
MODULE *IR_consensus*

This is a TLA+ specification of the *Inconsistent* Replication algorithm. (And a mechanically-checked proof of its correctness using *TLAPS*)

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

Constants

Constant parameters: Replicas: the set of all replicas (Replica *IDs*)

Clients: the set of all clients (Client *IDs*)

Quorums: the set of all quorums *SuperQuorums*: the set of all super quorums Results: the set of all possible result types *OperationBody*: the set of all possible operation bodies (with arguments, etc. - can be infinite)

S: shard *id* of the shard *Replicas* constitute

f: maximum number of failures allowed (half of *n*)

Constants used to bound variables, for model checking (*Nat* is bounded) *max_vc*: maximum number of View-Changes allowed for each replicas *max_req*: maximum number of *op* requests performed by clients

CONSTANTS *Replicas*, *Clients*, *Quorums*, *SuperQuorums*, *Results*, *OpBody*,
AppClientFail, *AppReplicaFail*,
SuccessfulInconsistentOp(-, -, -), *SuccessfulConsensusOp*(-, -, -, -),
Merge(-, -, -),
Sync(-),
ExecInconsistent(-),
ExecConsensus(-),
Decide(-),
f,
S, *Shards*, *S* = shard *id*
max_vc, *max_req*

ASSUME *IsFiniteSet*(*Replicas*)

ASSUME *QuorumAssumption* \triangleq
 \wedge *Quorums* \subseteq SUBSET *Replicas*
 \wedge *SuperQuorums* \subseteq SUBSET *Replicas*
 $\wedge \forall Q1, Q2 \in \textit{Quorums} : Q1 \cap Q2 \neq \{\}$
 $\wedge \forall Q \in \textit{Quorums}, R1, R2 \in \textit{SuperQuorums} :$
 $Q \cap R1 \cap R2 \neq \{\}$

ASSUME *FailuresAssumption* \triangleq
 $\forall Q \in \textit{Quorums} : \textit{Cardinality}(Q) > f$

The possible states of a replica and the two types of operations currently defined by *IR*.

ReplicaState \triangleq {"NORMAL", "FAILED", "RECOVERING", "VIEW-CHANGING"}
ClientState \triangleq {"NORMAL", "FAILED"}

$OpType \triangleq \{ \text{"Inconsistent"}, \text{"Consensus"} \}$
 $OpStatus \triangleq \{ \text{"TENTATIVE"}, \text{"FINALIZED"} \}$

Definition of operation space

$MessageId \triangleq [cid : Clients, msgid : Nat]$
 $Operations \triangleq [type : OpType, body : OpBody] \cup \{ \langle \rangle \}$

Message is defined to be the set of all possible messages

Assume unique message ids

Assume no more than f replica failures

We use shard to specify for what shard this message was
(we share the variables)

$Message \triangleq$
 $[type : \{ \text{"PROPOSE"} \},$
 $id : MessageId,$
 $op : Operations,$
 $v : Nat]$
 $\cup [type : \{ \text{"REPLY"} \},$ reply no result
 $id : MessageId,$
 $v : Nat,$
 $src : Replicas]$
 \cup
 $[type : \{ \text{"REPLY"} \},$ reply with result
 $id : MessageId,$
 $v : Nat,$
 $res : Results,$
 $src : Replicas]$
 $v = \text{view num.}$
 \cup
 $[type : \{ \text{"START-VIEW-CHANGE"} \},$
 $v : Nat,$
 $src : Replicas]$
 \cup
 $[type : \{ \text{"DO-VIEW-CHANGE"} \},$
 $r : \text{SUBSET } ([msgid : MessageId,$
 $op : Operations,$
 $res : Results,$
 $status : OpStatus]$
 $\cup [msgid : MessageId,$
 $op : Operations,$
 $status : OpStatus]),$
 $v : Nat,$
 $lv : Nat,$

$$\begin{aligned}
& \text{src} : \text{Replicas}, \\
& \text{dst} : \text{SUBSET } \text{Replicas}] \\
\cup & \\
& [\text{type} : \{ \text{"START-VIEW"} \}, \\
& \text{v} : \text{Nat}, \\
& \text{r} : \text{SUBSET } ([\text{msgid} : \text{MessageId}, \\
& \quad \text{op} : \text{Operations}, \\
& \quad \text{res} : \text{Results}, \\
& \quad \text{status} : \text{OpStatus}] \\
& \quad \cup [\text{msgid} : \text{MessageId}, \\
& \quad \text{op} : \text{Operations}, \\
& \quad \text{status} : \text{OpStatus}]), \\
& \text{src} : \text{Replicas}] \\
\cup & \\
& [\text{type} : \{ \text{"FINALIZE"} \}, \text{finalize with no result} \\
& \text{id} : \text{MessageId}, \\
& \text{op} : \text{Operations}, \\
& \text{res} : \text{Results}, \\
& \text{v} : \text{Nat}] \\
\cup & \\
& [\text{type} : \{ \text{"FINALIZE"} \}, \text{finalize with result} \\
& \text{id} : \text{MessageId}, \\
& \text{op} : \text{Operations}, \\
& \text{res} : \text{Results}, \\
& \text{v} : \text{Nat}] \\
\cup & \\
& [\text{type} : \{ \text{"CONFIRM"} \}, \\
& \text{v} : \text{Nat}, \\
& \text{id} : \text{MessageId}, \\
& \text{op} : \text{Operations}, \\
& \text{res} : \text{Results}, \\
& \text{src} : \text{Replicas}]
\end{aligned}$$

