

---

MODULE *SemiSync*

---

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

CONSTANT *Client*, *Replica*, *nil*

VARIABLE

*zk\_leader*, *zk\_epoch*, *zk\_leader\_epoch*, *zk\_status*,  
*old\_leaders*, *zk\_catchup\_index*,  
*db*, *db\_leader*, *db\_replicated*, *db\_epoch*, *db\_status*,  
*next\_req*, *client\_leader*, *client\_success*,  
*pending*, *pending\_db*, *client\_leader\_epoch*,  
*healer\_status*, *healer\_epoch*, *healer\_replicas*

*vars*  $\triangleq$  {  
*zk\_leader*, *zk\_epoch*, *zk\_leader\_epoch*, *zk\_status*,  
*old\_leaders*, *zk\_catchup\_index*,  
*db*, *db\_leader*, *db\_replicated*, *db\_epoch*, *db\_status*,  
*next\_req*, *client\_leader*, *client\_success*,  
*pending*, *pending\_db*, *client\_leader\_epoch*,  
*healer\_status*, *healer\_epoch*, *healer\_replicas*  
}

*zk\_vars*  $\triangleq$  {*zk\_leader*, *zk\_epoch*, *zk\_leader\_epoch*, *zk\_status*,  
*old\_leaders*, *zk\_catchup\_index*}

*db\_vars*  $\triangleq$  {*db*, *db\_leader*, *db\_replicated*, *db\_epoch*, *db\_status*}

*client\_vars*  $\triangleq$  {  
*next\_req*, *client\_leader*, *client\_success*,  
*pending*, *pending\_db*, *client\_leader\_epoch*}

*healer\_vars*  $\triangleq$  {*healer\_status*, *healer\_epoch*, *healer\_replicas*}

*max\_next\_req*  $\triangleq$  4

*max\_change\_leader*  $\triangleq$  4

*ReqSet*  $\triangleq$  60 .. (60 + *max\_next\_req*)

*Epoch*  $\triangleq$  0 .. 20

*NullReqSet*  $\triangleq$  *ReqSet*  $\cup$  {*nil*}

*NullReplica*  $\triangleq$  *Replica*  $\cup$  {*nil*}

*LogOffset*  $\triangleq$  0 .. 20

*NullLogOffset*  $\triangleq$  *LogOffset*  $\cup$  {*nil*}

*Range*(*f*)  $\triangleq$  {*f*[*x*] : *x*  $\in$  DOMAIN *f*}

$replication\_factor \triangleq 2$

$Quorum \triangleq \{x \in \text{SUBSET } Replica : \text{Cardinality}(x) = replication\_factor\}$

$TypeOK \triangleq$

$\wedge zk\_leader \in Replica$   
 $\wedge zk\_epoch \in 1 \dots 30$   
 $\wedge zk\_leader\_epoch \in 1 \dots max\_change\_leader$   
 $\wedge zk\_status \in \{ "Normal", "ChangingLeader", "WaitReplicaLog" \}$   
 $\wedge old\_leaders \subseteq Replica$   
 $\wedge zk\_catchup\_index \in LogOffset$   
  
 $\wedge db \in [Replica \rightarrow Seq(ReqSet)]$   
 $\wedge db\_leader \in [Replica \rightarrow Replica]$   
 $\wedge db\_replicated \in [Replica \rightarrow [Replica \rightarrow LogOffset]]$   
 $\wedge db\_epoch \in [Replica \rightarrow Epoch]$   
 $\wedge db\_status \in [Replica \rightarrow \{ "Writable", "Replica", "Frozen" \}]$   
  
 $\wedge next\_req \in ReqSet$   
 $\wedge client\_leader \in [Client \rightarrow Replica]$   
 $\wedge client\_success \in [Client \rightarrow Seq(ReqSet)]$   
 $\wedge pending \in [Client \rightarrow NullReqSet]$   
 $\wedge pending\_db \in [Client \rightarrow NullReplica]$   
 $\wedge client\_leader\_epoch \in [Client \rightarrow Epoch]$   
  
