles a write request. Note that the writes are blocking and processed by one primary one at a time. Only after a write finishes, does the primary start another one.

$PrimaryWrite(r, \; o, \; w, \; ver) \; \triangleq$   changes $\langle data, \; channel, \; stat, \; resps \rangle$, sends $Reply$
   $\wedge \; cache[r][o].prim = r$    $r$ believes itself to be the primary
   $\wedge \; stat[r].lock[o] = \text{"rdy"}$    the primary doesn't have other writes in progress on object $o$

   $\wedge \;\; \vee \; \wedge \; PrimaryContinuesWrite(r, \; o, \; w, \; \{cache[r][o].prim\} \cup cache[r][o].sec, \; cache[r][o].version)$
                     EITHER: perform the write

     $\vee \; \wedge \; DropWrite(r, \; o, \; w) \backslash *$ OR: drop the request alltogether
       $\wedge \; Reply(r, \; \text{"wrFinished"}, \; [ack \mapsto \text{"error"}, \; object \mapsto o, \; w \mapsto w])$

$VersionBased\_GetWriteResponse(r, \; o, \; w, \; version, \; sender) \; \triangleq$
   IF $(cache[r][o].version \leq version)$ THEN    good version
     $[type \mapsto \text{"wrResp"}, \; sender \mapsto r, \; object \mapsto o,$
     $w \mapsto w, \; ack \mapsto \text{"ok"}, \; version \mapsto cache[r][o].version]$
   ELSE    bad version
     $[type \mapsto \text{"wrResp"}, \; sender \mapsto r, \; object \mapsto o,$
     $w \mapsto w, \; ack \mapsto \text{"badver"}, \; version \mapsto cache[r][o].version]$

Secondary continues the write request from prim. It pushes the data to its local store (always considered to succeed from this at this point, b/c anyway I provide for general request dropping & version has been or doesn't need to be checked.

$SecondaryContinuesWrite(r, \; o, \; w, \; ver, \; sender) \; \triangleq$    changes $\langle data, \; channel \rangle$
   LET
      $wr\_resp \; \triangleq \; GetWriteResponse(r, \; o, \; w, \; ver, \; sender)$
   IN
     IF $wr\_resp.type = \text{"wrResp"}$ THEN
       $\wedge \; wr\_resp.ack = \text{"ok"} \Rightarrow StoreUpdate(r, \; o, \; w)$   store the update persistenly locally
       $\wedge \; wr\_resp.ack \neq \text{"ok"} \Rightarrow$ UNCHANGED $data$   don't store, there's an error

$$\land \ channel' = [channel \ \text{EXCEPT}$$
$$![sender] = @ \circ \langle wr\_resp \rangle, \quad \text{send response}$$
$$![r] = Tail(@)] \quad \text{remove from channel}$$
$$\text{ELSE} \quad \land \ channel' = [channel \ \text{EXCEPT} \ ![r] = Tail(@)] \quad \text{no response should be sent back}$$
$$\land \ \text{UNCHANGED} \ data$$

A secondary replica processes a write request from primary. The secondary fully executes it and sends a reply back to sender. We don't simulate dropping of the request by *sec.* because there's no need to, given the general *TransLoss*.

$SecondaryWrite(r, o, w, ver, sender) \triangleq$ changes $\langle data, channel, stat, resps \rangle$
$\quad \land \ \lor SecondaryContinuesWrite(r, o, w, ver, sender)$

$\qquad$ No need to drop the write at the secondaries. This drop will happen due to
$\qquad$ *TransLoss()*.
$\qquad \lor DropWrite(r, o\ w) \setminus$ * OR: drop the request alltogether
$\quad \land \ \text{UNCHANGED} \ \langle stat, resps \rangle$

$A(n \ ex-)$primary processes a response to a write request from replica sender. The response might already be too late (it has already timed out), or it might be or might be in-time.

