─────── MODULE *voldemort* ───────

/*Replicated storage protocol with clienside routing.

/*Debaditya Basak, 11 Nov 2016

EXTENDS *Integers*, *Sequences*, *FiniteSets*, *TLC*

CONSTANTS $N$, $C$, $STOP$, $ReadQ$, $WriteQ$, $FAILNUM$

ASSUME $N = 5 \wedge C = 1 \wedge STOP < 10 \wedge 1 \leq ReadQ \wedge ReadQ \leq 3$
$\wedge 1 \leq WriteQ \wedge WriteQ \leq 3 \wedge 0 \leq FAILNUM \wedge FAILNUM \leq 2$

$Nodes \triangleq 1 .. N$

$Clients \triangleq N + 1 .. N + C$  /*should give different *ID* space to Client

**--algorithm** *voldemort* **{**
   **variable** $FailNum = FAILNUM$,
        $state =$ "Reading",
        $state1 =$ "InProcess",
        $wQ = WriteQ$,
        $up = [n \in Nodes \mapsto \text{TRUE}]$,  /*Initially all nodes are up
        $db = [n \in Nodes \mapsto [ver \mapsto 0, val \mapsto 0]]$;
          /*All nodes have database, wherein $[ver = 0, val = 0]$ stored for the item

   **define** {
   $UpNodes \triangleq \{i \in Nodes : up[i] = \text{TRUE}\}$
   $ReturnReadQ \triangleq$ CHOOSE $i \in$ SUBSET $(UpNodes) : Cardinality(i) = ReadQ$
   $ReturnWriteQ \triangleq$ CHOOSE $i \in$ SUBSET $(UpNodes) : Cardinality(i) = WriteQ$
   $NodeR \triangleq$ CHOOSE $i \in ReturnReadQ : \forall j \in ReturnReadQ : db[i].ver \geq db[j].ver$
   $NodeW(Q) \triangleq$ CHOOSE $i \in Q : \forall j \in Q : i \leq j$
   $ClientNotWriting \triangleq state1 =$ "InProcess" $\wedge state =$ "Reading"
   $ClientAtomicWriteDone \triangleq state =$ "Writing" $\wedge state1 =$ "WriteEnd"
   }

   **fair process** ( $c \in Clients$ )
   **variable** $cntr = 0$, $hver = 0$, $Q = \{\}$;
   {
  $CL$: **while** ( $cntr \leq STOP$ ) {
      $state :=$ "Reading" ;
      $state1 :=$ "InProcess" ;
      $cntr := cntr + 1$ ;
      $hver := db[NodeR].ver + 1$ ;
      $Q := ReturnWriteQ$ ;
       /*Nodes can fail or come back up between atomic states *CL* and *CL1*
      $CL1$: **while** ( $Q \neq \{\}$ ) {
         $state :=$ "Writing" ;
         $db[NodeW(Q)].ver := hver \parallel db[NodeW(Q)].val := cntr$ ;
         $Q := Q \setminus \{NodeW(Q)\}$ ;
         **if** ( $Q = \{\}$ ) {
          $state1 :=$ "WriteEnd" ;
          }

```
                    }
                }
            }

    fair process ( n ∈ Nodes )
    {
     NODE : while ( TRUE ∧ ClientNotWriting ) {
        /*To make Clients-process WRITE atomic
        /*Nodes state change only when ClientNotWriting is TRUE
            if ( FailNum > 0 ∧ up[self] = TRUE ) {   /*Storage node can fail
                up[self] := FALSE ;
                FailNum := FailNum − 1 ;
                }
            else if ( up[self] = FALSE ) {   /*Or recover
                up[self] := TRUE ;
                FailNum := FailNum + 1 ;
                }
            }
        }

}
```

**********************************************************************

BEGIN TRANSLATION

VARIABLES $FailNum$, $state$, $state1$, $wQ$, $up$, $db$, $pc$

define statement

$UpNodes \triangleq \{i \in Nodes : up[i] = \text{TRUE}\}$

$ReturnReadQ \triangleq \text{CHOOSE } i \in \text{SUBSET } (UpNodes) : Cardinality(i) = ReadQ$

$ReturnWriteQ \triangleq \text{CHOOSE } i \in \text{SUBSET } (UpNodes) : Cardinality(i) = WriteQ$

$NodeR \triangleq \text{CHOOSE } i \in ReturnReadQ : \forall j \in ReturnReadQ : db[i].ver \geq db[j].ver$