 $\wedge healer\_status \in \{ "Init", "UpdatingLeader", "WaitReplica" \}$   
 $\wedge healer\_epoch \in Epoch$   
 $\wedge healer\_replicas \in [Replica \rightarrow NullLogOffset]$

$Init \triangleq$

$\wedge zk\_leader \in Replica$   
 $\wedge zk\_epoch = 1$   
 $\wedge zk\_leader\_epoch = 1$   
 $\wedge zk\_status = "Normal"$   
 $\wedge old\_leaders = \{ \}$   
 $\wedge zk\_catchup\_index = 0$   
  
 $\wedge db = [r \in Replica \mapsto \langle \rangle]$   
 $\wedge db\_leader = [r \in Replica \mapsto zk\_leader]$   
 $\wedge db\_replicated = [r \in Replica \mapsto [r1 \in Replica \mapsto 0]]$   
 $\wedge db\_epoch = [r \in Replica \mapsto zk\_epoch]$   
 $\wedge db\_status = [r \in Replica \mapsto \text{IF } zk\_leader = r \text{ THEN "Writable" ELSE "Replica"}]$   
  
 $\wedge next\_req = 60$   
 $\wedge client\_leader = [c \in Client \mapsto zk\_leader]$   
 $\wedge client\_success = [c \in Client \mapsto \langle \rangle]$

$$\begin{aligned}
& \wedge \text{pending} = [c \in \text{Client} \mapsto \text{nil}] \\
& \wedge \text{pending\_db} = [c \in \text{Client} \mapsto \text{nil}] \\
& \wedge \text{client\_leader\_epoch} = [c \in \text{Client} \mapsto \text{zk\_leader\_epoch}] \\
& \wedge \text{healer\_status} = \text{"init"} \\
& \wedge \text{healer\_epoch} = 1 \\
& \wedge \text{healer\_replicas} = [r \in \text{Replica} \mapsto \text{nil}]
\end{aligned}$$

$$\begin{aligned}
\text{StartRequest}(c) & \triangleq \\
& \wedge \text{pending}[c] = \text{nil} \\
& \wedge \text{next\_req} < 60 + \text{max\_next\_req} \\
& \wedge \text{db\_status}[\text{client\_leader}[c]] = \text{"Writable"} \\
& \wedge \text{next\_req}' = \text{next\_req} + 1 \\
& \wedge \text{pending}' = [\text{pending} \text{ EXCEPT } ![c] = \text{next\_req}'] \\
& \wedge \text{pending\_db}' = [\text{pending\_db} \text{ EXCEPT } ![c] = \text{client\_leader}[c]] \\
& \wedge \text{LET } \text{leader} \triangleq \text{client\_leader}[c] \text{ IN} \\
& \quad \wedge \text{db}' = [\text{db} \text{ EXCEPT } ![\text{leader}] = \text{Append}(@, \text{next\_req}')] \\
& \quad \wedge \text{db\_replicated}' = [ \\
& \quad \quad \text{db\_replicated} \text{ EXCEPT } ![\text{leader}][\text{leader}] = \text{Len}(\text{db}'[\text{leader}])] \\
& \wedge \text{UNCHANGED } \langle \text{client\_leader}, \text{client\_success}, \text{client\_leader\_epoch} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{db\_leader}, \text{db\_epoch}, \text{db\_status} \rangle \\
& \wedge \text{UNCHANGED } \text{zk\_vars} \\
& \wedge \text{UNCHANGED } \text{healer\_vars}
\end{aligned}$$