Variables and State Predicates

\ * Variables: 1. State at each replica:

- $rState$ = Denotes current replica state. Either:
 - *NORMAL* (processing operations)
 - *VIEW-CHANGING* (participating in recovery)
- $rRecord$ = Unordered set of operations and their results
- $rViewNumber$ = current view number

2. State of communication medium: $sentMsg$ = sent (but not yet received) messages

3. State at client: $cCurrentOperation$ = crt operation requested by the client

$cMmessageCounter$ = the message I must use for the next operation

* \

VARIABLES $rState, rRecord, rCrtView, rLastView, rViewReplies,$
 $rViewOnDisk,$
 $rNonce, sentMsg, cCrtOp,$
 $cCrtOpToFinalize, cMsgCounter, cCrtOpReplies, cCrtOpConfirms,$
 $cState, aSuccessful, arRecord, aVisibility, gViewChangesNo$

Defining these tuples makes it easier to express which variables remain unchanged

Replica variables.

$rVars \triangleq \langle rState, rRecord, rCrtView, rViewOnDisk, rLastView,$
 $rViewReplies, rNonce \rangle$

Client variables.

$cVars \triangleq \langle cCrtOp,$ current operation at a client
 $cCrtOpToFinalize,$
 $cCrtOpReplies,$ current operation replies
 $cCrtOpConfirms,$
 $cMsgCounter,$
 $cState \rangle$

Application variables.

$aVars \triangleq \langle aSuccessful, arRecord, aVisibility \rangle$ we use them to write invariants

Other variables.

$oVars \triangleq \langle sentMsg, gViewChangesNo \rangle$

All variables.

$vars \triangleq \langle rVars, cVars, aVars, oVars \rangle$

$TypeOK \triangleq$

$\wedge rState[S] \in [Replicas \rightarrow ReplicaState]$
 $\wedge rRecord[S] \in [Replicas \rightarrow \text{SUBSET} ([msgid : MessageId,$
 $op : Operations,$
 $res : Results,$
 $status : OpStatus]$
 $\cup [msgid : MessageId,$
 $op : Operations,$
 $status : OpStatus])]$
 $\wedge rCrtView[S] \in [Replicas \rightarrow Nat]$
 $\wedge rViewOnDisk[S] \in [Replicas \rightarrow Nat]$
 $\wedge rLastView[S] \in [Replicas \rightarrow Nat]$
 $\wedge rViewReplies[S] \in [Replicas \rightarrow \text{SUBSET} ([type : \{ "start-view-change" \},$
 $v : Nat,$
 $src : Replicas]$
 $\cup [type : \{ "do-view-change" \},$
 $v : Nat,$
 $lv : Nat,$
 $r : \text{SUBSET} ([msgid : MessageId,$

$$\begin{aligned}
&\wedge cCrtOp = [c \in Clients \mapsto \langle \rangle] \\
&\wedge cCrtOpToFinalize = [c \in Clients \mapsto \langle \rangle] \\
&\wedge cCrtOpReplies = [c \in Clients \mapsto \{\}] \\
&\wedge cCrtOpConfirms = [c \in Clients \mapsto \{\}] \\
&\wedge cMsgCounter = [c \in Clients \mapsto 0] \\
&\wedge cState = [c \in Clients \mapsto \text{"NORMAL"}] \\
&\wedge aSuccessful = \{\} \\
&\wedge aVisibility = [o \in MessageId \mapsto \{\}] \\
&\wedge arRecord = [r \in Replicas \mapsto \{\}] \\
&\wedge gViewChangesNo = 0
\end{aligned}$$

Actions

$$\begin{aligned}
Send(m) &\triangleq sentMsg' = [sentMsg \text{ EXCEPT } ![S] = @ \cup \{m\}] \\
AmLeader(r, v) &\triangleq r = (v \% Cardinality(Replicas)) + 1 \\
IsLeader(r, v) &\triangleq AmLeader(r, v) \\
NotIsLeader(r, v) &\triangleq r \neq (v \% Cardinality(Replicas)) + 1 \\
LeaderOf(v) &\triangleq \text{CHOOSE } x \in Replicas : IsLeader(x, v)
\end{aligned}$$

