

EXTENDS *Naturals, Sequences, FiniteSets, TLC*

The set of *Paxos* replicas

CONSTANT *Replicas*

The set of *Paxos* clients

CONSTANT *Clients*

The set of possible values

CONSTANT *Values*

An empty value

CONSTANT *Nil*

Request/response types

CONSTANTS

MClientRequest,
MClientReply,
MRepairRequest,
MRepairReply,
MAbortRequest,
MAbortReply,
MViewChangeRequest,
MViewChangeReply,
MStartViewRequest

Replica statuses

CONSTANTS

SNormal,
SAborting,
SViewChange

Entry types

CONSTANTS

TValue,
TNoOp

Message schemas

$ViewIDs \triangleq [viewID \mapsto n \in (1 \dots)]$

ClientRequest

[src $\mapsto c \in Clients$,
dest $\mapsto r \in Replicas$,
type $\mapsto MClientRequest$,
viewID $\mapsto i \in ViewIDs$,
value $\mapsto v \in Values$,
seqNum $\mapsto s \in (1 \dots)$]

$timestamp \mapsto t \in (1 \dots)$

ClientReply

[$src \mapsto r \in Replicas,$
 $dest \mapsto c \in Clients,$
 $type \mapsto MClientReply,$
 $viewID \mapsto i \in ViewIDs,$
 $value \mapsto v \in Values,$
 $seqNum \mapsto s \in (1 \dots),$
 $index \mapsto i \in (1 \dots),$
 $succeeded \mapsto TRUE \vee FALSE]$

RepairRequest

[$src \mapsto r \in Replicas,$
 $dest \mapsto r \in Replicas,$
 $type \mapsto MRepairRequest,$
 $viewID \mapsto i \in ViewIDs,$
 $client \mapsto c \in Clients,$
 $seqNum \mapsto s \in (1 \dots),$
 $timestamp \mapsto t \in (1 \dots)]$

RepairReply

[$src \mapsto r \in Replicas,$
 $dest \mapsto r \in Replicas,$
 $type \mapsto MRepairReply,$
 $viewID \mapsto i \in ViewIDs,$
 $client \mapsto c \in Clients,$
 $seqNum \mapsto s \in (1 \dots),$
 $value \mapsto v \in Values,$
 $timestamp \mapsto t \in (1 \dots)]$

AbortRequest

[$src \mapsto r \in Replicas,$
 $dest \mapsto r \in Replicas,$
 $type \mapsto MAbortRequest,$
 $viewID \mapsto i \in ViewIDs,$
 $client \mapsto c \in Clients,$
 $seqNum \mapsto s \in (1 \dots),$
 $timestamp \mapsto t \in (1 \dots)]$

AbortReply

[$src \mapsto r \in Replicas,$
 $dest \mapsto r \in Replicas,$
 $type \mapsto MAbortReply,$
 $viewID \mapsto i \in ViewIDs,$
 $client \mapsto c \in Clients,$
 $seqNum \mapsto s \in (1 \dots)]$

ViewChangeRequest

[$src \mapsto r \in Replicas,$
 $dest \mapsto r \in Replicas,$
 $type \mapsto MViewChangeRequest,$
 $viewID \mapsto i \in ViewIDs]$

```

ViewChangeReply
[src      ↦ r ∈ Replicas,
 dest     ↦ r ∈ Replicas,
 type     ↦ MViewChangeReply,
 viewID   ↦ i ∈ ViewIDs,
 lastViewID ↦ i ∈ (1 ..),
 log      ↦ [c ∈ Clients ↦ [ i ∈ 1 .. (1 ..) ↦ [
                      value ↦ v ∈ Values, timestamp ↦ (1 .. )]]]]

StartViewRequest
[src      ↦ r ∈ Replicas,
 dest     ↦ r ∈ Replicas,
 type     ↦ MStartViewRequest,
 viewID   ↦ i ∈ ViewIDs,
 timestamp ↦ t ∈ (1 ..), log      ↦ [c ∈ Clients ↦ [
                      i ∈ 1 .. (1 ..) ↦ [ value ↦ v ∈ Values, timestamp ↦ (1 .. )]]]]

