CHAIN REPLICATION

EXTENDS *Sequences*, *Naturals*, *TLC*

CONSTANT *FirstReplica*,   first replica, the one that's up at *Init*
          *InitialChainLgth*,
             *Adr*,
             *Rep*,
             *Val*,
             *Object*,
             *Reply*(_, _, _),
             *NoReply*

VARIABLES *master*,
             *cache*,
             *data*,
             *channel*,
             *stat*

---

$NoVal \triangleq$ CHOOSE $v : v \notin Val$
$InvalidVal \triangleq$ CHOOSE $v : v \notin Val \wedge v \neq NoVal$
$NoRep \triangleq$ CHOOSE $r : r \notin Rep$
$InvalidRep \triangleq$ CHOOSE $x : x \notin (Rep \cup \{NoRep\})$

---

$WrReqMsg \triangleq$    Write requests
   $[type : \{\text{"wrReq"}, \text{"cliWrReq"}\}, adr : Adr, w : Val]$
$Messages \triangleq$
   $WrReqMsg$

$TypeInvariant \triangleq$
   $\wedge\ master \in$
       $[chains : [Object \rightarrow Seq(Rep \cup \{NoRep\})],$
           $health : [Rep \rightarrow \{\text{"dead"}, \text{"alive"}\}]]$

   $\wedge\ cache \in [Rep \rightarrow [Object \rightarrow$
           $[left : Rep \cup \{NoRep, InvalidRep\}, right : Rep \cup \{NoRep, InvalidRep\}, in\_chain : \{\text{TRUE}, \text{FALSE}\}]]]$
   $\wedge\ data \in [Rep \rightarrow [Object \rightarrow [Adr \rightarrow Val \cup \{NoVal, InvalidVal\}]]]$

   $\wedge\ channel \in [Rep \rightarrow [Object \rightarrow [in : Seq(Messages), out : Seq(Messages)]]]$

   $\wedge\ stat \in [Rep \rightarrow [phase : \{\text{"dead"}, \text{"alive"}, \text{"recover"}, \text{"reconfig"}\}]]$

$Init \triangleq$
   $\wedge\ master = [chains \mapsto [o \in Object \mapsto \langle NoRep, FirstReplica, NoRep \rangle],$
                        $health \mapsto [r \in Rep \mapsto$ IF $r = FirstReplica$ THEN $\text{"alive"}$ ELSE $\text{"dead"}]]$

$\wedge\ cache = [r \in Rep \mapsto [o \in Object \mapsto$
$\qquad\qquad\qquad$ IF $r = FirstReplica$ THEN
$\qquad\qquad\qquad\quad [left\quad \mapsto NoRep,\ right \mapsto NoRep,\ in\_chain \mapsto$ TRUE$]$
$\qquad\qquad\qquad$ ELSE $[left \mapsto InvalidRep,\ right \mapsto InvalidRep,\ in\_chain \mapsto$ FALSE$]]]$
$\wedge\ data = [r \in Rep \mapsto [o \in Object \mapsto [a \in Adr \mapsto$
$\qquad\qquad\qquad$ IF $r = FirstReplica$ THEN $NoVal$ ELSE $InvalidVal]]]$

$\wedge\ channel = [r \in Rep \mapsto [o \in Object \mapsto [in \mapsto \langle\rangle,\ out \mapsto \langle\rangle]]]$

$\wedge\ stat = [r \in Rep \mapsto [phase \mapsto$ IF $r = FirstReplica$ THEN "alive" ELSE "dead"$]]$

---

## Useful functions

$Hd(o) \triangleq$ IF $master.chains[o] = \langle NoRep,\ NoRep\rangle$ THEN
$\qquad\quad NoRep$
$\qquad$ ELSE
$\qquad\quad master.chains[o][2]$
$Tl(o) \triangleq$ IF $master.chains[o] = \langle NoRep,\ NoRep\rangle$ THEN
$\qquad\ NoRep$
$\qquad$ ELSE
$\qquad\ master.chains[o][Len(master.chains[o]) - 1]$

---

## MASTER

$DelElem(seq,\ el) \triangleq$
$\quad$ IF $(\exists\, i \in 2\,..\,(Len(seq) - 1) : seq[i] = el)$ THEN
$\qquad$ LET $idx \triangleq$ CHOOSE $i \in 1\,..\,Len(seq) : seq[i] = el$ IN
$\qquad\quad SubSeq(seq,\ 1,\ idx - 1) \circ SubSeq(seq,\ idx + 1,\ Len(seq))$
$\quad$ ELSE $\ seq$
$InSeq(seq,\ el) \triangleq$
$\quad \exists\, i \in 2\,..\,(Len(seq) - 1) : seq[i] = el$