Client Actions

Note: CHOOSE does not introduce nondeterminism (the same value is chosen each time)

$$\begin{aligned}
&\backslash * \text{ Client sends a request} \\
ClientRequest(c, op) &\triangleq \\
&\wedge cCrtOp[S][c] = \langle \rangle \quad \text{the client is not waiting for a result} \\
&\quad \text{of another operation} \\
&\wedge cCrtOpToFinalize[S][c] = \langle \rangle \quad \text{the client is not waiting} \\
&\quad \text{to finalize operation} \\
&\wedge cMsgCounter' = [cMsgCounter \text{ EXCEPT } ![S][c] = @ + 1] \\
&\wedge cCrtOp' = [cCrtOp \text{ EXCEPT } ![S][c] = op] \\
&\wedge Send([type \mapsto \text{"PROPOSE"}, \\
&\quad id \mapsto [cid \mapsto c, msgid \mapsto cMsgCounter[S][c] + 1], \\
&\quad op \mapsto op, \\
&\quad v \mapsto 0]) \\
&\wedge \text{UNCHANGED } \langle rVars, aVars, cCrtOpReplies, cCrtOpToFinalize, \\
&\quad cCrtOpConfirms, cState, gViewChangesNo \rangle \\
&\wedge cMsgCounter[S][c] < max_req \quad \text{BOUND the number of requests a client can make} \\
&\quad \text{(useful for model checking)} \\
&\backslash * \text{ Client received a reply} \\
ClientReceiveReply(c) &\triangleq
\end{aligned}$$

$$\begin{aligned}
& \wedge \vee \wedge cCrtOp[S][c].type = \text{"Inconsistent"} \\
& \quad \wedge \text{__matchingViewNumbers}(Q, c) \\
& \quad \wedge aSuccessful' = aSuccessful \cup \\
& \quad \quad \{ [mid \mapsto [cid \mapsto c, \\
& \quad \quad \quad msgid \mapsto cMsgCounter[S][c], \\
& \quad \quad \quad op \mapsto cCrtOp[S][c]] \} \\
& \quad \wedge SuccessfulInconsistentOp(c, S, cCrtOp[S][c]) \\
& \quad \wedge Send([type \mapsto \text{"FINALIZE"}, \\
& \quad \quad id \mapsto [cid \mapsto c, msgid \mapsto cMsgCounter[S][c], \\
& \quad \quad \quad op \mapsto cCrtOp[S][c], \\
& \quad \quad \quad v \mapsto 0]) \\
& \quad \wedge \text{UNCHANGED } \langle cCrtOpToFinalize \rangle \\
& \vee \wedge cCrtOp[S][c].type = \text{"Consensus"} \\
& \quad \wedge \text{__matchingViewNumbers}(Q, c) \\
& \quad \wedge \text{LET } res \stackrel{\Delta}{=} \text{IF } \text{__matchingViewNumbersAndResults}(Q, c) \\
& \quad \quad \text{THEN} \\
& \quad \quad \quad \text{CHOOSE } result \in \\
& \quad \quad \quad \quad \{ res \in Results : \\
& \quad \quad \quad \quad \quad \exists reply \in cCrtOpReplies[S][c] : \\
& \quad \quad \quad \quad \quad \quad \wedge reply.src \in Q \\
& \quad \quad \quad \quad \quad \quad \wedge reply.res = res \} : \text{TRUE} \\
& \quad \quad \text{ELSE} \\
& \quad \quad \quad \text{Decide}(cCrtOpReplies[S][c]) \\
& \quad \text{IN} \\
& \quad \quad \wedge Send([type \mapsto \text{"FINALIZE"}, \\
& \quad \quad \quad id \mapsto [cid \mapsto c, msgid \mapsto cMsgCounter[S][c], \\
& \quad \quad \quad \quad op \mapsto cCrtOp[S][c], \\
& \quad \quad \quad \quad \quad res \mapsto res, \\
& \quad \quad \quad \quad \quad \quad v \mapsto 0]) \\
& \quad \quad \wedge cCrtOpToFinalize' = [cCrtOp \text{ EXCEPT } ![S][c] = cCrtOp[S][c]] \\
& \quad \quad \wedge \text{UNCHANGED } \langle aSuccessful \rangle \\
& \vee \exists SQ \in SuperQuorums : \\
& \quad \text{II. The IR Client got super quorum of responses} \\
& \quad \wedge \forall r \in SQ : \\
& \quad \quad \exists reply \in cCrtOpReplies[S][c] : reply.src = r \\
& \quad \wedge cCrtOp[S][c].type = \text{"Consensus"} \quad \text{only care if consensus } op \\
& \quad \wedge \text{__matchingViewNumbersAndResults}(SQ, c) \\
& \quad \wedge \text{LET } res \stackrel{\Delta}{=} \text{CHOOSE } result \in \\
& \quad \quad \{ res \in Results : \\
& \quad \quad \quad \exists reply \in cCrtOpReplies[S][c] : \\
& \quad \quad \quad \quad \wedge reply.src \in SQ \\
& \quad \quad \quad \quad \wedge reply.res = res \} : \text{TRUE} \\
& \quad \text{IN}
\end{aligned}$$