```

A sequence of replicas used for deterministic primary election

VARIABLE *replicas*

$globalVars \triangleq \langle replicas \rangle$

The set of all messages on the network

VARIABLE *messages*

The total number of messages sent

VARIABLE *messageCount*

$messageVars \triangleq \langle messages, messageCount \rangle$

Local client state

Strictly increasing representation of synchronized time

VARIABLE *cTime*

The highest known view *ID* for a client

VARIABLE *cViewID*

The current sequence number for a client

VARIABLE *cSeqNum*

A client response buffer

VARIABLE *cReps*

A set of all commits - used for model checking

VARIABLE *cCommits*

$clientVars \triangleq \langle cTime, cViewID, cSeqNum, cReps, cCommits \rangle$

Local replica state

The current status of a replica
 VARIABLE $rStatus$

The current view ID for a replica
 VARIABLE $rViewID$

A replica's commit log
 VARIABLE $rLog$

The current log index for a replica
 VARIABLE $rIndex$

The current sequence number for each session
 VARIABLE $rSeqNum$

The highest known timestamp for all sessions
 VARIABLE $rTimestamp$

The last known normal view
 VARIABLE $rLastViewID$

The set of received view change responses
 VARIABLE $rViewChanges$

The point ($client + \text{sequence number}$) in the log currently being aborted
 VARIABLE $rAbortPoint$

The set of abort responses received
 VARIABLE $rAbortReps$

$replicaVars \triangleq \langle rStatus, rViewID, rLog, rIndex, rSeqNum, rTimestamp, rLastViewID, rViewChanges, rAbortPoint, rAbortReps \rangle$

$vars \triangleq \langle globalVars, messageVars, clientVars, replicaVars \rangle$

This section provides helpers for the spec.

Creates a sequence from set 'S'
 RECURSIVE $SeqFromSet(-)$
 $SeqFromSet(S) \triangleq$
 IF $S = \{\}$ THEN
 $\langle \rangle$
 ELSE LET $x \triangleq \text{CHOOSE } x \in S : \text{TRUE}$
 IN $\langle x \rangle \circ SeqFromSet(S \setminus \{x\})$

Selects an element of set 'S'
 $Pick(S) \triangleq \text{CHOOSE } s \in S : \text{TRUE}$

RECURSIVE $SetReduce(-, -, -)$
 $SetReduce(Op(-, -), S, value) \triangleq$
 IF $S = \{\}$ THEN
 $value$
 ELSE
 LET $s \triangleq Pick(S)$
 IN $SetReduce(Op, S \setminus \{s\}, Op(s, value))$

Computes the greatest vlue in set 'S'
 $Max(S) \triangleq \text{CHOOSE } x \in S : \forall y \in S : x \geq y$

Computes the sum of numbers in set 'S'
 $Sum(S) \triangleq \text{LET } _op(a, b) \triangleq a + b$
 IN $SetReduce(_op, S, 0)$

A boolean indicating whether the given set is a quorum
 $IsQuorum(s) \triangleq Cardinality(s) * 2 \geq Cardinality(Replicas)$

The set of all quorums
 $Quorums \triangleq \{r \in \text{SUBSET } Replicas : IsQuorum(r)\}$

The primary for view 'v'
 $Primary(v) \triangleq replicas[(v \% Len(replicas)) + (\text{IF } v \geq Len(replicas) \text{ THEN } 1 \text{ ELSE } 0)]$

A boolean indicating whether replica 'r' is the primary for the current view
 $IsPrimary(r) \triangleq Primary(rViewID[r]) = r$

This section models the network.

Send a set of messages
 $Sends(ms) \triangleq$
 $\wedge \text{messages}' = messages \cup ms$
 $\wedge \text{messageCount}' = \text{messageCount} + Cardinality(ms)$

Send a message
 $Send(m) \triangleq Sends(\{m\})$

Reply to a message with a set of responses
 $Replies(req, reps) \triangleq$
 $\wedge \text{messages}' = (messages \cup reps) \setminus \{req\}$
 $\wedge \text{messageCount}' = \text{messageCount} + Cardinality(reps)$

Reply to a message
 $Reply(req, resp) \triangleq Replies(req, \{resp\})$

Discard a message
 $Discard(m) \triangleq$
 $\wedge \text{messages}' = messages \setminus \{m\}$

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

This section models client requests.