$\_RemoveRep(r) \triangleq$
$\ $ LET $alivepred(r1) \triangleq r1 \neq NoRep \wedge stat[r1].phase \neq$ "recover"
$\ $ IN
$\quad \wedge\ master.health[r] \neq$ "dead"
$\quad \wedge\ \forall\, o \in Object :$
$\qquad InSeq(master.chains[o],\ r) \Rightarrow Len(SelectSeq(master.chains[o],\ alivepred)) > 1$ <span style="background:#ccc">the two *NoRep* plus at leas</span>
$\quad \wedge\ master' = [health \mapsto [master.health$ EXCEPT $![r] =$ "dead"$],$
$\qquad\qquad\qquad chains \mapsto [o \in Object \mapsto$
$\qquad\qquad\qquad\quad DelElem(master.chains[o],\ r)]]$
$\quad \wedge\ stat' = [stat$ EXCEPT $![r].phase =$ "dead"$]$ <span style="background:#ccc">ASSUMPTION : FAIL-STOP, or no partitions</span>
$\qquad\qquad\qquad\qquad\qquad\qquad$ <span style="background:#ccc">I do it in-line here, *s.t.* I avoid transition</span>

$RemoveRep \triangleq$
$\ \exists\, r \in Rep :$

2

$\land$ $\_RemoveRep(r)$

$InitReplica(r,\ o)\ \triangleq$
   $\land\ stat' = [stat\ \text{EXCEPT}\ ![r].phase = \text{"recover"}]$   revive the replica (just $s.t.$ I don't have one more transition)
   $\land\ channel' = [channel\ \text{EXCEPT}\ ![r][o] = [in \mapsto \langle\rangle,\ out \mapsto \langle\rangle]]$
   $\land\ data' = [data\ \text{EXCEPT}\ ![r][o] = [a \in Adr \mapsto InvalidVal]]$
   $\land\ cache' = [cache\ \text{EXCEPT}\ ![r][o] = [cache[r][o]\ \text{EXCEPT}\ !.in\_chain = \text{FALSE}]]$   make it

$\_AddRep(r,\ o)\ \triangleq$
   $\land\ \neg InSeq(master.chains[o],\ r)$   $r$ is not yet in chain for $o$
   $\land\ master' = [master\ \text{EXCEPT}\ !.chains[o] = SubSeq(@,\ 1,\ Len(@)-1) \circ \langle r,\ NoRep\rangle,$
                                  $!.health[r]\ = \text{"alive"}]$   append $r$ to the chain
   $\land\ InitReplica(r,\ o)$

$AddRep\ \triangleq$
  $\exists\,r \in Rep:$
    $\land\ master.health[r] = \text{"dead"}$   the replica is not dead or restarting

    Everybody have $ACKed$ that they know that $r$ is dead
    $\land\ \forall\,r1 \in Rep: \forall\,o \in Object:$
      $cache[r1][o].in\_chain \land stat[r1].phase \neq \text{"dead"}$
      $\Rightarrow cache[r1][o].right \neq r \land cache[r1][o].left \neq r$

    There is an object to which I can add it.
    $\land\ \exists\,o \in Object:$
      $\land\ stat[Tl(o)].phase \neq \text{"recover"}$   the tail finished recovery
      $\land\ \_AddRep(r,\ o)$

$MasterActions\ \triangleq$
  $\lor\ RemoveRep \land \text{UNCHANGED}\ \langle cache,\ channel,\ data\rangle$   either remove a replica (kill it)
  $\lor\ AddRep$   $\land \text{UNCHANGED}\ \langle cache,\ stat,\ channel,\ data\rangle$    or add it to some chain

---

**REPLICA**

$\_RecvUpdateConfig(r,\ o)\ \triangleq$   replica rreceives an update message on the configuration from master
  $\land\ InSeq(master.chains[o],\ r)$   $r$ should currently be in the chain for $o$
                                  late $updateConfig$ messages are not supported $TODO$
  $\land\ \text{LET}\ idx\ \triangleq\ \text{CHOOSE}\ i \in 2\,..\,(Len(master.chains[o])-1):\ master.chains[o][i] = r\,\text{IN}$
    $\land$   $\lor\ cache[r][o].in\_chain = \text{FALSE}$   $r$ doesn't know it's in chain
       OR $r$ doesn't know its right left/right neighbors
      $\lor\ cache[r][o].left \neq master.chains[o][idx-1]$
      $\lor\ cache[r][o].right \neq master.chains[o][idx+1]$

