

```

/*Replicated storage protocol with chain replication.
EXTENDS Integers, Sequences, FiniteSets, TLC
CONSTANTS N, C, STOP, FAILNUM
ASSUME  $N - FAILNUM \geq 1 \wedge STOP < 5 \wedge 0 \leq FAILNUM \wedge FAILNUM \leq 2$ 
Nodes  $\triangleq 1 \dots N$ 
Clients  $\triangleq N + 1 \dots N + C$  /*should give different ID space to Client
Procs  $\triangleq 1 \dots N + C$ 
Configurator  $\triangleq N + C + 1$  /*Configurator is unfallable

--algorithm voldchain{
  variable FailNum = FAILNUM,
    msg = [ $j \in Procs \mapsto \langle \text{"X"}, -1, -1, \text{"X"}, -1 \rangle$ ], /*default message
    up = [ $n \in Nodes \mapsto \text{TRUE}$ ],
    db = [ $n \in Nodes \mapsto [ver \mapsto -1, val \mapsto -1, cli \mapsto -1]$ ],
    /*db is single record only
    chain =  $\langle 1 \rangle$ ,
    /*chain is a sequence initially empty
    status = "Writing" ;

  define {
    UpNodes  $\triangleq \{n \in Nodes : up[n] = \text{TRUE}\}$ 
    InChain(s)  $\triangleq \exists i \in 1 \dots Len(chain) : chain[i] = s$ 
    ChainNodes  $\triangleq \{chain[j] : j \in 1 \dots Len(chain)\}$ 
    FreeUpNode  $\triangleq \text{IF } (UpNodes \setminus ChainNodes) \neq \{\} \text{ THEN}$ 
      CHOOSE  $i \in (UpNodes \setminus ChainNodes) : i > 0$ 
    ELSE 0
    GetIndex(s)  $\triangleq \text{CHOOSE } i \in 1 \dots Len(chain) : chain[i] = s$ 
    /*Assume InChain(s), returns index of s in chain
    GetNext(s)  $\triangleq chain[GetIndex(s) + 1]$ 
    /*Assume InChain(s), returns successor of s in chain
    IsUp(s)  $\triangleq up[s] = \text{TRUE}$ 
  }

  fair process (  $c \in Clients$  ) /*Client process
  variable ctr = -1, hver = -1 ;
  {
    C0: await ( $Len(chain) > 0$ ) ;
    CL: while (  $ctr \leq STOP$  ) {
      status := "Reading" ;
      CLR: hver :=  $db[chain[Len(chain)]] . ver + 1$  ;
      ctr :=  $ctr + 1$  ;
      CLW: while (  $msg[self][1] \neq \text{"ACK"} \vee$ 
         $msg[self][2] \neq hver \vee$ 
         $msg[self][3] \neq ctr \vee$ 
         $msg[self][4] \neq \text{"UPDATE"}$  ) {

```

```

        msg[chain[1]] := ⟨ "SYN", hver, ctr, "UPDATE", self ⟩ ;
    } ;
    status := "Writing" ;
    msg[self] := ⟨ "X", -1, -1, "X", -1 ⟩ ;
} ;

fair process ( n ∈ Nodes ) /*Storage node
{
    ND: while ( TRUE ) {
        either
        NM: { if ( msg[self][1] = "SYN" ∧ Len(chain) > 0 ) {
            /*react to message
            if ( self = chain[Len(chain)] ∧ msg[self][4] = "QUERY" ) {
                msg[msg[self][5]] := ⟨ "ACK", db[self].ver, db[self].val, "QUERY", self ⟩ ||
                msg[self] := ⟨ "X", -1, -1, "X", -1 ⟩ ;
            }
            else if ( self = chain[Len(chain)] ∧ msg[self][4] = "UPDATE" ) {
                db[self].ver := msg[self][2] || db[self].val := msg[self][3] ;
                msg[msg[self][5]] := ⟨ "ACK", db[self].ver, db[self].val, "UPDATE", self ⟩ ||
                msg[self] := ⟨ "X", -1, -1, "X", -1 ⟩ ;
            }
            else if ( self ≠ chain[Len(chain)] ∧
                msg[self][4] = "UPDATE" ∧ InChain(self) = TRUE ) {
                db[self].ver := msg[self][2] || db[self].val := msg[self][3] ;
                msg[GetNext(self)] :=
                ⟨ "SYN", msg[self][2], msg[self][3], "UPDATE", msg[self][5] ⟩ ||
                msg[self] := ⟨ "X", -1, -1, "X", -1 ⟩ ;
            }
        }
        } or
        NDF: {
            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 ;
                msg[self] := ⟨ "X", -1, -1, "X", -1 ⟩ ;
                FailNum := FailNum + 1 ; }
        }
    } ;
}

fair process ( p = Configurator ) /*Maintain the chain
variables un = 0 ;