$$\begin{aligned}
\text{Replicate}(r) & \triangleq \\
& \wedge r \neq \text{db\_leader}[r] \\
& \wedge \text{db\_status}[r] = \text{"Replica"} \\
& \wedge \text{LET } \text{leader\_data} \triangleq \text{db}[\text{db\_leader}[r]] \\
& \quad \text{new\_len} \triangleq \text{Len}(\text{db}[r]) + 1 \\
& \quad \text{leader} \triangleq \text{db\_leader}[r] \\
& \text{ IN} \\
& \quad \wedge \text{Len}(\text{db}[r]) < \text{Len}(\text{leader\_data}) \\
& \quad \wedge \text{db}' = [\text{db} \text{ EXCEPT } ![r] = \text{Append}(@, \text{leader\_data}[\text{new\_len}])] \\
& \quad \wedge \text{db\_replicated}' = [\text{db\_replicated} \text{ EXCEPT } ![\text{leader}][r] = \text{new\_len}] \\
& \wedge \text{UNCHANGED } \langle \text{db\_leader}, \text{db\_epoch}, \text{db\_status} \rangle \\
& \wedge \text{UNCHANGED } \text{client\_vars} \\
& \wedge \text{UNCHANGED } \text{zk\_vars} \\
& \wedge \text{UNCHANGED } \text{healer\_vars}
\end{aligned}$$

$$\text{new\_repl} \triangleq [r \in \text{Replica} \mapsto 0]$$

$$\begin{aligned}
\text{initReplicated}(r) & \triangleq \\
& \wedge \text{db\_replicated}' = [ \\
& \quad \text{db\_replicated} \text{ EXCEPT } ![r] = [
\end{aligned}$$

$new\_repl \text{ EXCEPT } ![r] = Len(db[r])$   
 $]]$

$dbStatusFromZK(r) \triangleq$   
 IF  $zk\_status \in \{ \text{"Normal"}, \text{"WaitReplicaLog"} \} \wedge \neg(r \in old\_leaders)$   
   THEN IF  $zk\_leader = r$   
     THEN "Writable"  
     ELSE "Replica"  
   ELSE "Frozen"

$DBUpdateLeader(r) \triangleq$   
 $\wedge db\_epoch[r] < zk\_epoch$   
 $\wedge db\_epoch' = [db\_epoch \text{ EXCEPT } ![r] = zk\_epoch]$   
 $\wedge db\_leader' = [db\_leader \text{ EXCEPT } ![r] = zk\_leader]$   
 $\wedge db\_status' = [db\_status \text{ EXCEPT } ![r] = dbStatusFromZK(r)]$   
 $\wedge$  IF  $db\_leader[r] \neq zk\_leader$   
   THEN  $initReplicated(r)$   
   ELSE UNCHANGED  $db\_replicated$   
 $\wedge$  UNCHANGED  $\langle db \rangle$   
 $\wedge$  UNCHANGED  $zk\_vars$   
 $\wedge$  UNCHANGED  $client\_vars$   
 $\wedge$  UNCHANGED  $healer\_vars$

$minOfSet(S) \triangleq \text{CHOOSE } x \in S : \forall x1 \in S : x \leq x1$

$replicatedSet(r, Q) \triangleq \{ db\_replicated[r][r1] : r1 \in Q \}$

$minReplicate(r, Q) \triangleq minOfSet(replicatedSet(r, Q))$

$DBResponse(c) \triangleq$   
 $\wedge pending[c] \neq nil$   
 $\wedge$  LET  $leader \triangleq pending\_db[c]$  IN  
    $\wedge \exists index \in \text{DOMAIN } db[leader], Q \in Quorum :$   
      $\wedge pending[c] = db[leader][index]$   
      $\wedge index \leq minReplicate(leader, Q)$   
 $\wedge pending' = [pending \text{ EXCEPT } ![c] = nil]$   
 $\wedge pending\_db' = [pending\_db \text{ EXCEPT } ![c] = nil]$   
 $\wedge client\_success' = [client\_success \text{ EXCEPT } ![c] = Append(@, pending[c])]$   
 $\wedge$  UNCHANGED  $db\_vars$   
 $\wedge$  UNCHANGED  $next\_req$   
 $\wedge$  UNCHANGED  $\langle client\_leader, client\_leader\_epoch \rangle$   
 $\wedge$  UNCHANGED  $zk\_vars$   
 $\wedge$  UNCHANGED  $healer\_vars$

