
MODULE *MCBlueRep*

EXTENDS *bluerep*

VARIABLES *prevread*,
 for mapping to regular model
last_read_val,
last_committed_write

mcbblue_newvars $\triangleq \langle prevread, last_read_val, last_committed_write \rangle$
mcbbluevars $\triangleq mcvvars \circ mcbblue_newvars$

MCBlue_TypeInvariant \triangleq
 $\wedge Blue_TypeInvariant$
 $\wedge Print("M C typeInvar", TRUE)$

$\wedge prevread \in [BObject \rightarrow Val \cup \{NoVal\}]$
 $\wedge last_read_val \in [BObject \rightarrow Val \cup \{NoVal\}]$
 $\wedge last_committed_write \in [BObject \rightarrow [val : Val \cup \{NoVal\}, version : Nat]]$

MCBlue_Init \triangleq
 $\wedge Blue_Init$

$\wedge prevread = [o \in BObject \mapsto NoVal]$
 $\wedge last_read_val = [c \in BObject \mapsto NoVal]$
 $\wedge last_committed_write = [c \in BObject \mapsto [val \mapsto NoVal, version \mapsto 0]]$

MCBlue_Reply(*r*, *type*, *reply*) \triangleq
 $\vee \wedge type = "wrFinished"$
 All finished writes in Blue are guaranteed to be committed, even if they are committed to less than threshold replicas.
 $\wedge last_committed_write' = [last_committed_write \text{ EXCEPT } ![reply.object] =$
 IF $reply.w.version \geq @.version$ THEN $reply.w$ ELSE $@$]
 $\wedge prevread' = [prevread \text{ EXCEPT } ![reply.object] = NoVal]$ invalidate the previous read
 $\wedge UNCHANGED last_read_val$

$\vee \wedge type = "rd"$
 $\wedge prevread' = [prevread \text{ EXCEPT } ![reply.object] = reply.val]$
 $\wedge last_read_val' = [last_read_val \text{ EXCEPT } ![reply.object] = reply.val]$
 $\wedge UNCHANGED \langle last_committed_write \rangle$

MCBlue_NoReply \triangleq
 UNCHANGED *mcbblue_newvars*

MCBlue_ReplicaActions \triangleq
 $\vee ReplicaDeath \wedge UNCHANGED mcbblue_newvars$

$\vee \text{ReadVersion} \wedge \text{UNCHANGED } mcb\text{blue_newvars}$
 $\vee \text{ProcessMessage} \wedge \text{UNCHANGED } mcb\text{blue_newvars}$
 $\vee \text{PrimaryKillsFailedRep} \wedge \text{UNCHANGED } mcb\text{blue_newvars}$
 $\vee \text{FinishWrite}$ a replica finishes a write it has started

$MC\text{Blue_Next} \triangleq$
 $\vee \text{Blue_MasterActions} \wedge \text{UNCHANGED } mcb\text{blue_newvars}$
 $\vee \text{TimeActions} \wedge \text{UNCHANGED } mcb\text{blue_newvars}$
 $\vee \text{MCBlue_ReplicaActions}$
 $\vee \text{Blue_ClientActions}$

$MC\text{Blue_Spec} \triangleq MC\text{Blue_Init} \wedge \Box[MC\text{Blue_Next}]_{mcb\text{bluevars}}$

Invariants

This invariant should only hold if all reads are to the primary or if a replica never falsely declares another replica dead.

$\text{ReadNowLastCommitted} \triangleq$
 $\forall o \in B\text{Object} :$
 $\forall r \in \text{Rep} :$
 $(\wedge \text{cache}[r][o].\text{prim} = r$ r thinks it's primary
 $\wedge \text{stat}[r].\text{phase} = \text{"alive"}$ r is alive
 $\wedge \text{stat}[r].\text{lock}[o] = \text{"rdy"}$
 $\wedge \text{master.health}[r] = \text{"alive"}$ ASSUMPTION : I never read from a stale replica (this ensured by read leases)
 $)$ r is ready to answer query
 \Rightarrow
 $\text{data}[r][o] = \text{last_committed_write}[o].\text{val}$

$\text{ReadLastCommitted} \triangleq$
 $\forall o \in B\text{Object} :$
 $\vee \text{prevread}[o] = \text{No Val}$ there's been no read since the last commit
 $\vee \text{prevread}[o] = \text{last_committed_write}[o].\text{val}$ OR the read value is that of the last committed write

