

Adapted from *Ryan's spec*.

Replication protocol with a constant number of replicas. In the stable state, there is one primary replica, to which all requests should be sent. Reads are handled locally by the primary and do not involve the other replicas. To process a write, the primary writes the new value to its disk and then sends *RPC* write requests to the other replicas. The number of replicas that must report success in order for the primary to report success to the client is configurable.

EXTENDS *MCPrimSecRep*

VARIABLES

CONSTANTS *IPrim*, initial primary replica (an element of *Rep*)
MinOk, min number of successful replica writes for write to succeed
HdlReconciliation(-, -), * (*o*, *v*) – handles reconcilliation at *MC* level
BObject blue object type

NoVal \triangleq CHOOSE *v* : *v* \notin *Val*

Blue_ErrorTypes \triangleq {"error", error at secondary while writing
"badver"} primary has the bad version

what do I do w/ staled???

Blue_DataType \triangleq *Val* \cup {*NoVal*} only one address!!!

Blue_WriteType \triangleq *Val* \cup {*NoVal*} only one address supported

Blue_NoWrite \triangleq *NoVal*

Blue_InitDataVal \triangleq *NoVal*

Blue_InitPrim(*o*) \triangleq *IPrim*

Blue_InitSec(*o*) \triangleq *Rep* \ {*IPrim*} all the other reps are secondaries

Blue_PrepateWr(*r*, *o*, *w*) \triangleq *w*

Blue_TypeInvariant \triangleq

\wedge *TypeInvariant*

Blue_Init \triangleq

\wedge *MCInit*

MASTER functioning

master decrees that read or modify the global state. They are always sent from a replica in the system, which must be alive.

gsm \triangleq *master* * renamed

A primary replica reports that one of its secondaries is dead.

PrimReportsDeath(*o*, *reporter*, *reportee*) \triangleq primary reports that a secondary is dead
 \wedge *reporter* = *cache*[*reporter*][*o*].*prim* reporter thinks it's the primary
 \wedge *reporter* = *master.objects*[*o*].*prim* master also thinks reporter is the primary
 \wedge *reportee* \in *cache*[*reporter*][*o*].*sec* primary thinks the *reportee* is a secondary

$\wedge \text{reportee} \in \text{master.objects}[o].\text{sec}$ the master knows *reportee* is among the living secondaries
 $\wedge \text{master}' = [\text{master} \text{ EXCEPT } !.\text{health}[\text{reportee}] = \text{"dead"} ,$ *reportee is dead*
 $\quad !.\text{objects}[o].\text{sec} = @ \setminus \{\text{reportee}\},$ no longer a *sec*
 $\quad !.\text{objects}[o].\text{version} = @ + 1]$ get on to the next *version*!
 $\wedge \text{cache}' = [\text{cache} \text{ EXCEPT } ![\text{reporter}][o] = \text{master'.objects}[o]]$ refresh cache at end of protocol
 $\wedge \text{UNCHANGED } \langle \text{data}, \text{stat}, \text{channel}, \text{resps} \rangle$
 $\wedge \text{NoReply}$

Before a new primary can become a prim and serve reads/writes, it must first reconcile *w/* the other replicas. This is modeled simply as another write request, which is started by the new primary, *w/* the value of its local disk. This action will block the new prim. This write will suffer the same possible failures. If this write fails, then the replicas who've failed are removed. If this new prim is removed by another replica, then that replica will do the same thing.

$\text{Reconcile}(o, \text{newp}) \triangleq$ A new primary reconciles *w/* the other replicas
 $\wedge \text{PrimaryContinuesWrite}(\text{newp}, o, \text{data}[\text{newp}][o], \{\text{cache}'[\text{newp}][o].\text{prim}\} \cup \text{cache}'[\text{newp}][o].\text{sec}, \text{cache}'[o].\text{prim})$

A secondary replica reports that its primary replica is dead. This action will only be enabled if *reportee* is indeed the primary replica. Because I don't model read leases, I ensure that a primary can only be declared dead when it really is dead. Because I already ensure that no replica can be stale at end of write, the master agrees *w/* making any secondary the new primary.

$\text{PrimReportedDead}(o, \text{reporter}, \text{reportee}) \triangleq$
 $\quad \wedge \text{stat}[\text{reportee}].\text{phase} = \text{"dead"} \setminus *$ ASSUMPTION : ensure primary is not falsely declared dead
 $\wedge \text{master.health}[\text{reporter}] = \text{"alive"}$ master must think *sec* is alive
 $\wedge \text{reporter} \in \text{master.objects}[o].\text{sec}$ reporter should be a live secondary
 $\wedge \text{reportee} = \text{cache}[\text{reporter}][o].\text{prim}$ the *sec* must believe reporter is primary
 $\wedge \text{reportee} = \text{master.objects}[o].\text{prim}$ master thinks reported replica is primary
 $\wedge \text{master}' = [\text{master} \text{ EXCEPT } !.\text{health}[\text{reportee}] = \text{"dead"} ,$
 $\quad !.\text{objects}[o] = [\text{prim} \mapsto \text{reporter},$ make reporter the primary
 $\quad \quad \text{sec} \mapsto @.\text{sec} \setminus \{\text{reporter}\},$ reporter is no longer a *sec*
 $\quad \quad \text{version} \mapsto @.\text{version} + 1]]$ get on to the next version
 $\wedge \text{cache}' = [\text{cache} \text{ EXCEPT } ![\text{reporter}][o] = \text{master'.objects}[o]]$
 \quad refresh cache at end of protocol
 $\wedge \text{Reconcile}(o, \text{reporter})$ reconcile new primary with all the rest of the replicas
 $\wedge \text{UNCHANGED } \langle \text{stat}, \text{channel}, \text{resps} \rangle$