$$\begin{aligned}
\text{ClientUpdateLeader}(c) &\triangleq \\
&\wedge \text{client\_leader\_epoch}[c] < \text{zk\_leader\_epoch} \\
&\wedge \text{client\_leader\_epoch}' = [\text{client\_leader\_epoch} \text{ EXCEPT } ![c] = \text{zk\_leader\_epoch}] \\
&\wedge \text{client\_leader}' = [\text{client\_leader} \text{ EXCEPT } ![c] = \text{zk\_leader}] \\
&\wedge \text{pending}' = [\text{pending} \text{ EXCEPT } ![c] = \text{nil}] \\
&\wedge \text{pending\_db}' = [\text{pending\_db} \text{ EXCEPT } ![c] = \text{nil}] \\
&\wedge \text{UNCHANGED } \text{zk\_vars} \\
&\wedge \text{UNCHANGED } \text{healer\_vars} \\
&\wedge \text{UNCHANGED } \text{db\_vars} \\
&\wedge \text{UNCHANGED } \langle \text{next\_req}, \text{client\_success} \rangle
\end{aligned}$$

$$\begin{aligned}
\text{ReadyToChangeZKLeader} &\triangleq \\
&\wedge \text{zk\_leader\_epoch} < \text{max\_change\_leader} \\
&\wedge \text{Cardinality}(\text{Replica} \setminus \text{old\_leaders}) > \text{replication\_factor} \\
&\wedge \text{zk\_status} = \text{"Normal"} \\
&\wedge \text{zk\_status}' = \text{"ChangingLeader"} \\
&\wedge \text{zk\_epoch}' = \text{zk\_epoch} + 1 \\
&\wedge \text{old\_leaders}' = \text{old\_leaders} \cup \{\text{zk\_leader}\} \\
&\wedge \text{UNCHANGED } \langle \text{zk\_leader}, \text{zk\_leader\_epoch}, \text{zk\_catchup\_index} \rangle \\
&\wedge \text{UNCHANGED } \text{client\_vars} \\
&\wedge \text{UNCHANGED } \text{db\_vars} \\
&\wedge \text{UNCHANGED } \text{healer\_vars}
\end{aligned}$$

$$\begin{aligned}
\text{zkStatusToHealerStatus} &\triangleq \\
&\text{IF } \text{zk\_status} = \text{"ChangingLeader"} \\
&\quad \text{THEN } \text{"UpdatingLeader"} \\
&\quad \text{ELSE IF } \text{zk\_status} = \text{"WaitReplicaLog"} \\
&\quad \quad \text{THEN } \text{"WaitReplica"} \\
&\quad \quad \text{ELSE } \text{"Init"}
\end{aligned}$$

$$\begin{aligned}
\text{HealerUpdateState} &\triangleq \\
&\wedge \text{healer\_epoch} < \text{zk\_epoch} \\
&\wedge \text{healer\_epoch}' = \text{zk\_epoch} \\
&\wedge \text{healer\_replicas}' = [r \in \text{Replica} \mapsto \text{nil}] \\
&\wedge \text{healer\_status}' = \text{zkStatusToHealerStatus} \\
&\wedge \text{UNCHANGED } \text{zk\_vars} \\
&\wedge \text{UNCHANGED } \text{client\_vars} \\
&\wedge \text{UNCHANGED } \text{db\_vars}
\end{aligned}$$

$$\begin{aligned}
\text{HealerGetDBLog}(r) &\triangleq \\
&\wedge \text{healer\_status} = \text{"UpdatingLeader"}
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{healer\_replicas}[r] = \text{nil} \\
& \wedge \neg(r \in \text{old\_leaders}) \\
& \wedge \text{db\_epoch}[r] = \text{healer\_epoch} \\
& \wedge \text{healer\_replicas}' = [\text{healer\_replicas} \text{ EXCEPT } ![r] = \text{Len}(\text{db}[r])] \\
& \wedge \text{UNCHANGED } \langle \text{healer\_status}, \text{healer\_epoch} \rangle \\
& \wedge \text{UNCHANGED } \text{db\_vars} \\
& \wedge \text{UNCHANGED } \text{client\_vars} \\
& \wedge \text{UNCHANGED } \text{zk\_vars}
\end{aligned}$$

$$\text{collectedDB} \triangleq \{r \in \text{Replica} : \text{healer\_replicas}[r] \neq \text{nil}\}$$