$NodeW(Q) \triangleq \text{CHOOSE } i \in Q : \forall j \in Q : i \leq j$

$ClientNotWriting \triangleq state1 = \text{"InProcess"} \quad \wedge state = \text{"Reading"}$

$ClientAtomicWriteDone \triangleq state = \text{"Writing"} \wedge state1 = \text{"WriteEnd"}$

VARIABLES $cntr$, $hver$, $Q$

$vars \triangleq \langle FailNum, state, state1, wQ, up, db, pc, cntr, hver, Q \rangle$

$ProcSet \triangleq (Clients) \cup (Nodes)$

$Init \triangleq$   Global variables
$\wedge FailNum = FAILNUM$
$\wedge state = \text{"Reading"}$
$\wedge state1 = \text{"InProcess"}$
$\wedge wQ = WriteQ$
$\wedge up = [n \in Nodes \mapsto \text{TRUE}]$
$\wedge db = [n \in Nodes \mapsto [ver \mapsto 0, val \mapsto 0]]$

2

Process $c$
$\land cntr = [self \in Clients \mapsto 0]$
$\land hver = [self \in Clients \mapsto 0]$
$\land Q = [self \in Clients \mapsto \{\}]$
$\land pc = [self \in ProcSet \mapsto \text{CASE } self \in Clients \to \text{"CL"}$
$\qquad\qquad\qquad\qquad\quad \Box \quad self \in Nodes \to \text{"NODE"}]$

$CL(self) \triangleq \land pc[self] = \text{"CL"}$
$\qquad\qquad \land \text{IF } cntr[self] \leq STOP$
$\qquad\qquad\qquad \text{THEN } \land state' = \text{"Reading"}$
$\qquad\qquad\qquad\qquad\quad \land state1' = \text{"InProcess"}$
$\qquad\qquad\qquad\qquad\quad \land cntr' = [cntr \text{ EXCEPT } ![self] = cntr[self] + 1]$
$\qquad\qquad\qquad\qquad\quad \land hver' = [hver \text{ EXCEPT } ![self] = db[NodeR].ver + 1]$
$\qquad\qquad\qquad\qquad\quad \land Q' = [Q \text{ EXCEPT } ![self] = ReturnWriteQ]$
$\qquad\qquad\qquad\qquad\quad \land pc' = [pc \text{ EXCEPT } ![self] = \text{"CL1"}]$
$\qquad\qquad\qquad \text{ELSE } \land pc' = [pc \text{ EXCEPT } ![self] = \text{"Done"}]$
$\qquad\qquad\qquad\qquad\quad \land \text{UNCHANGED } \langle state, state1, cntr, hver, Q \rangle$
$\qquad\qquad \land \text{UNCHANGED } \langle FailNum, wQ, up, db \rangle$

$CL1(self) \triangleq \land pc[self] = \text{"CL1"}$
$\qquad\qquad\quad \land \text{IF } Q[self] \neq \{\}$
$\qquad\qquad\qquad \text{THEN } \land state' = \text{"Writing"}$
$\qquad\qquad\qquad\qquad\quad \land db' = [db \text{ EXCEPT } ![NodeW(Q[self])].ver = hver[self],$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad ![NodeW(Q[self])].val = cntr[self]]$
$\qquad\qquad\qquad\qquad\quad \land Q' = [Q \text{ EXCEPT } ![self] = Q[self] \setminus \{NodeW(Q[self])\}]$
$\qquad\qquad\qquad\qquad\quad \land \text{IF } Q'[self] = \{\}$
$\qquad\qquad\qquad\qquad\qquad\quad \text{THEN } \land state1' = \text{"WriteEnd"}$
$\qquad\qquad\qquad\qquad\qquad\quad \text{ELSE } \land \text{TRUE}$
$\qquad\qquad\qquad\qquad\qquad\qquad\quad \land \text{UNCHANGED } state1$
$\qquad\qquad\qquad\qquad\quad \land pc' = [pc \text{ EXCEPT } ![self] = \text{"CL1"}]$
$\qquad\qquad\qquad \text{ELSE } \land pc' = [pc \text{ EXCEPT } ![self] = \text{"CL"}]$
$\qquad\qquad\qquad\qquad\quad \land \text{UNCHANGED } \langle state, state1, db, Q \rangle$
$\qquad\qquad\quad \land \text{UNCHANGED } \langle FailNum, wQ, up, cntr, hver \rangle$