    $\land\ cache' = [cache\ \text{EXCEPT}\ ![r][o].left = master.chains[o][idx-1],$
                                $![r][o].right = master.chains[o][idx+1],$
                                $!\,[r][o].in\_chain = \text{TRUE}]$

3

$\wedge$ IF ( $\wedge$ $cache[r][o].right \neq cache'[r][o].right$   right neighbor has changed
        $\wedge$ $stat[r].phase \neq$ "recover" )   replica is not recovering
    THEN
        mark state as $reconfig$
        $stat' = [stat$ EXCEPT $![r].phase =$ "reconfig"]
    ELSE UNCHANGED $stat$
$\wedge$ UNCHANGED $\langle master,\ data,\ channel \rangle$

$RecvUpdateConfig \triangleq$
  $\exists\, r \in Rep :$
    $\wedge$ $stat[r].phase \in \{$ "alive", "reconfig", "recover" $\}$   $r$ is alive or reconfiguring, but not dead or recovering
    $\wedge\, \exists\, o \in Object : \_RecvUpdateConfig(r,\ o)$

$IsNotClientReq(m) \triangleq$
  $m.type \neq$ "cliWrReq"

$FwdWrite(r,\ o,\ q) \triangleq$   while in $reconfig$ state, fwd a part of the writes already processed at $r$
                                to $q$,. who's $r$'s right neighbor
LET $index\_unknown \triangleq Len(SelectSeq(channel[q][o].in,\ IsNotClientReq)) + Len(channel[q][o].out) + 1$ IN
  $\wedge$ $Len(channel[r][o].out) \geq index\_unknown$     there are still writes to fwd
  $\wedge$ $channel' = [channel$ EXCEPT $![q][o].in = @ \circ \langle channel[r][o].out[index\_unknown] \rangle]$
  $\wedge$ $NoReply$

$SendAcks(r,\ o) \triangleq$   while in $reconfig$ state, $ACK$ ALL unacked committed writes to the client
  $\wedge Len(channel[r][o].out) > 0$
  $\wedge$ LET $wmsg \triangleq Head(channel[r][o].out)$ IN
      $Reply(r,$ "wr", $[object \mapsto o,\ adr \mapsto wmsg.adr,\ val \mapsto wmsg.w,\ ack \mapsto$ "ok"$])$
  $\wedge channel' = [r1 \in Rep \mapsto [channel[r1]$ EXCEPT $![o] =$
                      $[in \mapsto channel[r1][o].in,$
                       $out \mapsto$ IF $Len(channel[r1][o].out) > 0$ THEN $Tail(channel[r1][o].out)$
                                ELSE $\langle\rangle]]]$

$ResendNext(r,\ o) \triangleq$   replica $r$ reconciles with its right neighbor by re-sending to it the "sent" writes
  $\wedge stat[r].phase =$ "reconfig"   reconfiguring state
  $\wedge (cache[r][o].right = NoRep \Rightarrow SendAcks(r,\ o))$   send acks to client
  $\wedge (cache[r][o].right \neq NoRep \Rightarrow FwdWrite(r,\ o,\ cache[r][o].right))$
  $\wedge$ UNCHANGED $\langle stat,\ data,\ master,\ cache \rangle$

$FinishReconfig(r,\ o) \triangleq$   finish a reconfiguration state
  $\wedge stat[r].phase =$ "reconfig"
  no more things to fwd
  $\wedge cache[r][o].right = NoRep \Rightarrow Len(channel[r][o].out) = 0$
  $\wedge cache[r][o].right \neq NoRep \Rightarrow Len(channel[r][o].out) =$
      $Len(SelectSeq(channel[cache[r][o].right][o].in,\ IsNotClientReq))$
      $+ Len(channel[cache[r][o].right][o].out)$

$\land stat' = [stat \text{ EXCEPT } ![r].phase = \text{"alive"}]$
$\land \text{UNCHANGED } \langle data, channel, master, cache \rangle$

$Reconfiguration \triangleq$ sends reply
  $\lor \exists r \in Rep, o \in Object :$
      $\lor ResendNext(r, o)$ either advance reconfiguration
      $\lor FinishReconfig(r, o) \land NoReply$ or finish reconfiguration
  $\lor RecvUpdateConfig \land NoReply$