```

```

{
  P: while ( TRUE ) {
    if ( Len(chain) = 0 ) {
      chain := Append(chain, FreeUpNode); }
    else if ( Len(chain) < 3 ∧ FreeUpNode > 0 ) {
      un := FreeUpNode;
      db[un].ver := db[chain[Len(chain)]] . ver
      || db[un].val := db[chain[Len(chain)]] . val
      || db[un].cli := db[chain[Len(chain)]] . cli;
      chain := Append(SelectSeq(chain, IsUp), FreeUpNode); }
    else if ( Len(chain) > 0 ) {
      chain := SelectSeq(chain, IsUp); }
  } ;
}
}

*****
BEGIN TRANSLATION
VARIABLES FailNum, msg, up, db, chain, status, pc

define statement
UpNodes  $\triangleq$  {n ∈ Nodes : up[n] = TRUE}
InChain(s)  $\triangleq$  ∃ i ∈ 1 .. Len(chain) : chain[i] = s
ChainNodes  $\triangleq$  {chain[j] : j ∈ 1 .. Len(chain)}
FreeUpNode  $\triangleq$  IF (UpNodes \ ChainNodes) ≠ {} THEN
  CHOOSE i ∈ (UpNodes \ ChainNodes) : i > 0
  ELSE 0
GetIndex(s)  $\triangleq$  CHOOSE i ∈ 1 .. Len(chain) : chain[i] = s
GetNext(s)  $\triangleq$  chain[GetIndex(s) + 1]
IsUp(s)  $\triangleq$  up[s] = TRUE

VARIABLES ctr, hver, un

vars  $\triangleq$  ⟨FailNum, msg, up, db, chain, status, pc, ctr, hver, un⟩

ProcSet  $\triangleq$  (Clients) ∪ (Nodes) ∪ {Configurator}

Init  $\triangleq$  Global variables
  ∧ FailNum = FAILNUM
  ∧ msg = [j ∈ Procs ↦ ⟨“X”, -1, -1, “X”, -1⟩]
  ∧ up = [n ∈ Nodes ↦ TRUE]
  ∧ db = [n ∈ Nodes ↦ [ver ↦ -1, val ↦ -1, cli ↦ -1]]
  ∧ chain = ⟨1⟩
  ∧ status = “Writing”
  Process c

```

$$\begin{aligned}
& \wedge ctr = [self \in Clients \mapsto -1] \\
& \wedge hver = [self \in Clients \mapsto -1] \\
& \text{Process } p \\
& \wedge un = 0 \\
& \wedge pc = [self \in ProcSet \mapsto \text{CASE } self \in Clients \rightarrow \text{"C0"} \\
& \quad \quad \quad \square \quad self \in Nodes \rightarrow \text{"ND"} \\
& \quad \quad \quad \square \quad self = Configurator \rightarrow \text{"P"}]
\end{aligned}$$