$c(self) \triangleq CL(self) \lor CL1(self)$

$NODE(self) \triangleq \land pc[self] = \text{"NODE"}$
$\qquad\qquad\qquad \land \text{IF TRUE } \land ClientNotWriting$
$\qquad\qquad\qquad\qquad \text{THEN } \land \text{IF } FailNum > 0 \land up[self] = \text{TRUE}$
$\qquad\qquad\qquad\qquad\qquad\qquad\quad \text{THEN } \land up' = [up \text{ EXCEPT } ![self] = \text{FALSE}]$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \land FailNum' = FailNum - 1$
$\qquad\qquad\qquad\qquad\qquad\qquad\quad \text{ELSE } \land \text{IF } up[self] = \text{FALSE}$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \text{THEN } \land up' = [up \text{ EXCEPT } ![self] = \text{TRUE}]$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \land FailNum' = FailNum + 1$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \text{ELSE } \land \text{TRUE}$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \land \text{UNCHANGED } \langle FailNum, up \rangle$

3

$$\qquad\qquad\qquad \wedge\ pc' = [pc \text{ EXCEPT } ![self] = \text{``NODE''}]$$
$$\text{ELSE}\quad \wedge\ pc' = [pc \text{ EXCEPT } ![self] = \text{``Done''}]$$
$$\wedge\ \text{UNCHANGED}\ \langle FailNum,\ up\rangle$$
$$\wedge\ \text{UNCHANGED}\ \langle state,\ state1,\ wQ,\ db,\ cntr,\ hver,\ Q\rangle$$

$$n(self)\ \triangleq\ NODE(self)$$

$$Next\ \triangleq\ (\exists\, self\ \in\ Clients : c(self))$$
$$\vee\ (\exists\, self\ \in\ Nodes : n(self))$$
$$\vee\ \boxed{\text{Disjunct to prevent deadlock on termination}}$$
$$((\forall\, self\ \in\ ProcSet : pc[self] = \text{``Done''})\ \wedge\ \text{UNCHANGED}\ vars)$$

$$Spec\ \triangleq\ \wedge\ Init \wedge \Box[Next]_{vars}$$
$$\wedge\ \forall\, self\ \in\ Clients : \text{WF}_{vars}(c(self))$$
$$\wedge\ \forall\, self\ \in\ Nodes : \text{WF}_{vars}(n(self))$$

$$Termination\ \triangleq\ \Diamond(\forall\, self\ \in\ ProcSet : pc[self] = \text{``Done''})$$

END TRANSLATION

*InvP* is the Invariant function which checks *db* consistency after every atomic Client WRITE ends.
$$InvP\ \triangleq\ \forall\, p\ \in\ Clients : (ClientAtomicWriteDone = \text{TRUE} \Rightarrow$$
$$(db[NodeR].ver\quad = hver[p] \wedge db[NodeR].val = cntr[p]))$$

---

\* Observation: Model checking the system with $FAILNUM = ReadQ$ and $ReadQ = WriteQ$
\*            causes it to violate the Invariant. This is evident, as, the size
\*            of $WriteQ/ReadQ$ is then not enough to overcome the node failure
\*            count in case $FAILNUM$ number of nodes fail. So say $ReadQ/WriteQ = 2$
\*            and $FAILNUM = 2$. In this case the Client process atomically Reads
\*            the highest version from 2 $ReadQ$ nodes, and fetches the $WriteQ$ to
\*            write the latest entries into. But if now the system decides to
\*            fail the two nodes which were selected as the $WriteQ$, the next
\*            invocation of the Client process will write data to FAILED nodes.
\*          - This is where the Invariant will fail because it will try to get
\*            the values from a $ReadQ$ which is not consistent with nodes to
\*            which the Client process is writing the latest updates to.
\*
\*          - The Invariant is also violated if $ReadQ < = FAILNUM$     and
\*            $WriteQ > FAILNUM$. This is because even though $WriteQ$ selects safe
\*            number of up-nodes to perform the writes, the $ReadQ$ is only
\*            returning the highest version from a subset of nodes in $WriteQ$.
\*            Thereby failing to return the true highest version that exists in
\*            the up-nodes.
\*
\*          - So, Invariant will only be satisfied if $FAILNUM < ReadQ$ and
\*            $FAILNUM < WriteQ$. This is the pre-requisite that guarantees single
\*            -copy consistency in the system.
\*
\* *teamMember*: Debaditya *Basak*, 50206177