$$\begin{aligned}
\text{HealerUpdateLeader} & \triangleq \\
& \wedge \text{healer\_status} = \text{"UpdatingLeader"} \\
& \wedge \text{healer\_epoch} = \text{zk\_epoch} \\
& \wedge \text{Cardinality}(\text{collectedDB}) \geq \text{replication\_factor} \\
& \wedge \exists r \in \text{collectedDB} : \\
& \quad \wedge \forall r1 \in \text{collectedDB} : \text{healer\_replicas}[r] \geq \text{healer\_replicas}[r1] \\
& \quad \wedge \text{zk\_leader}' = r \\
& \quad \wedge \text{zk\_catchup\_index}' = \text{healer\_replicas}[r] \\
& \wedge \text{zk\_status}' = \text{"WaitReplicaLog"} \\
& \wedge \text{zk\_epoch}' = \text{zk\_epoch} + 1 \\
& \wedge \text{zk\_leader\_epoch}' = \text{zk\_leader\_epoch} + 1 \\
& \wedge \text{UNCHANGED } \text{old\_leaders} \\
& \wedge \text{UNCHANGED } \text{healer\_vars} \\
& \wedge \text{UNCHANGED } \text{db\_vars} \\
& \wedge \text{UNCHANGED } \text{client\_vars}
\end{aligned}$$

$$\begin{aligned}
\text{HealerUpdateToNormal} & \triangleq \\
& \wedge \text{healer\_status} = \text{"WaitReplica"} \\
& \wedge \text{zk\_status} = \text{"WaitReplicaLog"} \\
& \wedge \text{healer\_epoch} = \text{zk\_epoch} \\
& \wedge \exists Q \in \text{Quorum} : \\
& \quad \wedge \neg(\text{old\_leaders} \subseteq Q) \\
& \quad \wedge \forall r \in Q : \text{Len}(\text{db}[r]) \geq \text{zk\_catchup\_index} \\
& \wedge \text{zk\_status}' = \text{"Normal"} \\
& \wedge \text{zk\_epoch}' = \text{zk\_epoch} + 1 \\
& \wedge \text{UNCHANGED } \langle \text{old\_leaders}, \text{zk\_leader\_epoch}, \text{zk\_leader}, \text{zk\_catchup\_index} \rangle \\
& \wedge \text{UNCHANGED } \text{healer\_vars} \\
& \wedge \text{UNCHANGED } \text{db\_vars} \\
& \wedge \text{UNCHANGED } \text{client\_vars}
\end{aligned}$$

$$\begin{aligned}
\text{JoinOldOlder}(r) & \triangleq \\
& \wedge r \in \text{old\_leaders} \\
& \wedge \text{zk\_status} = \text{"Normal"}
\end{aligned}$$

$$\begin{aligned}
& \wedge db\_epoch[r] = zk\_epoch \\
& \wedge Len(db[r]) \leq zk\_catchup\_index \\
& \wedge old\_leaders' = old\_leaders \setminus \{r\} \\
& \wedge zk\_epoch' = zk\_epoch + 1 \\
& \wedge \text{UNCHANGED } db\_vars \\
& \wedge \text{UNCHANGED } \langle zk\_leader\_epoch, zk\_leader, zk\_status, zk\_catchup\_index \rangle \\
& \wedge \text{UNCHANGED } client\_vars \\
& \wedge \text{UNCHANGED } healer\_vars
\end{aligned}$$

$$\begin{aligned}
TruncateOldLeader(r) & \triangleq \\
& \wedge r \in old\_leaders \\
& \wedge zk\_status = \text{"Normal"} \\
& \wedge db\_status[r] = \text{"Frozen"} \\
& \wedge db[r] \neq \langle \rangle \\
& \wedge db' = [db \text{ EXCEPT } ![r] = \langle \rangle] \\
& \wedge \text{UNCHANGED } \langle db\_epoch, db\_leader, db\_replicated, db\_status \rangle \\
& \wedge \text{UNCHANGED } zk\_vars \\
& \wedge \text{UNCHANGED } client\_vars \\
& \wedge \text{UNCHANGED } healer\_vars
\end{aligned}$$