$$\begin{aligned}
& op \mapsto cCrtOpToFinalize[S][c], \\
& res \mapsto reply.res \} \\
& \wedge SuccessfulConsensusOp(c, S, cCrtOp[S][c], reply.res) \text{ respond to app} \\
& \wedge cCrtOpToFinalize' = [cCrtOpToFinalize \text{ EXCEPT } ![S][c] = \langle \rangle] \\
& \wedge cCrtOpConfirms' = [cCrtOpConfirms \text{ EXCEPT } ![S][c] = \{\}] \\
& \wedge \text{UNCHANGED } \langle rVars, cCrtOp, cCrtOpReplies, \\
& \quad cMsgCounter, cState, arRecord, aVisibility, oVars \rangle \\
& \backslash * \text{ Client fails and loses all data} \\
& ClientFail(c) \triangleq \\
& \wedge cState' = [cState \text{ EXCEPT } ![S][c] = \text{"FAILED"}] \\
& \wedge cMsgCounter' = [cMsgCounter \text{ EXCEPT } ![S][c] = 0] \\
& \wedge cCrtOp' = [cCrtOp \text{ EXCEPT } ![S][c] = \langle \rangle] \\
& \wedge cCrtOpReplies' = [cCrtOpReplies \text{ EXCEPT } ![S][c] = \{\}] \\
& \wedge AppClientFail \\
& \wedge \text{UNCHANGED } \langle rVars, aVars, oVars \rangle \\
& \backslash * \text{ Client recovers} \\
& \backslash * \text{ Not implemented yet} \\
& ClientRecover(c) \triangleq \text{FALSE}
\end{aligned}$$

Replica Actions

$$\begin{aligned}
& \backslash * \text{ Replica sends a reply} \\
& ReplicaReceiveRequest(r) \triangleq \\
& \exists msg \in sentMsg[S] : \\
& \wedge msg.type = \text{"PROPOSE"} \\
& \wedge \text{not already replied for this } op \\
& \neg(\exists rec \in rRecord[S][r] : rec.msgid = msg.id) \\
& \wedge \vee \wedge msg.op.type = \text{"Inconsistent"} \\
& \quad \wedge Send([type \mapsto \text{"REPLY"}, \\
& \quad \quad id \mapsto msg.id, \\
& \quad \quad v \mapsto rCrtView[S][r], \\
& \quad \quad src \mapsto r]) \\
& \wedge rRecord' = [rRecord \text{ EXCEPT } ![S][r] = \\
& \quad @ \cup \{[msgid \mapsto msg.id, \\
& \quad \quad op \mapsto msg.op, \\
& \quad \quad status \mapsto \text{"TENTATIVE"}]\}] \\
& \vee \wedge msg.op.type = \text{"Consensus"} \\
& \wedge \text{LET } res \triangleq ExecConsensus(msg.op) \\
& \text{IN} \\
& \quad \wedge Send([type \mapsto \text{"REPLY"}, \\
& \quad \quad id \mapsto msg.id, \\
& \quad \quad v \mapsto rCrtView[S][r], \\
& \quad \quad res \mapsto res,
\end{aligned}$$

$$\begin{aligned}
& \text{src} \mapsto r]) \\
& \wedge rRecord' = [rRecord \text{ EXCEPT } ![S][r] = \\
& \quad @ \cup \{[msgid \mapsto msg.id, \\
& \quad \quad op \mapsto msg.op, \\
& \quad \quad res \mapsto res, \\
& \quad \quad status \mapsto \text{"TENTATIVE"}]\}] \\
& \wedge \text{UNCHANGED } \langle rState, rCrtView, rLastView, rViewOnDisk, \\
& \quad rViewReplies, rNonce, cVars, aVars, gViewChangesNo \rangle
\end{aligned}$$

```

\ * Replica receives a message from an IR Client to finalize an op
\ * For inconsistent operations the replica just
\ * executes the operation.
ReplicaReceiveFinalize(r)  $\triangleq$ 
 $\exists msg \in sentMsg[S] :$ 
   $\wedge msg.type = \text{"FINALIZE"}$ 
   $\wedge msg.v \geq rCrtView[S][r]$ 
   $\wedge \text{LET } recs \triangleq \{rec \in rRecord[S][r] : \text{Must be only 1 record}$ 
     $\wedge rec.msgid = msg.id$ 
     $\wedge rec.op = msg.op\}$ 
  IN
     $\vee \wedge msg.op.type = \text{"Inconsistent"}$ 
     $\wedge \text{IF}$ 
       $\text{Replica knows of this op}$ 
       $recs \neq \{\}$ 
      THEN
        IF  $\forall rec \in recs : rec.status \neq \text{"FINALIZED"}$ 
          THEN  $ExecInconsistent(msg.op)$ 
          ELSE TRUE
        ELSE
           $\text{Replica didn't hear of this op}$ 
           $ExecInconsistent(msg.op)$ 
       $\wedge rRecord' = [rRecord \text{ EXCEPT } ![S][r] = (@ \setminus recs) \cup$ 
         $\{[msgid \mapsto msg.id,$ 
         $\quad op \mapsto msg.op,$ 
         $\quad status \mapsto \text{"FINALIZED"}]\}]$ 
       $\wedge \text{UNCHANGED } \langle sentMsg \rangle$ 
     $\vee \wedge msg.op.type = \text{"Consensus"}$ 
     $\wedge rRecord' = [rRecord \text{ EXCEPT } ![S][r] = (@ \setminus recs) \cup$ 
       $\{[msgid \mapsto msg.id,$ 
       $\quad op \mapsto msg.op,$ 
       $\quad res \mapsto msg.res,$ 
       $\quad status \mapsto \text{"FINALIZED"}]\}]$ 
     $\wedge Send([type \mapsto \text{"CONFIRM"},$ 
       $\quad v \mapsto rCrtView[S][r],$ 
       $\quad id \mapsto msg.id,$ 

