

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 \langle
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*
 \rangle

zk_vars \triangleq \langle *zk_leader*, *zk_epoch*, *zk_leader_epoch*, *zk_status*,
old_leaders, *zk_catchup_index* \rangle

db_vars \triangleq \langle *db*, *db_leader*, *db_replicated*, *db_epoch*, *db_status* \rangle

client_vars \triangleq \langle
next_req, *client_leader*, *client_success*,
pending, *pending_db*, *client_leader_epoch* \rangle

healer_vars \triangleq \langle *healer_status*, *healer_epoch*, *healer_replicas* \rangle

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

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

replication_factor \triangleq 2

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

$TypeOK \triangleq$

$\wedge zk_leader \in Replica$
 $\wedge zk_epoch \in 1 \dots 30$
 $\wedge zk_leader_epoch \in 1 \dots \text{max_change_leader}$
 $\wedge zk_status \in \{\text{"Normal"}, \text{"ChangingLeader"}, \text{"WaitReplicaLog"}\}$
 $\wedge old_leaders \subseteq Replica$
 $\wedge zk_catchup_index \in (\text{LogOffset} \cup \{nil\})$

 $\wedge db \in [Replica \rightarrow Seq(ReqSet)]$
 $\wedge db_leader \in [Replica \rightarrow Replica]$
 $\wedge db_replicated \in [Replica \rightarrow [Replica \rightarrow \text{LogOffset}]]$
 $\wedge db_epoch \in [Replica \rightarrow Epoch]$
 $\wedge db_status \in [Replica \rightarrow \{\text{"Writable"}, \text{"Replica"}, \text{"Frozen"}\}]$

 $\wedge next_req \in ReqSet$
 $\wedge client_leader \in [Client \rightarrow Replica]$
 $\wedge client_success \in [Client \rightarrow Seq(ReqSet)]$
 $\wedge pending \in [Client \rightarrow NullReqSet]$
 $\wedge pending_db \in [Client \rightarrow NullReplica]$
 $\wedge client_leader_epoch \in [Client \rightarrow Epoch]$

 $\wedge healer_status \in \{\text{"Init"}, \text{"UpdatingLeader"}, \text{"WaitReplica"}\}$
 $\wedge healer_epoch \in Epoch$
 $\wedge healer_replicas \in [Replica \rightarrow \text{LogOffset} \cup \{nil\}]$

$Init \triangleq$

$\wedge zk_leader \in Replica$
 $\wedge zk_epoch = 1$
 $\wedge zk_leader_epoch = 1$
 $\wedge zk_status = \text{"Normal"}$
 $\wedge old_leaders = \{\}$
 $\wedge zk_catchup_index = nil$

 $\wedge db = [r \in Replica \mapsto \langle \rangle]$
 $\wedge db_leader = [r \in Replica \mapsto zk_leader]$
 $\wedge db_replicated = [r \in Replica \mapsto [r1 \in Replica \mapsto 0]]$
 $\wedge db_epoch = [r \in Replica \mapsto zk_epoch]$
 $\wedge db_status = [r \in Replica \mapsto \text{IF } zk_leader = r \text{ THEN "Writable" ELSE "Replica"}]$

 $\wedge next_req = 60$
 $\wedge client_leader = [c \in Client \mapsto zk_leader]$
 $\wedge client_success = [c \in Client \mapsto \langle \rangle]$
 $\wedge pending = [c \in Client \mapsto nil]$
 $\wedge pending_db = [c \in Client \mapsto nil]$

$$\wedge \text{client_leader_epoch} = [c \in \text{Client} \mapsto \text{zk_leader_epoch}]$$

$$\wedge \text{healer_status} = \text{"Init"}$$

$$\wedge \text{healer_epoch} = 1$$

$$\wedge \text{healer_replicas} = [r \in \text{Replica} \mapsto \text{nil}]$$

$$\text{StartRequest}(c) \triangleq$$

$$\wedge \text{pending}[c] = \text{nil}$$

$$\wedge \text{next_req} < 60 + \text{max_next_req}$$

$$\wedge \text{db_status}[\text{client_leader}[c]] = \text{"Writable"}$$

$$\wedge \text{next_req}' = \text{next_req} + 1$$

$$\wedge \text{pending}' = [\text{pending} \text{ EXCEPT } ![c] = \text{next_req}']$$

$$\wedge \text{pending_db}' = [\text{pending_db} \text{ EXCEPT } ![c] = \text{client_leader}[c]]$$

$$\wedge \text{LET } \text{leader} \triangleq \text{client_leader}[c] \text{ IN}$$

$$\wedge \text{db}' = [\text{db} \text{ EXCEPT } ![\text{leader}] = \text{Append}(@, \text{next_req}')]$$