$$\begin{aligned}
C0(self) \triangleq & \wedge pc[self] = \text{"C0"} \\
& \wedge (Len(chain) > 0) \\
& \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"CL"}] \\
& \wedge \text{UNCHANGED } \langle FailNum, msg, up, db, chain, status, ctr, hver, un \rangle
\end{aligned}$$

$$\begin{aligned}
CL(self) \triangleq & \wedge pc[self] = \text{"CL"} \\
& \wedge \text{IF } ctr[self] \leq STOP \\
& \quad \text{THEN } \wedge status' = \text{"Reading"} \\
& \quad \quad \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"CLR"}] \\
& \quad \text{ELSE } \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"Done"}] \\
& \quad \quad \wedge \text{UNCHANGED } status \\
& \wedge \text{UNCHANGED } \langle FailNum, msg, up, db, chain, ctr, hver, un \rangle
\end{aligned}$$

$$\begin{aligned}
CLR(self) \triangleq & \wedge pc[self] = \text{"CLR"} \\
& \wedge hver' = [hver \text{ EXCEPT } ![self] = db[chain[Len(chain)]] . ver + 1] \\
& \wedge ctr' = [ctr \text{ EXCEPT } ![self] = ctr[self] + 1] \\
& \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"CLW"}] \\
& \wedge \text{UNCHANGED } \langle FailNum, msg, up, db, chain, status, un \rangle
\end{aligned}$$

$$\begin{aligned}
CLW(self) \triangleq & \wedge pc[self] = \text{"CLW"} \\
& \wedge \text{IF } msg[self][1] \neq \text{"ACK"} \vee \\
& \quad msg[self][2] \neq hver[self] \vee \\
& \quad msg[self][3] \neq ctr[self] \vee \\
& \quad msg[self][4] \neq \text{"UPDATE"} \\
& \quad \text{THEN } \wedge msg' = [msg \text{ EXCEPT } ![chain[1]] = \langle \text{"SYN"}, hver[self], ctr[self], \text{"UPDATE"}, se \rangle \\
& \quad \quad \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"CLW"}] \\
& \quad \quad \wedge \text{UNCHANGED } status \\
& \quad \text{ELSE } \wedge status' = \text{"Writing"} \\
& \quad \quad \wedge msg' = [msg \text{ EXCEPT } ![self] = \langle \text{"X"}, -1, -1, \text{"X"}, -1 \rangle] \\
& \quad \quad \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"CL"}] \\
& \wedge \text{UNCHANGED } \langle FailNum, up, db, chain, ctr, hver, un \rangle
\end{aligned}$$

$$c(self) \triangleq C0(self) \vee CL(self) \vee CLR(self) \vee CLW(self)$$

$$\begin{aligned}
ND(self) \triangleq & \wedge pc[self] = \text{"ND"} \\
& \wedge \vee \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"NM"}] \\
& \quad \vee \wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"NDF"}] \\
& \wedge \text{UNCHANGED } \langle FailNum, msg, up, db, chain, status, ctr, hver, un \rangle
\end{aligned}$$