```

$$\begin{aligned}
& op \mapsto msg.op, \\
& res \mapsto msg.res, \\
& src \mapsto r]) \\
& \wedge \text{UNCHANGED } \langle rState, rCrtView, rLastView, rViewReplies, rViewOnDisk, \\
& \quad rNonce, cVars, aVars, gViewChangesNo \rangle
\end{aligned}$$

\setminus * A recovering replica starts the view change procedure

$$\begin{aligned}
& ReplicaSendDoViewChange(r) \triangleq \\
& \wedge \vee \wedge rState[S][r] = \text{"NORMAL"} \vee rState[S][r] = \text{"VIEW-CHANGING"} \\
& \quad \wedge rCrtView' = [rCrtView \text{ EXCEPT } ![S][r] = @ + 1] \\
& \quad \wedge rViewOnDisk' = [rViewOnDisk \text{ EXCEPT } ![S][r] = rCrtView[S][r] + 1] \\
& \quad \wedge rState' = [rState \text{ EXCEPT } ![S][r] = \text{"VIEW-CHANGING"}] \\
& \quad \wedge Send([type \mapsto \text{"DO-VIEW-CHANGE"}, \\
& \quad \quad v \mapsto rCrtView[S][r] + 1, \\
& \quad \quad lv \mapsto rLastView[S][r], \\
& \quad \quad r \mapsto rRecord[S][r], \\
& \quad \quad src \mapsto r, \\
& \quad \quad dst \mapsto Replicas]) \\
& \vee \wedge rState[S][r] = \text{"FAILED"} \\
& \quad \wedge rState' = [rState \text{ EXCEPT } ![S][r] = \text{"RECOVERING"}] \\
& \quad \wedge rCrtView' = [rCrtView \text{ EXCEPT } ![S][r] = rViewOnDisk[S][r]] \\
& \quad \wedge Send([type \mapsto \text{"DO-VIEW-CHANGE"}, \\
& \quad \quad v \mapsto rViewOnDisk[S][r], \\
& \quad \quad lv \mapsto rLastView[S][r], \\
& \quad \quad r \mapsto rRecord[S][r], \\
& \quad \quad src \mapsto r, \\
& \quad \quad dst \mapsto Replicas \setminus \{x \in Replicas : IsLeader(x, rViewOnDisk[S][r])\}]) \\
& \quad \wedge \text{UNCHANGED } \langle rViewOnDisk \rangle \\
& \vee \wedge rState[S][r] = \text{"RECOVERING"} \\
& \quad \wedge rCrtView' = [rCrtView \text{ EXCEPT } ![S][r] = @ + 1] \\
& \quad \wedge Send([type \mapsto \text{"DO-VIEW-CHANGE"}, \\
& \quad \quad v \mapsto rCrtView[S][r] + 1, \\
& \quad \quad lv \mapsto rLastView[S][r], \\
& \quad \quad r \mapsto rRecord[S][r], \\
& \quad \quad src \mapsto r, \\
& \quad \quad dst \mapsto Replicas \setminus \{x \in Replicas : IsLeader(x, rCrtView[S][r] + 1)\}]) \\
& \quad \wedge \text{UNCHANGED } \langle rViewOnDisk, rState \rangle \\
& \wedge \text{UNCHANGED } \langle cVars, rLastView, rViewReplies, rRecord, \\
& \quad rNonce, aVars \rangle \\
& \wedge gViewChangesNo[S] < max_vc \quad \text{BOUND on number of view changes} \\
& \wedge gViewChangesNo' = [gViewChangesNo \text{ EXCEPT } ![S] = @ + 1]
\end{aligned}$$

\setminus * Replica received DO-VIEW-CHANGE message

$$\begin{aligned}
& ReplicaReceiveDoViewChange(r) \triangleq \\
& \wedge \exists msg \in sentMsg[S] :
\end{aligned}$$

