

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

CONSTANTS *COMMANDS*, set of records like [key: STRING, value: STRING]
 REPLICAS set of replicas

ASSUME *IsFiniteSet*(*REPLICAS*)

VARIABLES *msgPool*, messages in transit
 leader, current leader
 epoch, current epoch
 specPools, *specPool* of each replica
 unsyncedCmds, unsynced commands (backend)
 syncedCmds, synced commands (backend)
 requestedCmds, client requested commands
 committedCmds commands that client believes are stable in the *RSM*

vars \triangleq \langle *msgPool*,
 leader,
 epoch,
 specPools,
 syncedCmds,
 unsyncedCmds,
 requestedCmds,
 committedCmds \rangle

The initial state of the system.

Init \triangleq
 \wedge *msgPool* = {}
 \wedge *leader* \in *REPLICAS*
 \wedge *epoch* = 1
 \wedge *specPools* = [*r* \in *REPLICAS* \mapsto {}]
 \wedge *unsyncedCmds* = {}
 \wedge *syncedCmds* = {}
 \wedge *requestedCmds* = {}
 \wedge *committedCmds* = {}

Helper function for converting a set to a sequence.

SetToSeqs(*set*) \triangleq
 LET *len* \triangleq 1 .. *Cardinality*(*set*)
 seqs \triangleq {*f* \in [*len* \rightarrow *set*] : $\forall i, j \in \text{len} : i \neq j \Rightarrow f[i] \neq f[j]$ }
 IN *seqs*

Helper function for checking if a target exists in a sequence.

ExistInSeq(*seq*, *target*) \triangleq
 $\exists i \in 1 \dots \text{Len}(\text{seq}) : \text{seq}[i] = \text{target}$

SuperQuorums: In $N = 2 * f + 1$ replicas, a *SuperQuorum* is a set of replicas that contains at least $f + (f + 1) / 2 + 1$ replicas (including the leader).

The client can consider a command as committed if and only if it receives positive responses from a *SuperQuorum*.

This defines the set consisting of all *SuperQuorums*.

SuperQuorums \triangleq

$$\begin{aligned} &\{q \in \text{SUBSET } REPLICAS : \\ &\quad \wedge \text{Cardinality}(q) \geq (\text{Cardinality}(REPLICAS) * 3) \div 4 + 1 \\ &\quad \wedge \text{leader} \in q\} \end{aligned}$$

LeastQuorums: In $N = 2 * f + 1$ replicas, a *LeastQuorum* is a set of replicas that contains at least $(f + 1) / 2 + 1$ replicas.

When a replica becomes a leader, it must recover the command if and only if the command is a *LeastQuorum* of replicas' *specPool*.

This defines the set consisting of all *LeastQuorums*.

LeastQuorums \triangleq

$$\begin{aligned} &\{q \in \text{SUBSET } REPLICAS : \\ &\quad \text{Cardinality}(q) \geq \text{Cardinality}(REPLICAS) \div 4 + 1\} \end{aligned}$$

RecoveryQuorums: In $N = 2 * f + 1$ replicas, a *RecoveryQuorum* is a set of replicas that contains at least $f + 1$ replicas.

A replica must gather a *RecoveryQuorum* of replicas' *specPool* to recover the commands that need to be recovered.

This defines the set consisting of all *RecoveryQuorums*.

RecoveryQuorums \triangleq

$$\begin{aligned} &\{q \in \text{SUBSET } REPLICAS : \\ &\quad \text{Cardinality}(q) \geq \text{Cardinality}(REPLICAS) \div 2 + 1\} \end{aligned}$$

Client sends a request to all replicas.

ClientSendRequest(cmd) \triangleq

$$\begin{aligned} &\wedge \text{requestedCmds}' = \text{requestedCmds} \cup \{cmd\} \\ &\wedge \text{msgPool}' = \text{msgPool} \cup \\ &\quad [type : \{\text{"request"}\}, \\ &\quad \text{cmd} : \{cmd\}, \\ &\quad \text{dst} : REPLICAS] \\ &\wedge \text{UNCHANGED } \langle \text{leader}, \\ &\quad \text{epoch}, \\ &\quad \text{specPools}, \\ &\quad \text{syncdCmds}, \\ &\quad \text{unsyncdCmds}, \\ &\quad \text{committedCmds} \rangle \end{aligned}$$

Replica receives a request from the client.

If there is no conflict command (command on the same key) in the *specPool*, the replica adds the command to its *specPool*.

The leader will always add the command to *unsyncedCmds*.