A replica detects that another replica is dead, which can happen for any reason (or no reason at all). This simulates a replica detecting another replica dead based on heartbeats. Aside from that, allowing death detection to occur at any time is a conservative way to model and check all possible paths of execution. This decree is split into two cases, which are handled above.

$\text{DeathDecree}(o, \text{reporter}, \text{reportee}) \triangleq$ reporter reports that *reportee* is dead
 $\quad \wedge \text{reporter} \neq \text{reportee}$
 $\quad \wedge \text{master.health}[\text{reportee}] = \text{"alive"}$ the *reportee* must have been believed alive
 $\quad \wedge \text{master.objects}[o].\text{version} = \text{cache}[\text{reporter}][o].\text{version}$ reporter must be up-to-date

$$\begin{aligned}
& \wedge \vee \text{PrimReportsDeath}(o, \text{reporter}, \text{reportee}) \\
& \vee \text{PrimReportedDead}(o, \text{reporter}, \text{reportee}) \\
\text{DeathDecree} & \triangleq \\
& \exists o \in \text{BObject} : \\
& \quad \exists \text{reporter} \in \text{Rep}, \text{reportee} \in \text{Rep} : \\
& \quad \vee \text{DeathDecree}(o, \text{reporter}, \text{reportee}) \\
\text{Blue_MasterActions} & \triangleq \text{DeathDecree}
\end{aligned}$$

REPLICAS

Store an update. If the advertised version is older then the local vers reply *w/ badversion* and do not update store. Otherwise, update store & reply *w/ ok*.

$$\begin{aligned}
\text{Blue_StoreUpdate}(r, o, w, \text{version}) & \triangleq \\
\text{data}' & = [\text{data} \text{ EXCEPT } ![r][o] = w]
\end{aligned}$$

$$\begin{aligned}
\text{SelectSubset}(\text{set}, \text{func}(_)) & \triangleq \\
\text{LET } F[s \in \text{SUBSET}(\text{set})] & \triangleq \\
& \text{IF } s = \{\} \text{ THEN } \{\} \\
& \text{ELSE IF } \exists r \in s : \text{func}(r) \text{ THEN} \\
& \quad \text{LET } r \triangleq \text{CHOOSE } r \in s : \text{func}(r) \text{ IN} \\
& \quad F[s \setminus \{r\}] \cup \{r\} \\
& \text{ELSE } \{\} \\
\text{IN } F[\text{set}]
\end{aligned}$$

The primary kills all stale replicas who've failed to perform write. This action is only enabled if the primary hasn't received any *badver*. The primary goes to the master and requests that the failed resp be killed. The prim can only remove secondaries, not itself. The action is not enabled when ONLY the primary has failed writing. When this happens, the primary will eventually get removed by some replica. Until this happens, though, the primary cannot respond to new writes/reads because it's "busy" and it will never finish the current write.

$$\begin{aligned}
\text{PrimaryKillsFailedRep}(\text{reporter}, o) & \triangleq \\
\text{LET } \text{WriteFinished} & \triangleq \\
& \forall r \in \text{Rep} : \\
& \quad \wedge \text{resps}[\text{reporter}][o][r] \notin \{\text{"waiting"}, \text{"badver"}\} \\
\text{ExistReplicasToKill} & \triangleq \\
& \wedge \exists r \in \text{Rep} : \text{exists a replica that should be killed} \\
& \quad \wedge r \neq \text{reporter} \text{ the reporter cannot kill itself} \\
& \quad \wedge \text{resps}[\text{reporter}][o][r] \notin \{\text{"ok"}, \text{"n/a"}\} \text{ } r \text{ has failed (timeout/error)} \\
\text{SelectOk}(r) & \triangleq \text{resps}[\text{reporter}][o][r] = \text{"ok"} \\
\text{IN} \\
& \wedge \text{reporter} = \text{cache}[\text{reporter}][o].\text{prim} \text{ reporter thinks it's the primary} \\
& \wedge \text{stat}[\text{reporter}].\text{lock}[o] = \text{"busy"} \text{ reporter is mid-write} \\
& \wedge \text{WriteFinished} \text{ have responses from all live replicas, and all agree on version number} \\
& \wedge \text{ExistReplicasToKill} \text{ there exist some replicas to kill} \\
& \text{Go to master and kill the replicas that have failed}
\end{aligned}$$