$$\begin{aligned}
& \wedge msg.type = \text{"DO-VIEW-CHANGE"} \\
& \wedge r \in msg.dst \\
& \wedge \vee \wedge rState[S][r] = \text{"NORMAL"} \\
& \quad \wedge msg.v > rCrtView[S][r] \\
& \quad \vee \wedge rState[S][r] = \text{"VIEW-CHANGING"} \\
& \quad \wedge msg.v \geq rCrtView[S][r] \\
& \wedge rState' = [rState \text{ EXCEPT } ![S][r] = \text{"VIEW-CHANGING"}] \\
& \quad \text{keep only the one with the higher view } (v) \\
& \wedge \vee \wedge IsLeader(r, msg.v) \\
& \quad \wedge \text{LET} \\
& \quad \quad existingRecord \triangleq \{x \in rViewReplies[S][r] : \\
& \quad \quad \quad \wedge x.type = \text{"do-view-change"} \\
& \quad \quad \quad \wedge x.src = msg.src\} \text{ should only be one item in set} \\
& \quad \text{IN} \\
& \quad \text{IF } \forall x \in existingRecord : x.v < msg.v \\
& \quad \quad \text{THEN } rViewReplies' = [rViewReplies \text{ EXCEPT } ![S][r] = \\
& \quad \quad \quad (@ \setminus existingRecord) \cup \\
& \quad \quad \quad \{[type \mapsto \text{"do-view-change"}, \\
& \quad \quad \quad \quad v \mapsto msg.v, \\
& \quad \quad \quad \quad lv \mapsto msg.lv, \\
& \quad \quad \quad \quad r \mapsto msg.r, \\
& \quad \quad \quad \quad src \mapsto msg.src]\}] \\
& \quad \quad \text{ELSE FALSE} \\
& \quad \vee \text{UNCHANGED } \langle rViewReplies \rangle \\
& \wedge rCrtView' = [rCrtView \text{ EXCEPT } ![S][r] = msg.v] \\
& \wedge rViewOnDisk' = [rViewOnDisk \text{ EXCEPT } ![S][r] = msg.v] \\
& \wedge Send([type \mapsto \text{"DO-VIEW-CHANGE"}, \\
& \quad v \mapsto msg.v, \\
& \quad lv \mapsto rLastView[S][r], \\
& \quad r \mapsto rRecord[S][r], \\
& \quad src \mapsto r, \\
& \quad dst \mapsto Replicas]) \\
& \wedge \text{UNCHANGED } \langle cVars, rLastView, rRecord, rNonce, \\
& \quad aVars, gViewChangesNo \rangle
\end{aligned}$$

Note: Assume one reply for view change per replica in Q

(in *ReplicaReceiveDoViewChange* we keep only the most recent reply)

$ReplicaSendStartView(r) \triangleq$

$\wedge \exists Q \in Quorums :$

$\wedge \forall r1 \in Q :$

$\wedge \forall r2 \in Q : \exists rr, pr \in rViewReplies[S][r] :$

$\wedge rr.type = \text{"do-view-change"}$

$\wedge pr.type = \text{"do-view-change"}$

$\wedge rr.src = r1$

$\wedge pr.src = r2$

$$\begin{aligned}
& \wedge rr.v = pr.v \\
& \wedge rr.v \geq rCrtView[S][r] \\
& \text{received at a least a quorum of replies} \\
\wedge \text{LET} \\
& A \triangleq \\
& \quad \text{set of all do-view-change replies from } Q, \\
& \quad \{x \in rViewReplies[S][r] : \wedge x.src \in Q \\
& \quad \quad \wedge x.type = \text{"do-view-change"}\} \\
& B \triangleq \\
& \quad \text{keep only the replies from the maximum view} \\
& \quad \{x \in A : \forall y \in A : y.lv \leq x.lv\} \\
& C \triangleq \\
& \quad \text{set of all records received in replies in } B \\
& \quad \text{UNION } \{x.r : x \in B\} \\
& recoveredConsensusOps_R \triangleq \\
& \quad \text{any finalized consensus operation (in at least one record,} \\
& \quad \text{in the maximum latest view)} \\
& \quad \{[msgid \mapsto y.msgid, op \mapsto y.op, res \mapsto y.res] : y \in \\
& \quad \{x \in C : \\
& \quad \quad \wedge x.op.type = \text{"Consensus"} \\
& \quad \quad \wedge x.status = \text{"FINALIZED"}\}\} \\
& recoveredConsensusOps_d \triangleq \\
& \quad \text{any consensus operation found in at least a majority of a Quorum} \\
& \quad \{[msgid \mapsto y.msgid, op \mapsto y.op, res \mapsto y.res] : y \in \\
& \quad \{x \in C : \\
& \quad \quad \wedge x.op.type = \text{"Consensus"} \\
& \quad \quad \wedge x.status = \text{"TENTATIVE"} \\
& \quad \quad \wedge \exists P \in SuperQuorums : \\
& \quad \quad \quad \forall replica \in Q \cap P : \\
& \quad \quad \quad \exists reply \in B : \\
& \quad \quad \quad \quad \wedge reply.src = replica \\
& \quad \quad \quad \quad \wedge x \in reply.r\}\} \setminus recoveredConsensusOps_R \\
& recoveredConsensusOps_u \triangleq \\
& \quad \text{the rest of consensus ops found in at least one record} \\
& \quad \text{(discard the result)} \\
& \quad \{[msgid \mapsto z.msgid, op \mapsto z.op] : z \in \\
& \quad ((\{[msgid \mapsto y.msgid, op \mapsto y.op, res \mapsto y.res] : y \in \\
& \quad \quad \{x \in C : x.op.type = \text{"Consensus"}\}\} \\
& \quad \quad \setminus recoveredConsensusOps_d) \setminus recoveredConsensusOps_R)\} \\
& recoveredInconsistentOps_R \triangleq \\
& \quad \text{any inconsistent operation found in any received record}
\end{aligned}$$