Client 'c' sends value 'v' to all replicas

Client requests are ordered globally using physical timestamps and locally (within the client) using client sequence numbers. Sequence numbers are sequential and unique within each view.

When the client sends a request it generates a new timestamp. Physical timestamps are modeled here as a strictly increasing global clock simulating synchronized system clocks. The sequence number for the client is also incremented and sent with the request.

$$\begin{aligned} \text{ClientRequest}(c, v) &\triangleq \\ &\wedge cTime' = cTime + 1 \\ &\wedge cSeqNum' = [cSeqNum \text{ EXCEPT } ![c] = cSeqNum[c] + 1] \\ &\wedge \text{Sends}(\{[src \mapsto c, \\ &\quad dest \mapsto r, \\ &\quad type \mapsto MClientRequest, \\ &\quad viewID \mapsto cViewID[c], \\ &\quad seqNum \mapsto cSeqNum'[c], \\ &\quad value \mapsto v, \\ &\quad timestamp \mapsto cTime'] : r \in Replicas\}) \\ &\wedge \text{UNCHANGED } \langle globalVars, replicaVars, cViewID, cReps, cCommits \rangle \end{aligned}$$

Client 'c' handles a response 'm' from replica 'r'

When a response is received by the client, if the client is still in the request view it can process the response. The client is responsible for determining commitment by counting responses for each sequence number. Once a response is received from a majority of the replicas including the primary replica, the response is committed. Committed responses are stored in a history variable for checking against invariants.

$$\begin{aligned} \text{HandleClientReply}(c, r, m) &\triangleq \\ &\wedge \vee \wedge m.viewID = cViewID[c] \\ &\wedge cReps' = [cReps \text{ EXCEPT } ![c] = cReps[c] \cup \{m\}] \\ &\wedge \text{LET} \\ &\quad okReps \triangleq \{n \in cReps'[c] : n.seqNum = m.seqNum \wedge n.succeeded\} \\ &\quad hasPrimary \triangleq \exists n \in okReps : Primary(n.viewID) = n.src \\ &\text{IN} \\ &\quad \text{IF } hasPrimary \text{ THEN} \\ &\quad \quad \text{LET} \\ &\quad \quad \quad primaryRep \triangleq \text{CHOOSE } n \in okReps : n.src = Primary(n.viewID) \\ &\quad \quad \quad matchReps \triangleq \{n \in okReps : n.index = primaryRep.index\} \\ &\quad \quad \quad isCommitted \triangleq \{n.src : n \in matchReps\} \in Quorums \\ &\quad \text{IN} \\ &\quad \quad \text{IF } isCommitted \text{ THEN} \\ &\quad \quad \quad cCommits' = [cCommits \text{ EXCEPT } ![c] = cCommits[c] \cup \{primaryRep\}] \end{aligned}$$

$$\begin{aligned}
& \text{ELSE} \\
& \quad \text{UNCHANGED } \langle cCommits \rangle \\
& \text{ELSE} \\
& \quad \text{UNCHANGED } \langle cCommits \rangle \\
& \quad \wedge \text{UNCHANGED } \langle cViewID, cSeqNum \rangle \\
& \quad \vee \wedge m.viewID > cViewID[c] \\
& \quad \quad \wedge cViewID' = [cViewID \text{ EXCEPT } ![c] = m.viewID] \\
& \quad \quad \wedge cSeqNum' = [cSeqNum \text{ EXCEPT } ![c] = 0] \\
& \quad \quad \wedge cReps' = [cReps \text{ EXCEPT } ![c] = \{\}] \\
& \quad \quad \wedge \text{UNCHANGED } \langle cCommits \rangle \\
& \quad \vee \wedge m.viewID < cViewID[c] \\
& \quad \quad \wedge \text{UNCHANGED } \langle cViewID, cSeqNum, cReps, cCommits \rangle \\
& \wedge Discard(m) \\
& \wedge \text{UNCHANGED } \langle globalVars, replicaVars, cTime \rangle
\end{aligned}$$

This section models the replica protocol.