$Reconcile(r, o, a) \triangleq$ Replica $r$ helps its recovering right neighbor to reconcile
  $TODO$: missing some transitions in the real system!
  LET $IsForAddress(m) \triangleq (m.adr = a)$ IN
  $\land \; stat[r].phase = \text{"alive"}$
  $\land \; cache[r][o].right \neq NoRep$
  $\land \; stat[cache[r][o].right].phase = \text{"recover"}$ the right neighbor is in recovery state
  $\land \; data[cache[r][o].right][o][a] \neq data[r][o][a]$ hasn't yet reconciled for this address

  $\land \; SelectSeq(channel[r][o].out, IsForAddress) = \langle\rangle$ no messages already in the "sent" that await fwding
      or have already been fwded to the right neighbor
      Such updates already contain the value that I would send now
      and some other updates, as well.

  send it the value at address a
  $\land data' = [data \text{ EXCEPT } ![cache[r][o].right][o][a] = data[r][o][a]]$

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

$FinishReconcile(r, o) \triangleq$ $r$ knows it has finished reconciliation and thus it can now start
      serving requests when its data is up-to-date and
      it knows its correct neighbors.
  $\land stat[r].phase = \text{"recover"}$ $r$ is recovering

  Reconciliation is finished when the left neighbor has sent $r$ all the data
  this simulates a FINISH RECONCILE message received from left neighbor
  $\land \forall a \in Adr : data[r][o][a] \neq InvalidVal$ the data was fully reconciled

  By the time I finish reconciliation, the new replica must have read the config from the master
  at least once. Otherwise, it can't start working, b/c it doesn't know its place in the chain
  $\land cache[r][o].in\_chain$

  $\land stat' = [stat \text{ EXCEPT } ![r].phase = \text{"alive"}]$ replica will start to respond to queries
  $\land \text{UNCHANGED } \langle master, cache, data, channel \rangle$

$Recovery \triangleq$ sends reply
  $\exists r \in Rep, o \in Object :$
    $\lor \exists a \in Adr : Reconcile(r, o, a)$
    $\lor FinishReconcile(r, o)$

$ReplicaDeath \triangleq$ FALSE  avoid this transition b/c it's fail-stop anyway

$HdlWrite(r,\ o,\ wmsg,\ ch) \triangleq$
  $\wedge\ cache[r][o].right \neq NoRep$  not tail
  fwd the write to the right neighbor
  $\wedge\ channel' = [ch$ EXCEPT $![r][o].out = @ \circ \langle wmsg \rangle,$
                           $![cache[r][o].right][o].in = @ \circ \langle wmsg \rangle]$

$FinishWrite(r,\ o,\ wmsg,\ ch) \triangleq$  sends a reply
  $\wedge\ cache[r][o].right = NoRep$  no right neighbor, so $r$ believes it's tail

  reply to client
  $\wedge\ Reply(r,\ \text{``wr''},\ [object \mapsto o,\ adr \mapsto wmsg.adr,\ val \mapsto wmsg.w,\ ack \mapsto \text{``ok"}])$
  $ACK$ to all replicas, to remove their – ATOMICALLY, b/c it's only for performance
  $\wedge\ channel' = [r1 \in Rep \mapsto [channel[r1]$ EXCEPT $![o] =$
                    $[in \mapsto ch[r1][o].in,$
                    $out \mapsto$ IF $Len(ch[r1][o].out) > 0$ THEN $Tail(ch[r1][o].out)$
                      ELSE $\langle\rangle]]]$

$ProcessWrite(r,\ o,\ wmsg,\ ch) \triangleq$  replica $r$ processes a message in its incoming $FIFO$ – Sends reply
  $\wedge\ stat[r].phase \in \{\text{``alive''},\ \text{``recover''}\}$  not reconfiguring, or dead
  $\wedge\ cache[r][o].in\_chain$  $r$ believes itself in the chain for $o$

  $\wedge\ data' = [data$ EXCEPT $![r][o][wmsg.adr] = wmsg.w]$  commit to disk
  $\wedge\ \vee\ HdlWrite\ \ (r,\ o,\ wmsg,\ ch) \wedge NoReply$
    $\vee\ FinishWrite(r,\ o,\ wmsg,\ ch)$
  $\wedge$ UNCHANGED $\langle master,\ cache,\ stat \rangle$

$\_ProcessMsg(r,\ o) \triangleq$  sends reply
  $\wedge\ channel[r][o].in \neq \langle\rangle$  the channel is not empty
  $\wedge\ stat[r].phase \in \{\text{``alive''},\ \text{``recover''}\}$