$$\begin{aligned}
& \{[msgid \mapsto y.msgid, op \mapsto y.op] : y \in \\
& \quad \{x \in C : x.op.type = \text{"Inconsistent"}\}\} \\
mergedRecordInconsistent & \triangleq \\
& \{x \in Merge(recoveredConsensusOps_R \\
& \quad \cup recoveredInconsistentOps_R, \\
& \quad recoveredConsensusOps_d, \\
& \quad recoveredConsensusOps_u) : x.op.type = \text{"Inconsistent"}\} \\
mergedRecordConsensus & \triangleq \\
& \{x \in Merge(recoveredConsensusOps_R \\
& \quad \cup recoveredInconsistentOps_R, \\
& \quad recoveredConsensusOps_d, \\
& \quad recoveredConsensusOps_u) : x.op.type = \text{"Consensus"}\} \\
masterRecord & \triangleq \\
& \{[msgid \mapsto x.msgid, \\
& \quad op \mapsto x.op, \\
& \quad status \mapsto \text{"FINALIZED"}] : \\
& \quad x \in mergedRecordInconsistent\} \\
& \cup \\
& \{[msgid \mapsto x.msgid, \\
& \quad op \mapsto x.op, \\
& \quad res \mapsto x.res, \\
& \quad status \mapsto \text{"FINALIZED"}] : \\
& \quad x \in mergedRecordConsensus\} \\
v_new & \triangleq \\
& \quad \text{the one decided by quorum } Q \\
& \quad \text{CHOOSE } v \in \{x.v : x \in A\} : \text{TRUE} \\
\text{IN} & \\
& \wedge rRecord' = [rRecord \text{ EXCEPT } ![S][r] = masterRecord] \\
& \wedge Sync(masterRecord) \\
& \wedge Send([type \mapsto \text{"START-VIEW"}, \\
& \quad v \mapsto v_new, \\
& \quad r \mapsto masterRecord, \\
& \quad src \mapsto r]) \\
& \wedge rCrtView' = [rCrtView \text{ EXCEPT } ![S][r] = v_new] \\
& \wedge rLastView' = [rLastView \text{ EXCEPT } ![S][r] = v_new] \\
& \quad \wedge \text{Assert}(Cardinality(masterRecord) = 0, \text{"Should fail -- ReplicaSendStartView"}) \\
& \wedge rState' = [rState \text{ EXCEPT } ![S][r] = \text{"NORMAL"}] \\
& \wedge rViewReplies' = [rViewReplies \text{ EXCEPT } ![S][r] = \{\}] \\
& \wedge \text{UNCHANGED } \langle rNonce, rViewOnDisk, cVars, aVars, gViewChangesNo \rangle \\
& \quad \backslash * \text{ A replica receives a start view message} \\
ReplicaReceiveStartView(r) & \triangleq
\end{aligned}$$

```

 $\wedge \exists msg \in sentMsg[S] :$ 
 $\wedge msg.type = \text{"START-VIEW"}$ 
 $\wedge \vee \wedge rState[S][r] = \text{"NORMAL"}$ 
 $\wedge msg.v > rCrtView[S][r]$ 
 $\vee \wedge \vee rState[S][r] = \text{"VIEW-CHANGING"}$ 
 $\vee rState[S][r] = \text{"RECOVERING"}$ 
 $\wedge msg.v \geq rCrtView[S][r]$ 
 $\wedge rCrtView' = [rCrtView \text{ EXCEPT } ![S][r] = msg.v]$ 
 $\wedge rLastView' = [rLastView \text{ EXCEPT } ![S][r] = msg.v]$ 
 $\wedge rRecord' = [rRecord \text{ EXCEPT } ![S][r] = msg.r]$ 
 $\wedge Sync(msg.r)$ 
 $\wedge$  Check if the operations received in the master record
must be added to the aSuccessful
LET
 $successfulOps \triangleq \{x \in msg.r : \exists Q \in Quorums :$ 
 $\vee r1 \in Q :$ 
 $\vee x \in rRecord[S][r1]$ 
 $\vee x \in arRecord[S][r1]$ 
 $\vee r1 = r\}$ 
IN
 $aSuccessful' = aSuccessful \cup$ 
 $\{[mid \mapsto x.msgid,$ 
 $op \mapsto x.op,$ 
 $res \mapsto x.res] :$ 
 $x \in$ 
 $\{y \in successfulOps :$ 
 $y.op.type = \text{"Consensus"}\}\}$ 
 $\cup$ 
 $\{[mid \mapsto x.msgid,$ 
 $op \mapsto x.op] :$ 
 $x \in$ 
 $\{y \in successfulOps :$ 
 $y.op.type = \text{"Inconsistent"}\}\}$ 
 $\wedge rViewOnDisk' = [rViewOnDisk \text{ EXCEPT } ![S][r] = msg.v + 1]$ 
 $\wedge rState' = [rState \text{ EXCEPT } ![S][r] = \text{"NORMAL"}]$ 
 $\wedge rViewReplies' = [rViewReplies \text{ EXCEPT } ![S][r] = \{\}]$ 
 $\wedge \text{UNCHANGED } \langle rNonce, cVars, arRecord, aVisibility, oVars \rangle$ 
\ * A replica fails and loses everything except the view number
\ * The view number has been written to disk
 $ReplicaFail(r) \triangleq$ 
 $\wedge rState' = [rState \text{ EXCEPT } ![S][r] = \text{"FAILED"}]$ 
 $\wedge rRecord' = [rRecord \text{ EXCEPT } ![S][r] = \{\}]$ 
 $\wedge arRecord' = [arRecord \text{ EXCEPT } ![S][r] = rRecord[S][r]]$ 
save record only for