$$\begin{aligned}
NM(self) &\triangleq \wedge pc[self] = \text{"NM"} \\
&\wedge \text{IF } msg[self][1] = \text{"SYN"} \wedge Len(chain) > 0 \\
&\quad \text{THEN } \wedge \text{IF } self = chain[Len(chain)] \wedge msg[self][4] = \text{"QUERY"} \\
&\quad \quad \text{THEN } \wedge msg' = [msg \text{ EXCEPT } ![msg[self][5]] = \langle \text{"ACK"}, db[self].ver, db[self] \\
&\quad \quad \quad ![self] = \langle \text{"X"}, -1, -1, \text{"X"}, -1 \rangle] \\
&\quad \quad \quad \wedge db' = db \\
&\quad \text{ELSE } \wedge \text{IF } self = chain[Len(chain)] \wedge msg[self][4] = \text{"UPDATE"} \\
&\quad \quad \text{THEN } \wedge db' = [db \text{ EXCEPT } ![self].ver = msg[self][2], \\
&\quad \quad \quad ![self].val = msg[self][3]] \\
&\quad \quad \quad \wedge msg' = [msg \text{ EXCEPT } ![msg[self][5]] = \langle \text{"ACK"}, db'[self] \\
&\quad \quad \quad \quad ![self] = \langle \text{"X"}, -1, -1, \text{"X"}, -1 \rangle] \\
&\quad \text{ELSE } \wedge \text{IF } self \neq chain[Len(chain)] \wedge \\
&\quad \quad msg[self][4] = \text{"UPDATE"} \wedge InChain(self) = \text{TRUE} \\
&\quad \quad \text{THEN } \wedge db' = [db \text{ EXCEPT } ![self].ver = msg[self][2] \\
&\quad \quad \quad ![self].val = msg[self][3]] \\
&\quad \quad \quad \wedge msg' = [msg \text{ EXCEPT } ![GetNext(self)] = \langle \text{"X"}, -1, \\
&\quad \quad \quad \quad ![self] = \langle \text{"X"}, -1, \\
&\quad \quad \quad \quad db \rangle] \\
&\quad \text{ELSE } \wedge \text{TRUE} \\
&\quad \quad \wedge \text{UNCHANGED } \langle msg, db \rangle \\
&\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"ND"}] \\
&\wedge \text{UNCHANGED } \langle FailNum, up, chain, status, ctr, hver, un \rangle \\
\\
NDF(self) &\triangleq \wedge pc[self] = \text{"NDF"} \\
&\wedge \text{IF } FailNum > 0 \wedge up[self] = \text{TRUE} \\
&\quad \text{THEN } \wedge up' = [up \text{ EXCEPT } ![self] = \text{FALSE}] \\
&\quad \quad \wedge FailNum' = FailNum - 1 \\
&\quad \quad \wedge msg' = msg \\
&\quad \text{ELSE } \wedge \text{IF } up[self] = \text{FALSE} \\
&\quad \quad \text{THEN } \wedge up' = [up \text{ EXCEPT } ![self] = \text{TRUE}] \\
&\quad \quad \quad \wedge msg' = [msg \text{ EXCEPT } ![self] = \langle \text{"X"}, -1, -1, \text{"X"}, -1 \rangle] \\
&\quad \quad \quad \wedge FailNum' = FailNum + 1 \\
&\quad \text{ELSE } \wedge \text{TRUE} \\
&\quad \quad \wedge \text{UNCHANGED } \langle FailNum, msg, up \rangle \\
&\wedge pc' = [pc \text{ EXCEPT } ![self] = \text{"ND"}] \\
&\wedge \text{UNCHANGED } \langle db, chain, status, ctr, hver, un \rangle \\
\\
n(self) &\triangleq ND(self) \vee NM(self) \vee NDF(self) \\
\\
P &\triangleq \wedge pc[Configurator] = \text{"P"} \\
&\wedge \text{IF } Len(chain) = 0 \\
&\quad \text{THEN } \wedge chain' = Append(chain, FreeUpNode) \\
&\quad \quad \wedge \text{UNCHANGED } \langle db, un \rangle \\
&\quad \text{ELSE } \wedge \text{IF } Len(chain) < 3 \wedge FreeUpNode > 0
\end{aligned}$$

```

THEN  $\wedge un' = FreeUpNode$ 
       $\wedge db' = [db \text{ EXCEPT } ![un'].ver = db[chain[Len(chain)]]].ver,$ 
       $![un'].val = db[chain[Len(chain)]]].val,$ 
       $![un'].cli = db[chain[Len(chain)]]].cli$ 
       $\wedge chain' = Append(SelectSeq(chain, IsUp), FreeUpNode)$ 
ELSE  $\wedge \text{IF } Len(chain) > 0$ 
      THEN  $\wedge chain' = SelectSeq(chain, IsUp)$ 
      ELSE  $\wedge \text{TRUE}$ 
       $\wedge chain' = chain$ 
       $\wedge \text{UNCHANGED } \langle db, un \rangle$ 
 $\wedge pc' = [pc \text{ EXCEPT } ![Configurator] = "P"]$ 
 $\wedge \text{UNCHANGED } \langle FailNum, msg, up, status, ctr, hver \rangle$ 

 $p \triangleq P$ 

 $Next \triangleq p$ 
       $\vee (\exists self \in Clients : c(self))$ 
       $\vee (\exists self \in Nodes : n(self))$ 

 $Spec \triangleq \wedge Init \wedge \square [Next]_{vars}$ 
       $\wedge \forall self \in Clients : WF_{vars}(c(self))$ 
       $\wedge \forall self \in Nodes : WF_{vars}(n(self))$ 
       $\wedge WF_{vars}(p)$ 

END TRANSLATION

 $Invariant\_1 \triangleq \forall cl \in Clients : (status = "Writing" \Rightarrow$ 
       $(hver[cl] = db[chain[Len(chain)]]].ver \wedge ctr[cl] = db[chain[Len(chain)]]].val)$ 
 $Invariant\_2 \triangleq \forall nl \in 1 .. Len(chain) - 1 : ((Len(chain) > 1) \Rightarrow (db[chain[nl]].ver \geq db[chain[nl + 1]].ver))$ 