$MC\text{Blue_AllInvariants} \triangleq$
 $\wedge \text{AllInvariants}$
 $\wedge \text{ReadNowLastCommitted}$
 $\wedge \text{ReadLastCommitted}$

Refinement mapping from Blue \rightarrow SimpleStore

$\text{NotFailed}(rsp) \triangleq$
 $\forall r \in \text{Rep} : \text{rsp}[r] \neq \text{"badver"}$ the only way to fail is to have a badversion response in Blue
 $\text{map_pending_wreq}(o) \triangleq$
 $\text{LET } \text{has_wreq}(r) \triangleq$

```

 $\wedge stat[r].lock[o] = \text{"busy"}$ 
 $\wedge stat[r].in\_progress[o].val \neq NoWrite$ 
 $\wedge stat[r].in\_progress[o].version \geq last\_committed\_write[o].version$ 
 $\wedge NotFailed(resps[r][o])$  the write is not yet failed

 $F[reps \in SUBSET (Rep)] \triangleq$ 
  IF  $(reps = \{\}) \vee (\forall r \in reps : \neg has\_wrreq(r))$ 
  THEN  $\langle \rangle$ 
  ELSE LET  $r\_minver \triangleq$  CHOOSE  $r \in reps :$ 
     $\wedge has\_wrreq(r)$ 
    pick the one  $w/$  min version number
     $\wedge \forall r1 \in reps : (has\_wrreq(r1) \Rightarrow stat[r].in\_progress[o].version \leq stat[r1].in\_progress[o].version)$ 
  IN  $\langle stat[r\_minver].in\_progress[o].val \circ F[reps \setminus \{r\_minver\}] \rangle$ 
IN (  $Print(\text{"wrreq"} = ", F[Rep])$ 

 $map\_pending\_wrresp(o) \triangleq$  map the overlapping writes
LET
   $has\_wrresp(r) \triangleq$ 
     $\wedge stat[r].lock[o] = \text{"busy"}$ 
     $\wedge stat[r].in\_progress[o].val \neq NoWrite$ 
     $\wedge stat[r].in\_progress[o].version < last\_committed\_write[o].version$ 
     $\wedge NotFailed(resps[r][o])$  the write is not yet failed

   $F[reps \in SUBSET (Rep), ret \in [Val \rightarrow Nat]] \triangleq$ 
    IF  $(reps = \{\}) \vee (\forall r \in reps : \neg has\_wrresp(r))$ 
    THEN  $ret$ 
    ELSE LET  $r \triangleq$  CHOOSE  $r \in reps : has\_wrresp(r)$ 
      IN  $F[reps \setminus \{r\}, [ret \text{ EXCEPT } ![stat[r].in\_progress[o].val] = @ + 1]]$ 
  IN (  $Print(\text{"wrresp"} = ", F[Rep, [v \in Val \mapsto 0]])$ 

 $map\_failed\_wr(o) \triangleq$ 
LET
   $G[reps \in SUBSET (Rep), ret \in [Val \rightarrow Nat]] \triangleq$ 
    IF  $reps = \{\} \vee \forall r \in reps : NotFailed(resps[r][o])$ 
    THEN  $ret$ 
    ELSE LET  $r \triangleq$  CHOOSE  $r \in reps : \neg NotFailed(resps[r][o])$   $r$  is failed
      IN  $G[reps \setminus \{r\}, [ret \text{ EXCEPT } ![stat[r].in\_progress[o].val] = @ + 1]]$ 
  IN (  $Print(\text{"failed\_wr"} = ", G[Rep, [v \in Val \mapsto 0]])$ 

 $map\_store(o) \triangleq$  map the store value
 $last\_committed\_write[o].val$ 

 $map\_last\_read\_val(o) \triangleq$ 
 $last\_read\_val[o]$ 

 $mapped\_object \triangleq$  CHOOSE  $o \in BObject : TRUE$ 
 $SS \triangleq$  INSTANCE  $simplestore\_quickrd$  WITH

```

```

pending_wrreq ← map_pending_wrreq(mapped_object),
pending_wrrsp ← map_pending_wrrsp(mapped_object),
failed_wr ← map_failed_wr(mapped_object),
last_read_val ← map_last_read_val(mapped_object),
store ← map_store(mapped_object),
pending_rd ← 0 0-stage reads, so I never buffer reads in the channel

```

$Blue_Implements_SS \triangleq SS!SS1_Spec$

THEOREM $MCBlue_Spec \Rightarrow MCBlue_AllInvariants$