$$\wedge \text{db_replicated}' = [$$

$$\text{db_replicated} \text{ EXCEPT } ![\text{leader}][\text{leader}] = \text{Len}(\text{db}'[\text{leader}]))$$

$$\wedge \text{UNCHANGED } \langle \text{client_leader}, \text{client_success}, \text{client_leader_epoch} \rangle$$

$$\wedge \text{UNCHANGED } \langle \text{db_leader}, \text{db_epoch}, \text{db_status} \rangle$$

$$\wedge \text{UNCHANGED } \text{zk_vars}$$

$$\wedge \text{UNCHANGED } \text{healer_vars}$$

$$\text{Replicate}(r) \triangleq$$

$$\wedge r \neq \text{db_leader}[r]$$

$$\wedge \text{db_status}[r] = \text{"Replica"}$$

$$\wedge \text{LET } \text{leader_data} \triangleq \text{db}[\text{db_leader}[r]]$$

$$\text{new_len} \triangleq \text{Len}(\text{db}[r]) + 1$$

$$\text{leader} \triangleq \text{db_leader}[r]$$

IN

$$\wedge \text{Len}(\text{db}[r]) < \text{Len}(\text{leader_data})$$

$$\wedge \text{db}' = [\text{db} \text{ EXCEPT } ![r] = \text{Append}(@, \text{leader_data}[\text{new_len}])]$$

$$\wedge \text{db_replicated}' = [\text{db_replicated} \text{ EXCEPT } ![\text{leader}][r] = \text{new_len}]$$

$$\wedge \text{UNCHANGED } \langle \text{db_leader}, \text{db_epoch}, \text{db_status} \rangle$$

$$\wedge \text{UNCHANGED } \text{client_vars}$$

$$\wedge \text{UNCHANGED } \text{zk_vars}$$

$$\wedge \text{UNCHANGED } \text{healer_vars}$$

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

$$\text{initReplicated}(r) \triangleq$$

$$\wedge \text{db_replicated}' = [$$

$$\text{db_replicated} \text{ EXCEPT } ![r] = [$$

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

]]

$$\begin{aligned}
& dbStatusFromZK(r) \triangleq \\
& \text{IF } zk_status \in \{ \text{"Normal"}, \text{"WaitReplicaLog"} \} \wedge \neg(r \in old_leaders) \\
& \quad \text{THEN IF } zk_leader = r \\
& \quad \quad \text{THEN "Writable"} \\
& \quad \quad \text{ELSE "Replica"} \\
& \quad \text{ELSE "Frozen"}
\end{aligned}$$

$$\begin{aligned}
& DBUpdateLeader(r) \triangleq \\
& \quad \wedge db_epoch[r] < zk_epoch \\
& \quad \wedge db_epoch' = [db_epoch \text{ EXCEPT } ![r] = zk_epoch] \\
& \quad \wedge db_leader' = [db_leader \text{ EXCEPT } ![r] = zk_leader] \\
& \quad \wedge db_status' = [db_status \text{ EXCEPT } ![r] = dbStatusFromZK(r)] \\
& \quad \wedge \text{IF } db_leader[r] \neq zk_leader \\
& \quad \quad \text{THEN } initReplicated(r) \\
& \quad \quad \text{ELSE UNCHANGED } db_replicated \\
& \quad \wedge \text{UNCHANGED } \langle db \rangle \\
& \quad \wedge \text{UNCHANGED } zk_vars \\
& \quad \wedge \text{UNCHANGED } client_vars \\
& \quad \wedge \text{UNCHANGED } healer_vars
\end{aligned}$$

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

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

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

$$\begin{aligned}
& DBResponse(c) \triangleq \\
& \quad \wedge pending[c] \neq nil \\
& \quad \wedge \text{LET } leader \triangleq pending_db[c] \text{ IN} \\
& \quad \quad \wedge \exists index \in \text{DOMAIN } db[leader], Q \in Quorum : \\
& \quad \quad \quad \wedge pending[c] = db[leader][index] \\
& \quad \quad \quad \wedge index \leq minReplicate(leader, Q) \\
& \quad \wedge pending' = [pending \text{ EXCEPT } ![c] = nil] \\
& \quad \wedge pending_db' = [pending_db \text{ EXCEPT } ![c] = nil] \\
& \quad \wedge client_success' = [client_success \text{ EXCEPT } ![c] = Append(@, pending[c])] \\
& \quad \wedge \text{UNCHANGED } db_vars \\
& \quad \wedge \text{UNCHANGED } next_req \\
& \quad \wedge \text{UNCHANGED } \langle client_leader, client_leader_epoch \rangle \\
& \quad \wedge \text{UNCHANGED } zk_vars \\
& \quad \wedge \text{UNCHANGED } healer_vars
\end{aligned}$$