$$checked\_max\_epoch \triangleq 14$$

$$\begin{aligned}
TerminateCond & \triangleq \\
& \wedge next\_req = 60 + max\_next\_req \\
& \wedge zk\_status = \text{"Normal"} \\
& \wedge zk\_leader\_epoch = max\_change\_leader \\
& \wedge zk\_epoch \geq checked\_max\_epoch - 1 \\
& \wedge \forall c \in Client : pending[c] = nil \wedge pending\_db[c] = nil \\
& \wedge \forall c \in Client : client\_leader\_epoch[c] = zk\_leader\_epoch \\
& \wedge \forall r \in Replica : db\_epoch[r] = zk\_epoch
\end{aligned}$$

$$\begin{aligned}
Terminated & \triangleq \\
& \wedge TerminateCond \\
& \wedge \text{UNCHANGED } vars
\end{aligned}$$

$$\begin{aligned}
Next & \triangleq \\
& \vee \exists c \in Client : \\
& \quad \vee StartRequest(c) \\
& \quad \vee DBResponse(c) \\
& \quad \vee ClientUpdateLeader(c) \\
& \vee \exists r \in Replica : \\
& \quad \vee Replicate(r) \\
& \quad \vee DBUpdateLeader(r) \\
& \quad \vee JoinOldOlder(r)
\end{aligned}$$

$$\begin{aligned}
& \vee \textit{TruncateOldLeader}(r) \\
& \vee \textit{ReadyToChangeZKLeader} \\
& \vee \textit{HealerUpdateState} \\
& \vee \exists r \in \textit{Replica} : \textit{HealerGetDBLog}(r) \\
& \vee \textit{HealerUpdateLeader} \\
& \vee \textit{HealerUpdateToNormal} \\
& \vee \textit{Terminated} \\
\\
\textit{Spec} & \triangleq \textit{Init} \wedge \Box[\textit{Next}]_{\textit{vars}} \\
\textit{FairSpec} & \triangleq \textit{Spec} \wedge \text{WF}_{\textit{vars}}(\textit{Next}) \\
\\
\textit{Consistent} & \triangleq \\
& \wedge \forall c \in \textit{Client} : \\
& \quad \wedge \textit{Len}(\textit{client\_success}[c]) \leq \textit{Len}(\textit{db}[\textit{zk\_leader}]) \\
& \quad \wedge \forall x \in \textit{Range}(\textit{client\_success}[c]) : x \in \textit{Range}(\textit{db}[\textit{zk\_leader}]) \\
\\
\textit{Perms} & \triangleq \textit{Permutations}(\textit{Replica}) \\
\\
\textit{Inv} & \triangleq \\
& \wedge \textit{zk\_epoch} < \textit{checked\_max\_epoch} \\
& \wedge \textit{zk\_leader\_epoch} \leq \textit{max\_change\_leader} \\
& \wedge (\textit{zk\_leader\_epoch} \geq 2) \Rightarrow (\forall c \in \textit{Client} : \textit{Len}(\textit{client\_success}[c]) < 4) \\
\\
\textit{Finish} & \triangleq \Diamond \textit{TerminateCond} \\
\\
\textit{CanRequestCond} & \triangleq \\
& \wedge \textit{zk\_leader\_epoch} \geq \textit{max\_change\_leader} \\
& \wedge \textit{next\_req} < 60 + \textit{max\_next\_req} \\
\\
\textit{CanStartReq}(c) & \triangleq \\
& \wedge \textit{pending}[c] = \textit{nil} \\
& \wedge \textit{db\_status}[\textit{client\_leader}[c]] = \text{"Writable"} \\
\\
\textit{CanRequest} & \triangleq \\
& \textit{CanRequestCond} \Rightarrow \exists c \in \textit{Client} : \textit{CanStartReq}(c) \\
\\
\textit{StillCanRequest} & \triangleq \Box \Diamond \textit{CanRequest} \\
\\
\textit{WeakFairnessOfStartReq} & \triangleq \text{WF}_{\textit{vars}}(\exists c \in \textit{Client} : \textit{StartRequest}(c))
\end{aligned}$$


---