$\wedge master.objects[o].version = cache[reporter][o].version$ master must agree on version number
 $\wedge master.objects[o].prim = reporter$ master believes reporter is prim (redundant)
 $\wedge master' = [master \text{ EXCEPT}$
 $\quad !.health = [r \in Rep \mapsto \text{IF } resps[reporter][o][r] \neq \text{"ok"}$
 $\quad \quad \text{THEN "dead" ELSE } @[r]],$
 $\quad !.objects[o].version = @ + 1,$
 $\quad \quad !.objects[o].sec = SelectSubset(@, SelectOk)]$
 $\wedge resps' = [resps \text{ EXCEPT } ![reporter][o] = [r \in Rep \mapsto \text{IF } @[r] \neq \text{"ok"}$
 $\quad \quad \text{THEN "n/a" ELSE } @[r]]$ mark as killed
 $\wedge cache' = [cache \text{ EXCEPT } ![reporter][o] = master'.objects[o]]$
 $\wedge \text{UNCHANGED } \langle stat, data, channel \rangle$
 $PrimaryKillsFailedRep \triangleq$
 $\exists r \in Rep, o \in BObject : \neg PrimaryKillsFailedRep(r, o)$

A helper function used to count the number of replicas that have successfully completed their part of a write request. Given a function with a subset of replicas as its domain, it counts the number of replicas that the function takes to "ok".

$NumOk(func) \triangleq$
 $\text{LET } F[d \in \text{SUBSET } Rep] \triangleq$
 $\quad \text{IF } \exists r \in d : func[r] = \text{"ok"}$
 $\quad \quad \text{THEN } 1 + F[d \setminus \{\text{CHOOSE } r \in d : func[r] = \text{"ok"}\}]$
 $\quad \quad \text{ELSE } 0$
 IN
 $F[\text{DOMAIN } func]$

R Replies to the client after the write is finished. A write can finish only when:

- all replicas have responded (no "waiting" *resps*) (checked in *PrimSecRep*)
- no replica has responded *w/ BadVer*
- no "error" in *resps* remains (all of those that were "error" were "killed")

$Blue_ReplyFinishedWrite(r, o, w) \triangleq$
 $\text{LET } NoBadVersion \triangleq$ no replica replied *badver*
 $\quad \forall r1 \in Rep : resps[r][o][r1] \neq \text{"badver"}$
 $NoStaleReplica \triangleq$ no stale replica has remained
 $\quad \quad \quad \text{stale means that a replica has failed to commit the write}$
 $\quad \forall r1 \in Rep : resps[r][o][r1] \in \{\text{"ok"}, \text{"n/a"}\}$
 $Success \triangleq$ the write succeeds if it has been committed to at least *MinOk* reps
 $\quad NumOk(resps[r][o]) \geq MinOk$
 IN
 $\wedge NoBadVersion$ *r* cannot finish a write if it received a bad version
 $\wedge NoStaleReplica$ no stale replicas have remaining
 $\wedge \text{IF } Success \text{ THEN } success$
 $\quad Reply(r, \text{"wrFinished"}, [ack \mapsto \text{"ok"}, object \mapsto o, w \mapsto w])$
 $\text{ELSE } failure$
 $\quad Reply(r, \text{"wrFinished"}, [ack \mapsto \text{"error"}, object \mapsto o, w \mapsto w])$

Perform the read on the local store of replica r .

$Blue_StoreRead(r, o) \triangleq$
 $[ack \mapsto \text{"ok"}, object \mapsto o, val \mapsto data[r][o]]$

$Blue_ReplicaActions \triangleq$
 $ReplicaActions \vee PrimaryKillsFailedRep$

CLIENT actions.

Read-from-primary policy of Blue.

$Blue_CliRead(o) \triangleq$
 $\exists r \in Rep :$
 $\wedge master.health[r] = \text{"alive"}$ ASSUMPTION : r is not a stale primary.
 This is ensured in the real protocol via the *ReadLease*.
 But I don't model read leases here.
 $\wedge stat[r].phase = \text{"alive"}$
 $\wedge r = cache[r][o].prim$ Read-from-primary policy
 $\wedge stat[r].lock[o] = \text{"rdy"}$ no writes in-progress at r
 $\wedge _CliRead(r, o)$

$Blue_ClientActions \triangleq$
 $\exists o \in Object :$
 $\vee \exists w \in WriteType : CliWrite(o, w)$
 $\vee Blue_CliRead(o)$

$Blue_Next \triangleq$
 $\vee Blue_MasterActions$
 $\vee TimeActions$
 $\vee Blue_ReplicaActions$
 $\vee Blue_ClientActions$

$Blue_Impl \triangleq Blue_Init \wedge \Box[Blue_Next]_{mcvars}$

THEOREM $Blue_Impl \Rightarrow \Box Blue_TypeInvariant$