ReplicaReceiveRequest(*r*, *msg*) \triangleq
 IF $\neg \text{ExistInSeq}(\text{syncdCmds}, \text{msg.cmd})$ THEN
 LET *conflict* $\triangleq (\exists \text{cmd} \in \text{specPools}[r] : \text{cmd.key} = \text{msg.cmd.key})$ IN
 $\wedge \text{specPools}' = [\text{specPools} \text{ EXCEPT } ![r] =$
 IF *conflict* THEN @ ELSE @ $\cup \{\text{msg.cmd}\}$
 $\wedge \text{unsyncdCmds}' =$
 IF *r* = *leader* THEN *Append*(*unsyncdCmds*, *msg.cmd*)
 ELSE *unsyncdCmds*
 $\wedge \text{msgPool}' = (\text{msgPool} \setminus \{\text{msg}\}) \cup$
 $\{[type \mapsto \text{"response"},$
 $\text{cmd} \mapsto \text{msg.cmd},$
 $ok \mapsto \neg \text{conflict},$
 $src \mapsto r]\}$
 $\wedge \text{UNCHANGED } \langle \text{leader},$
 epoch,
 syncdCmds,
 requestedCmds,
 committedCmds \rangle
 ELSE If the command is already synced, the replica does nothing.
 $\wedge \text{msgPool}' = \text{msgPool} \setminus \{\text{msg}\}$
 $\wedge \text{UNCHANGED } \langle \text{leader},$
 epoch,
 specPools,
 syncdCmds,
 unsyncdCmds,
 requestedCmds,
 committedCmds \rangle

Client receives a response from a replica.

If the client got positive responses from a *SuperQuorum*, the client considers the command as committed.

If the client got negative responses from a *LeastQuorum*, the command can never accepted by a *SuperQuorum*. Thus the client stops waiting for it.

ClientReceiveResponse(*msg*) \triangleq
 LET *sameCmdResp* \triangleq
 $\{\text{resp} \in \text{msgPool} : \text{resp.type} = \text{"response"} \wedge \text{resp.cmd} = \text{msg.cmd}\}$
 IN
 $\vee \wedge \{m.\text{src} : m \in \{\text{resp} \in \text{sameCmdResp} : \text{resp.ok}\}\} \in \text{SuperQuorums}$
 $\wedge \text{committedCmds}' = \text{committedCmds} \cup \{\text{msg.cmd}\}$
 $\wedge \text{msgPool}' = \text{msgPool} \setminus \text{sameCmdResp}$
 $\wedge \text{UNCHANGED } \langle \text{leader},$

$$\begin{aligned}
& \text{epoch,} \\
& \text{specPools,} \\
& \text{syncdCmds,} \\
& \text{unsyncdCmds,} \\
& \text{requestedCmds} \rangle \\
\vee \wedge \{m.\text{src} : m \in \{\text{resp} \in \text{sameCmdResp} : \neg \text{resp.ok}\}\} \in \text{LeastQuorums} \\
& \wedge \text{msgPool}' = \text{msgPool} \setminus \text{sameCmdResp} \\
& \wedge \text{UNCHANGED} \langle \text{leader,} \\
& \text{epoch,} \\
& \text{specPools,} \\
& \text{syncdCmds,} \\
& \text{unsyncdCmds,} \\
& \text{requestedCmds,} \\
& \text{committedCmds} \rangle
\end{aligned}$$

Client Actions

$$\begin{aligned}
\text{ClientAction} & \triangleq \\
& \vee \exists \text{cmd} \in (\text{COMMANDS} \setminus \text{requestedCmds}) : \text{ClientSendRequest}(\text{cmd}) \\
& \vee \exists \text{msg} \in \text{msgPool} : \wedge \text{msg.type} = \text{"response"} \\
& \wedge \text{ClientReceiveResponse}(\text{msg})
\end{aligned}$$

Replica Actions

$$\begin{aligned}
\text{ReplicaAction} & \triangleq \\
& \exists \text{msg} \in \text{msgPool} : \\
& \wedge \text{msg.type} = \text{"request"} \\
& \wedge \text{ReplicaReceiveRequest}(\text{msg.dst}, \text{msg})
\end{aligned}$$

Syncing an *unsyncdCmd*

$$\begin{aligned}
\text{SyncAction} & \triangleq \\
& \wedge \text{unsyncdCmds} \neq \langle \rangle \\
& \wedge \text{specPools}' = [r \in \text{REPLICAS} \mapsto \text{specPools}[r] \setminus \{\text{Head}(\text{unsyncdCmds})\}] \\
& \wedge \text{syncdCmds}' = \text{Append}(\text{syncdCmds}, \text{Head}(\text{unsyncdCmds})) \\
& \wedge \text{unsyncdCmds}' = \text{Tail}(\text{unsyncdCmds}) \\
& \wedge \text{UNCHANGED} \langle \text{msgPool,} \\
& \text{leader,} \\
& \text{epoch,} \\
& \text{requestedCmds,} \\
& \text{committedCmds} \rangle
\end{aligned}$$

Leader Change Action

The new leader should gather at least a *RecoveryQuorum* of replicas *specPool* to recover the commands.