$$\begin{aligned}
& ClientUpdateLeader(c) \triangleq \\
& \quad \wedge client_leader_epoch[c] < zk_leader_epoch \\
& \quad \wedge client_leader_epoch' = [client_leader_epoch \text{ EXCEPT } ![c] = zk_leader_epoch]
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{client_leader}' = [\text{client_leader} \text{ EXCEPT } ![c] = \text{zk_leader}] \\
& \wedge \text{pending}' = [\text{pending} \text{ EXCEPT } ![c] = \text{nil}] \\
& \wedge \text{pending_db}' = [\text{pending_db} \text{ EXCEPT } ![c] = \text{nil}] \\
& \wedge \text{UNCHANGED } \text{zk_vars} \\
& \wedge \text{UNCHANGED } \text{healer_vars} \\
& \wedge \text{UNCHANGED } \text{db_vars} \\
& \wedge \text{UNCHANGED } \langle \text{next_req}, \text{client_success} \rangle
\end{aligned}$$

$$\begin{aligned}
\text{ReadyToChangeZKLeader} & \triangleq \\
& \wedge \text{zk_leader_epoch} < \text{max_change_leader} \\
& \wedge \text{Cardinality}(\text{Replica} \setminus \text{old_leaders}) > \text{replication_factor} \\
& \wedge \text{zk_status} = \text{"Normal"} \\
& \wedge \text{zk_status}' = \text{"ChangingLeader"} \\
& \wedge \text{zk_epoch}' = \text{zk_epoch} + 1 \\
& \wedge \text{old_leaders}' = \text{old_leaders} \cup \{\text{zk_leader}\} \\
& \wedge \text{UNCHANGED } \langle \text{zk_leader}, \text{zk_leader_epoch}, \text{zk_catchup_index} \rangle \\
& \wedge \text{UNCHANGED } \text{client_vars} \\
& \wedge \text{UNCHANGED } \text{db_vars} \\
& \wedge \text{UNCHANGED } \text{healer_vars}
\end{aligned}$$

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

$$\begin{aligned}
\text{HealerUpdateState} & \triangleq \\
& \wedge \text{healer_epoch} < \text{zk_epoch} \\
& \wedge \text{healer_epoch}' = \text{zk_epoch} \\
& \wedge \text{healer_replicas}' = [r \in \text{Replica} \mapsto \text{nil}] \\
& \wedge \text{healer_status}' = \text{zkStatusToHealerStatus} \\
& \wedge \text{UNCHANGED } \text{zk_vars} \\
& \wedge \text{UNCHANGED } \text{client_vars} \\
& \wedge \text{UNCHANGED } \text{db_vars}
\end{aligned}$$

$$\begin{aligned}
\text{HealerGetDBLog}(r) & \triangleq \\
& \wedge \text{healer_status} = \text{"UpdatingLeader"} \\
& \wedge \text{healer_replicas}[r] = \text{nil} \\
& \wedge \neg(r \in \text{old_leaders}) \\
& \wedge \text{db_epoch}[r] = \text{healer_epoch}
\end{aligned}$$

$\wedge \text{healer_replicas}' = [\text{healer_replicas} \text{ EXCEPT } ![r] = \text{Len}(\text{db}[r])]$
 $\wedge \text{UNCHANGED } \langle \text{healer_status}, \text{healer_epoch} \rangle$
 $\wedge \text{UNCHANGED } \text{db_vars}$
 $\wedge \text{UNCHANGED } \text{client_vars}$
 $\wedge \text{UNCHANGED } \text{zk_vars}$

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

$\text{HealerUpdateLeader} \triangleq$
 $\wedge \text{healer_status} = \text{"UpdatingLeader"}$
 $\wedge \text{healer_epoch} = \text{zk_epoch}$
 $\wedge \text{Cardinality}(\text{collectedDB}) \geq \text{replication_factor}$
 $\wedge \exists r \in \text{collectedDB} :$
 $\quad \wedge \forall r1 \in \text{collectedDB} : \text{healer_replicas}[r] \geq \text{healer_replicas}[r1]$
 $\quad \wedge \text{zk_leader}' = r$
 $\quad \wedge \text{zk_catchup_index}' = \text{healer_replicas}[r]$
 $\wedge \text{zk_status}' = \text{"WaitReplicaLog"}$
 $\wedge \text{zk_epoch}' = \text{zk_epoch} + 1$
 $\wedge \text{zk_leader_epoch}' = \text{zk_leader_epoch} + 1$
 $\wedge \text{UNCHANGED } \text{old_leaders}$
 $\wedge \text{UNCHANGED } \text{healer_vars}$
 $\wedge \text{UNCHANGED } \text{db_vars}$
 $\wedge \text{UNCHANGED } \text{client_vars}$