```

(* Observations: This project is server side routing, which acted as a black box to client- systems. This is different to Project 1 where global variables were used to do client side routing. Also there were no configurator in *Project1* as all the responsibilities for data management was being done by the client. In this case (theoretically), a dedicated *PAXOS* based configurator is deployed to maintain the state of replicated chain. However in this project, configurator is not directly accessed by nodes, which posed its own challenges to define what constitutes an active tail or head node. We here, always consider current chain config while identifying head or tail.

1. As clients were not supposed to communicate with configurator for this project, deciding active head and tail was a big challenge.
2. The assumption that atleast one node will be up all the time was a huge help in calculation of head and tail.
3. We faced problems while initializing the system and came to conclusion that system should always initialize with n nodes such that $n + 1 > FAILNUM$ (At least one node should be up for system to sustain consistent behavior).
4. Initializing with 1 or two nodes in this case was violating consistency as both nodes can go down (because of *FAILNUM*) and entire service will go down and system will lose data.

5. In case of 2 clients problem arises when say *client_1* read *ver 3* and could not proceed to write. Meanwhile *client_2* kept on reading and writing till it finish up writing lets say 5. And if now *client_1* writes its data that will violate consistency for *client_2*.
6. Again after lot of test we found a case when configurator process is not getting picked; meanwhile a node went down and after a write cycle it comes back on (This causes the tail to fail and come back up with its *db* reset to default - while the configurator did not remove/re-configure nodes as of yet). Because of this one of the *invariant_1* is failing but we feel that this the limitation of this system design in the requirements.
7. A better design would be nodes communicating with configurator to get details about successors and copying etc.
8. *Voldchain* is more fail *safe(fault tolerant)* when compared with project 1, as it can still serve with single up node.
9. In *Voldchain*, tailnode is analogous to *readQ* (in project 1) and headnode is analogous to the *writeQ*. Server side routing is hence more fault tolerant with even lesser number of nodes required for maintaining single copy consistency.
10. In *Project1* failnum must be less than the size of *readQ*, but *voldchain* can work with 1 upnode as discussed above.
11. Also 2 clients scenario will fail *invariant_1* as *db* of node can be consistant with only one client at a time and not with all of them.

*)

\ * *teamMember_1*: Debaditya Basak, 50206177

\ * *teamMember_2*: Divyank Sharma, 50207091