Commands occurring in the *specPool* of a *LeastQuorum* of replicas need to be recovered.

$$\begin{aligned}
\text{LeaderChangeAction} &\triangleq \\
&\exists \text{newLeader} \in \text{REPLICAS} : \\
&\quad \wedge \text{newLeader} \neq \text{leader} \\
&\quad \wedge \exists \text{recoveryQuorum} \in \text{RecoveryQuorums} : \\
&\quad \quad \text{LET } \text{specPoolCmds} \triangleq \text{UNION } \{(\{ \\
&\quad \quad \quad [\text{cmd} \mapsto \text{cmd}, r \mapsto r] : \text{cmd} \in \text{specPools}[r] \\
&\quad \quad \quad \}) : r \in \text{recoveryQuorum}\} \\
&\quad \quad \text{filteredSpecPoolCmds} \triangleq \{c1 \in \text{specPoolCmds} : \{ \\
&\quad \quad \quad c.r : c \in \{c2 \in \text{specPoolCmds} : c1.\text{cmd} = c2.\text{cmd}\} \\
&\quad \quad \quad \} \in \text{LeastQuorums}\} \\
&\quad \quad \text{newSpecPool} \triangleq \{c.\text{cmd} : c \in \text{filteredSpecPoolCmds}\} \\
&\quad \text{IN} \\
&\quad \wedge \text{specPools}' = [\text{specPools} \text{ EXCEPT } ![\text{newLeader}] = \text{newSpecPool}] \\
&\quad \wedge \text{LET } \text{newUnsyncedCmds} \triangleq \\
&\quad \quad \text{CHOOSE } s \in \text{SetToSeqs}(\text{newSpecPool}) : \text{TRUE} \\
&\quad \quad \text{IN } \text{unsyncedCmds}' = \text{unsyncedCmds} \circ \text{newUnsyncedCmds} \\
&\quad \wedge \text{leader}' = \text{newLeader} \\
&\quad \wedge \text{epoch}' = \text{epoch} + 1 \\
&\quad \wedge \text{UNCHANGED } \langle \text{msgPool}, \\
&\quad \quad \text{syncedCmds}, \\
&\quad \quad \text{requestedCmds}, \\
&\quad \quad \text{committedCmds} \rangle
\end{aligned}$$

$$\begin{aligned}
\text{Next} &\triangleq \\
&\vee \text{LeaderChangeAction} \\
&\vee \text{ClientAction} \\
&\vee \text{ReplicaAction} \\
&\vee \text{SyncAction}
\end{aligned}$$

$$\text{Spec} \triangleq \text{Init} \wedge \Box[\text{Next}]_{\text{vars}} \wedge \text{WF}_{\text{vars}}(\text{Next})$$

Type Check

$$\begin{aligned}
\text{TypeOK} &\triangleq \\
&\wedge \text{msgPool} \subseteq \\
&\quad [\text{type} : \{\text{"request"}\}, \\
&\quad \text{cmd} : \text{COMMANDS}, \\
&\quad \text{dst} : \text{REPLICAS}] \cup \\
&\quad [\text{type} : \{\text{"response"}\}, \\
&\quad \text{cmd} : \text{COMMANDS}, \\
&\quad \text{ok} : \{\text{TRUE}, \text{FALSE}\}, \\
&\quad \text{src} : \text{REPLICAS}] \\
&\wedge \text{leader} \in \text{REPLICAS} \\
&\wedge \text{epoch} \in \text{Nat} \\
&\wedge \forall r \in \text{REPLICAS} : \\
&\quad \text{LET } \text{specPool} \triangleq \text{specPools}[r] \text{IN}
\end{aligned}$$

$$\begin{aligned}
& \wedge specPool \subseteq COMMANDS \\
& \wedge \forall cmd1, cmd2 \in specPool : \\
& \quad cmd1.key \neq cmd2.key \\
& \wedge syncedCmds \in \{SetToSeqs(s) : s \in SUBSET\ COMMANDS\} \\
& \wedge unsyncedCmds \in \{SetToSeqs(s) : s \in SUBSET\ COMMANDS\} \\
& \wedge requestedCmds \subseteq COMMANDS \\
& \wedge committedCmds \subseteq COMMANDS
\end{aligned}$$

Stability Property

If the client considers a command has been committed, the command must eventually be synced.

$Stability \triangleq$

$$\begin{aligned}
& \forall cmd \in COMMANDS : \\
& \quad cmd \in committedCmds \leadsto cmd \in \{syncedCmds[i] : i \in DOMAIN\ syncedCmds\}
\end{aligned}$$

THEOREM $Spec \Rightarrow (\Box TypeOK) \wedge Stability$