$\text{HealerUpdateToNormal} \triangleq$
 $\wedge \text{healer_status} = \text{"WaitReplica"}$
 $\wedge \text{zk_status} = \text{"WaitReplicaLog"}$
 $\wedge \exists Q \in \text{Quorum} :$
 $\quad \wedge \neg(\text{old_leaders} \subseteq Q)$
 $\quad \wedge \forall r \in Q : \text{Len}(\text{db}[r]) \geq \text{zk_catchup_index}$
 $\wedge \text{zk_status}' = \text{"Normal"}$
 $\wedge \text{zk_epoch}' = \text{zk_epoch} + 1$
 $\wedge \text{UNCHANGED } \langle \text{old_leaders}, \text{zk_leader_epoch}, \text{zk_leader}, \text{zk_catchup_index} \rangle$
 $\wedge \text{UNCHANGED } \text{healer_vars}$
 $\wedge \text{UNCHANGED } \text{db_vars}$
 $\wedge \text{UNCHANGED } \text{client_vars}$

$\text{RecoverOldLeader}(r) \triangleq$
 $\wedge r \in \text{old_leaders}$
 $\wedge \text{zk_status} = \text{"Normal"}$
 $\wedge \text{db}' = [\text{db} \text{ EXCEPT } ![r] = \langle \rangle]$
 $\wedge \text{db_status}' = [\text{db_status} \text{ EXCEPT } ![r] = \text{"Replica"}]$
 $\wedge \text{db_epoch}' = [\text{db_epoch} \text{ EXCEPT } ![r] = \text{zk_epoch}]$
 $\wedge \text{db_leader}' = [\text{db_leader} \text{ EXCEPT } ![r] = \text{zk_leader}]$

$$\begin{aligned}
& \wedge db_replicated' = [db_replicated \text{ EXCEPT } ![r] = new_repl] \\
& \wedge old_leaders' = old_leaders \setminus \{r\} \\
& \wedge \text{UNCHANGED } \langle zk_epoch, zk_leader, zk_leader_epoch, zk_status, zk_catchup_index \rangle \\
& \wedge \text{UNCHANGED } client_vars \\
& \wedge \text{UNCHANGED } healer_vars
\end{aligned}$$

$$\begin{aligned}
TerminateCond & \triangleq \\
& \wedge next_req = 60 + max_next_req \\
& \wedge zk_status = \text{"Normal"} \\
& \wedge zk_leader_epoch = max_change_leader \\
& \wedge zk_epoch \geq 10 \\
& \wedge \forall c \in Client : pending[c] = nil \wedge pending_db[c] = nil \\
& \wedge \forall c \in Client : client_leader_epoch[c] = zk_leader_epoch \\
& \wedge \forall r \in Replica : db_epoch[r] = zk_epoch
\end{aligned}$$

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

$$\begin{aligned}
Next & \triangleq \\
& \vee \exists c \in Client : \\
& \quad \vee StartRequest(c) \\
& \quad \vee DBResponse(c) \\
& \quad \vee ClientUpdateLeader(c) \\
& \vee \exists r \in Replica : \\
& \quad \vee Replicate(r) \\
& \quad \vee DBUpdateLeader(r) \\
& \quad \vee RecoverOldLeader(r) \\
& \vee ReadyToChangeZKLeader \\
& \vee HealerUpdateState \\
& \vee \exists r \in Replica : HealerGetDBLog(r) \\
& \vee HealerUpdateLeader \\
& \vee HealerUpdateToNormal \\
& \vee Terminated
\end{aligned}$$

$$Spec \triangleq Init \wedge \Box[Next]_{vars}$$

$$FairSpec \triangleq Spec \wedge WF_{vars}(Next)$$

$$\begin{aligned}
Consistent & \triangleq \\
& \wedge \forall c \in Client : \\
& \quad \wedge Len(client_success[c]) \leq Len(db[zk_leader])
\end{aligned}$$

$$\wedge \forall x \in \text{Range}(\text{client_success}[c]) : x \in \text{Range}(\text{db}[\text{zk_leader}])$$

$$\text{Perms} \triangleq \text{Permutations}(\text{Replica})$$

$$\begin{aligned} \text{Inv} &\triangleq \\ &\wedge \text{zk_epoch} < 11 \\ &\wedge \text{zk_leader_epoch} \leq \text{max_change_leader} \\ &\wedge (\text{zk_leader_epoch} \geq 2) \Rightarrow (\forall c \in \text{Client}: \text{Len}(\text{client_success}[c]) < 4) \end{aligned}$$

$$\text{Finish} \triangleq \Diamond \text{TerminateCond}$$