```


invariant purposes

$$\begin{aligned}
& \wedge rCrtView' = [rCrtView \text{ EXCEPT } ![S][r] = 0] \\
& \wedge rLastView' = [rLastView \text{ EXCEPT } ![S][r] = 0] \\
& \wedge rViewReplies' = [rViewReplies \text{ EXCEPT } ![S][r] = \{\}] \\
& \wedge \text{UNCHANGED } \langle rViewOnDisk, rNonce, cVars, aSuccessful, aVisibility, oVars \rangle \\
& \wedge \text{We assume less than } f \text{ replicas fail simultaneously} \\
& \quad Cardinality(\{re \in Replicas : \\
& \qquad \vee rState[S][re] = \text{"FAILED"} \\
& \qquad \vee rState[S][re] = \text{"RECOVERING"}\}) < f
\end{aligned}$$

High-Level Actions

$$\begin{aligned}
ClientAction(c) & \triangleq \\
& \vee \wedge cState[c] = \text{"NORMAL"} \\
& \quad \wedge \quad \vee ClientRequest(c) \setminus * \text{some client tries to replicate commit an operation} \\
& \quad \vee ClientReceiveReply(c) \quad \text{some client receives a reply from a replica} \\
& \quad \vee ClientReceiveConfirm(c) \quad \text{some client receives a confirm from a replica} \\
& \quad \vee ClientFail(c) \quad \setminus * \text{some client fails} \\
& \quad \vee ClientSendFinalize(c) \quad \text{an operation is successful at some client} \\
& \quad \vee ClientFinalizeOp(c) \quad \text{an operation was finalized at some client} \\
& \vee \wedge cState[c] = \text{"FAILED"} \\
& \quad \wedge \vee ClientRecover(c) \\
\\
ReplicaAction(r) & \triangleq \\
& \vee \wedge rState[S][r] = \text{"NORMAL"} \\
& \quad \wedge \vee ReplicaReceiveRequest(r) \quad \text{some replica sends a reply to a } PROPOSE \text{ msg} \\
& \quad \vee ReplicaReceiveFinalize(r) \\
& \quad \vee ReplicaSendDoViewChange(r) \\
& \quad \vee ReplicaReceiveDoViewChange(r) \\
& \quad \vee ReplicaReceiveStartView(r) \\
& \quad \vee ReplicaFail(r) \quad \text{some replica fails} \\
& \vee \wedge rState[S][r] = \text{"FAILED"} \\
& \quad \wedge \vee ReplicaSendDoViewChange(r) \quad \text{start view-change protocol} \\
& \vee \wedge rState[S][r] = \text{"RECOVERING"} \\
& \quad \wedge \vee ReplicaSendDoViewChange(r) \quad \text{re-start view-change protocol (assume a} \\
& \qquad \qquad \qquad \text{timeout and still no response from the new leader)} \\
& \quad \vee ReplicaReceiveStartView(r) \\
& \vee \wedge rState[S][r] = \text{"VIEW-CHANGING"} \\
& \quad \wedge \vee ReplicaSendDoViewChange(r) \\
& \quad \vee ReplicaReceiveDoViewChange(r) \\
& \quad \vee ReplicaSendStartView(r) \\
& \quad \vee ReplicaReceiveStartView(r) \\
& \quad \vee ReplicaFail(r) \\
\\
Next & \triangleq
\end{aligned}$$

$$\begin{aligned}
& \vee \exists c \in \text{Clients} : \text{ClientAction}(c) \\
& \vee \exists r \in \text{Replicas} : \text{ReplicaAction}(r) \\
& \vee \setminus * \text{ Avoid deadlock by termination} \\
& (\forall i \in 1 \dots \text{Cardinality}(\text{Replicas}) : r\text{LastView}[S][i] = \text{max_vc}) \vee \text{UNCHANGED } \langle \text{vars} \rangle
\end{aligned}$$

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

$$\begin{aligned}
& \text{FaultTolerance} \triangleq \\
& \wedge \forall \text{successfulOp} \in \text{aSuccessful}, Q \in \text{Quorums} : \\
& \quad (\forall r \in Q : r\text{State}[S][r] = \text{"NORMAL"} \vee r\text{State}[S][r] = \text{"VIEW-CHANGING"}) \\
& \Rightarrow (\exists p \in Q : \exists \text{rec} \in r\text{Record}[S][p] : \\
& \quad \wedge \text{successfulOp.mid} = \text{rec.msgid} \\
& \quad \wedge \text{successfulOp.op} = \text{rec.op}) \quad \text{Not necessarily same result}
\end{aligned}$$

\ * Modification History
\ * Last modified *Thu Feb 15 10:51:10 PST 2018* by *o80053527*
\ * Last modified *Fri Aug 28 10:58:02 PDT 2015* by *aaasz*
\ * Created *Fri Dec 12 17:42:14 PST 2014* by *aaasz*