  $\wedge$ LET $wmsg \triangleq Head(channel[r][o].in)$
        $ch \triangleq [channel$ EXCEPT $![r][o].in = Tail(@)]$
    IN
       $\wedge\ wmsg.type = \text{``wrReq''} \Rightarrow ProcessWrite(r,\ o,\ wmsg,\ ch)$
       $\wedge\ wmsg.type = \text{``cliWrReq''} \Rightarrow$
            (IF $(cache[r][o].in\_chain \wedge cache[r][o].left = NoRep)$  check $r$ is head
              THEN
                  $ProcessWrite(r,\ o,\ [wmsg$ EXCEPT $!.type = \text{``wrReq''}],\ ch)$  go ahead w/ the write
              ELSE  $r$ is not head, so just drop the writ
                  $\wedge\ channel' = ch$
                  $\wedge$ UNCHANGED $\langle master,\ cache,\ stat,\ data \rangle \wedge NoReply)$

$ProcessMsg \triangleq$
  $\exists\ r \in Rep,\ o \in Object : \_ProcessMsg(r,\ o)$

$HeadWrite(r,\ o,\ a,\ v) \triangleq$  sends reply

$\land\ cache[r][o].left = NoRep$  $r$ believes itself head
The line below is a hacky solution to having a 0-phase write and still being mappable to $SS$
It fails a write that goes to a head when the head is the only one in the chain
$\land\ cache[r][o].right \neq NoRep$

$\land\ ProcessWrite(r,\ o,\ [type \mapsto \text{``wrReq''},\ adr \mapsto a,\ w \mapsto v],\ channel)$

$TailRead(r,\ o,\ a)\ \triangleq$  sends reply
    $\land\ stat[r].phase \in \{\ \text{``alive''}\ ,\ \text{``reconfig''}\}$  not recovering, reconfiguring, or dead
    $\land\ cache[r][o].in\_chain$  $r$ believes itself in the chain for $o$
    $\land\ cache[r][o].right = NoRep$

    $\land\ Reply(r,\ \text{``rd''},\ [object \mapsto o,\ adr \mapsto a,\ val \mapsto data[r][o][a],\ ack \mapsto \text{``ok''}])$
    $\land\ \text{UNCHANGED}\ \langle data,\ master,\ cache,\ channel,\ stat\rangle$

$ReplicaActions\ \triangleq$
    $\lor\ Reconfiguration$
    $\lor\ Recovery$
    $\lor\ ReplicaDeath$

    $\lor\ ProcessMsg$

---

**CLIENT**

$\_CliRead(o,\ a)\ \triangleq$
    $\exists\ r \in Rep:$
       $TailRead(r,\ o,\ a)$
$CliRead\ \triangleq$
    $\exists\ o \in Object,\ a \in Adr: \_CliRead(o,\ a)$

$\_CliWrite(o,\ a,\ v)\ \triangleq$
    $\land\ \exists\ r \in Rep:$
        the line below is more of a hack $s.t.$ I don't add a huge $\#$ of states
        $HeadWrite(r,\ o,\ a,\ v)$
        Uncomment below and comment above to be most fair
        $channel' = [channel\ \text{EXCEPT}\ ![r][o].in = @ \circ \langle [type \mapsto \text{``cliWrReq''},\ adr \mapsto a,\ w \mapsto v]\rangle]$
    $\land\ \text{UNCHANGED}\ \langle master,\ cache,\ data,\ stat\rangle$
$CliWrite\ \triangleq$
    $\exists\ o \in Object,\ a \in Adr,\ v \in Val:$
       $\_CliWrite(o,\ a,\ v)$

$ClientActions\ \triangleq$
    $\lor\ CliRead$
    $\lor\ CliWrite \land NoReply$

---

Full specification of the Chain replication system

$Next \triangleq$
    $\lor MasterActions$
    $\lor ReplicaActions$
    $\lor ClientActions$

Full spec:

$chainvars \triangleq \langle master,\ cache,\ data,\ channel,\ stat \rangle$
$Spec \triangleq Init \land \Box[Next]_{chainvars}$

Invariants

$includedInSeq(smallseq,\ bigseq) \triangleq$
  $\forall\, i \in 1\,..\,Len(smallseq):$
     $smallseq[i] = bigseq[i]$
$UpdatePropagation \triangleq$  Update propagation invariant from paper
   $\forall\, o \in Object:$
    $\forall\, i \in (3\,..\,(Len(master.chains[o]) - 1)):$
      $includedInSeq(channel[master.chains[o][i]][o].out,\ channel[master.chains[o][i-1]][o].out)$

$AllInvariants \triangleq$
    $\land TypeInvariant$
    $\land UpdatePropagation$

Theorem

THEOREM $Spec \Rightarrow \Box AllInvariants$