$ProcessWriteResp(r, o, w, sender, ack) \triangleq$ changes $\langle data, channel, stat, resps \rangle$
$\quad \land \ resps' = [resps \ \text{EXCEPT} \ ![r][o][sender] =$
$\qquad\qquad \text{IF} \ @ = \text{"timeout"} \lor @ = \text{"n/a"} \ \text{THEN} \ @ \quad \text{resp came too late, it's already expired}$
$\qquad\qquad \text{ELSE} \ ack] \quad \text{either } ok \text{ or some system-specific error}$

$\quad \land \ channel' = [channel \ \text{EXCEPT} \ ![r] = Tail(@)]$
$\quad \land \ \text{UNCHANGED} \ \langle data, stat \rangle$

A replica processes a message in its incoming channel. Depending on the type of message (update req, update resp), replica acts according to the three functions above.

$\_ProcessMessage(r) \triangleq$
$\quad \land \ channel[r] \neq \langle \rangle \quad \text{I have a message to process}$

$\quad \land \ \text{LET} \ m \triangleq Head(channel[r]) \text{IN}$
$\qquad\qquad \land \ \text{IF} \ m.type = \text{"wrReq"} \ \text{THEN}$
$\qquad\qquad\qquad \text{IF} \ m.sender = Client \ \text{THEN} \quad \text{upate request from client to a prim}$
$\qquad\qquad\qquad\qquad \text{Enters here only when } Client \text{ writes are not a 0-stage action}$
$\qquad\qquad\qquad\qquad PrimaryWrite \ (r, m.object, m.w, m.version)$
$\qquad\qquad\qquad\qquad TODO: \ NB: \text{This function sends a } Reply, \text{while the others don't} - BUG$

$\qquad\qquad\qquad \text{ELSE} \quad \text{It's a } wrReq \text{ from a primary to a secondary}$
$\qquad\qquad\qquad\qquad SecondaryWrite(r, m.object, m.w, m.version, m.sender)$

$\qquad\qquad \text{ELSE} \ \text{IF} \ m.type = \text{"wrResp"} \ \text{THEN}$
$\qquad\qquad\qquad\qquad ProcessWriteResp(r, m.object, m.w, m.sender, m.ack)$
$\qquad\qquad \text{ELSE}$
$\qquad\qquad\qquad\qquad \land \ Print(\text{"BUGGGG!!!! Wrong message type!"}, m) \neq \langle \rangle$

$\quad \land \ \text{UNCHANGED} \ \langle master, cache \rangle$

$ProcessMessage \triangleq$
   $\exists\, r \in Rep :$
      $\wedge\ stat[r].phase =$ "alive"   $r$ must be alive to process a message
      $\wedge\ \_ProcessMessage(r)$

A replica (either a current primary or not) finishes a write that it started some time ago (when it considered itself to be a primary).

$\_FinishWrite(r,\,o) \triangleq$
LET
     $NotFinishedWrReqs(m) \triangleq m.object \neq o \vee m.sender \neq r$
                                  this message does not refer to a write on object $o$
                                  initiated by replica $r$. So, I don't declare it finished
     $NotFinishedWrResps(m) \triangleq m.object \neq o$
IN
   $\wedge\ stat[r].phase =$ "alive"  allow writes to finish even for those started by dead replicas
                                       *TODO: Shouldn*'t do this for Blue??????
   $\wedge\ stat[r].lock[o] =$ "busy"  $r$ has indeed a write started

   $\wedge\ \forall\, r1 \in Rep : resps[r][o][r1] \neq$ "waiting"  all the replicas have submitted
                                           their responses, or have timed-out
                                         or weren't supposed to answer

   $\wedge\ ReplyFinishedWrite(r,\,o,\,stat[r].in\_progress[o])$
                                should send the appropriate *Reply*, based on state of the *resps*
                                but not change any of master, data, cache, etc.!

  Re-initialize all write-related state, to prepare the replica for the next write
   $\wedge\ stat' = [stat$ EXCEPT $![r].lock[o] =$ "rdy",  $r$ is ready to accept new updates
                                           (if it's still prim, of course)
                      $![r].in\_progress[o] = [val \mapsto NoWrite,\ version \mapsto 0]]$
  re-init the responses, to prepare $r$ for the next write
   $\wedge\ resps' = [resps$ EXCEPT $![r][o] = [r1 \in Rep \mapsto$ "waiting"$]]$

  Remove ALL messages related to this object. Leave the rest intact
   $\wedge\ channel' = [r1 \in Rep \mapsto$ IF $r1 \neq r$ THEN $SelectSeq(channel[r1],\ NotFinishedWrReqs)$
                                ELSE  $SelectSeq(channel[r],\ NotFinishedWrResps)]$
               I think that the above is not actually cheating, b/c you can imagine
               an implementation where replicas would simply identify older
               updates (via a sequence number, which is anyway necessary for
               keeping the order of updates.)
   $\wedge$ UNCHANGED $\langle master,\ data,\ cache \rangle$

$FinishWrite \triangleq$
 $\exists\, r \in Rep,\ o \in Object : \_FinishWrite(r,\,o)$

$ReplicaActions \triangleq$
   $\vee\ ReplicaDeath$

$\vee\ ReadVersion$
$\vee\ ProcessMessage$
$\vee\ FinishWrite$ <span style="background:#ccc">a replica finishes a write it has started</span>

---

<span style="background:#ccc">CLIENT actions.</span>

<span style="background:#ccc">Client wants to perform write $w$ on object $o$. The request is performed in a 0-stage fashion. The method models caching of replica locations in the client – BUT IT'S NOT CORRECT !!!$TODO\ TODO$.</span>

$\_CliWrite(prim,\ o,\ w)\ \triangleq$
  $\wedge\ PrimaryWrite(prim,\ o,\ w,\ master.objects[o].version)$
  $\wedge\ \text{UNCHANGED}\ \langle master,\ cache \rangle$
$CliWrite(o,\ w)\ \triangleq$
  $\exists\, r \in Rep :$
    $\wedge\ stat[r].phase =\ \text{"alive"}$
    $\wedge\ \_CliWrite(r,\ o,\ w)$

<span style="background:#ccc">In the most general case, a read is typically performed on the local copy of some replica that stores that object (is either a $sec$ or a primary). How you choose that replica depends on the protocol. In $GFS$, it's any replica. In Blue, it's gotta be prm The semantic of the read is that I read a whole object.</span>

$\_CliRead(r,\ o)\ \triangleq$
  $\wedge\ \text{LET}\ val\ \triangleq\ StoreRead(r,\ o)$
    $\text{IN}\quad Reply(r,\ \text{"rd"},\ val)$
  $\wedge\ \text{UNCHANGED}\ \langle master,\ cache,\ data,\ stat,\ resps,\ channel \rangle$

$CliRead(o)\ \triangleq$
  $\exists\, r \in Rep :$
    $\wedge\ stat[r].phase =\ \text{"alive"}$

    $\wedge\ master.health[r] =\ \text{"alive"}$ <span style="background:#ccc">ASSUMPTION : $Uncomment$ this if you want to test read-last-successful-X ($X =$ write/append). This ensures I don't read from a live but stale replica.</span>
    $\wedge\ r \in cache[r][o].sec \cup \{cache[r][o].prim\}$
       <span style="background:#ccc">$r$ considers itself either a $sec$ or primary of $o$</span>
    $\wedge\ \_CliRead(r,\ o)$

$ClientActions\ \triangleq$
  $\exists\, o \in Object :$
    $\vee\ \exists\, w \in WriteType : CliWrite(o,\ w)$
    $\vee\ CliRead(o)$

---

$Next\ \triangleq\ \vee\ MasterActions$
          $\vee\ TimeActions$
          $\vee\ ReplicaActions$
          $\vee\ ClientActions$

$Spec \triangleq \land Init \land \Box[Next]_{vars}$

Invariants
$AllInvariants \triangleq$
$\quad \land TypeInvariant$

THEOREM $Spec \Rightarrow \Box AllInvariants$