Replica 'r' requests a repair of the client 'c' request 'm'

$$\begin{aligned}
Repair(r, c, m) & \triangleq \\
& \wedge Replies(m, \{[src \mapsto r, \\
& \quad \quad \quad dest \mapsto d, \\
& \quad \quad \quad type \mapsto MRepairRequest, \\
& \quad \quad \quad viewID \mapsto rViewID[r], \\
& \quad \quad \quad client \mapsto c, \\
& \quad \quad \quad seqNum \mapsto m.seqNum, \\
& \quad \quad \quad timestamp \mapsto m.timestamp] : d \in Replicas\})
\end{aligned}$$

Replica 'r' aborts the client 'c' request 'm'

$$\begin{aligned}
Abort(r, c, m) & \triangleq \\
& \wedge IsPrimary(r) \\
& \wedge rStatus[r] = SNormal \\
& \wedge rStatus' = [rStatus \text{ EXCEPT } ![r] = SAborting] \\
& \wedge rAbortReps' = [rAbortReps \text{ EXCEPT } ![r] = \{\}] \\
& \wedge rAbortPoint' = [rAbortPoint \text{ EXCEPT } ![r] = [client \mapsto c, seqNum \mapsto m.seqNum]] \\
& \wedge Replies(m, \{[src \mapsto r, \\
& \quad \quad \quad dest \mapsto d, \\
& \quad \quad \quad type \mapsto MAbortRequest, \\
& \quad \quad \quad viewID \mapsto rViewID[r], \\
& \quad \quad \quad client \mapsto c, \\
& \quad \quad \quad seqNum \mapsto m.seqNum, \\
& \quad \quad \quad timestamp \mapsto m.timestamp] : d \in Replicas\})
\end{aligned}$$

Replica 'r' handles client 'c' request 'm'

Client requests with a view *ID* not matching the replica's view are rejected.

Clients reset their sequence number of

For requests in the correct view, the request must be sequential and linear to be appended to the *log*. That is, the request must have a 'seqNum' that is 1 + the prior 'seqNum' for the client, and the 'timestamp' must be greater than all prior timestamps in the *log*. This is necessary to ensure the primary *log* does not change when requests are reordered. The client can retry requests that are reordered with a new sequence number and timestamp.

To maintain consistency within the *log*, a separate sequence is maintained for each session (client), and each sequence number is assigned to a unique position in the session *log*. Session logs are logically merged into a totally ordered *log* using the request timestamps.

When a sequence number is skipped, the primary must commit a *TNoOp* entry to the *log*. It does so by running the *AbortRequest* protocol.

When a sequence number is skipped on a non-primary replica, the replica attempts to recover the request using the *RepairRequest* protocol.

$$\begin{aligned}
& \text{HandleClientRequest}(r, c, m) \triangleq \\
& \quad \wedge r\text{Status}[r] = S\text{Normal} \\
& \quad \wedge \vee \wedge m.\text{viewID} = r\text{ViewID}[r] \\
& \quad \wedge \text{LET} \\
& \quad \quad \begin{aligned}
& \text{index} && \triangleq r\text{Index}[r] + 1 \\
& \text{lastTimestamp} && \triangleq r\text{Timestamp}[r] \\
& \text{isSequential} && \triangleq m.\text{seqNum} = r\text{SeqNum}[r][c] + 1 \\
& \text{isLinear} && \triangleq m.\text{timestamp} > \text{lastTimestamp} \\
& \text{entry} && \triangleq [\text{type} \mapsto T\text{Value}, \\
& && \text{index} \mapsto \text{index}, \\
& && \text{value} \mapsto m.\text{value}, \\
& && \text{timestamp} \mapsto m.\text{timestamp}] \\
& \text{append}(e) && \triangleq [r\text{Log} \text{ EXCEPT } ![r] = [r\text{Log}[r] \text{ EXCEPT} \\
& && \quad \quad \quad ![c] = \text{Append}(r\text{Log}[r][c], e))]
\end{aligned} \\
& \text{IN} \\
& \quad \vee \wedge \text{isSequential} \\
& \quad \wedge \text{isLinear} \\
& \quad \wedge r\text{Log}' = \text{append}(\text{entry}) \\
& \quad \wedge r\text{Index}' = [r\text{Index} \text{ EXCEPT } ![r] = \text{index}] \\
& \quad \wedge r\text{SeqNum}' = [r\text{SeqNum} \text{ EXCEPT } ![r] = \\
& \quad \quad \quad [r\text{SeqNum}[r] \text{ EXCEPT } ![c] = m.\text{seqNum}]] \\
& \quad \wedge r\text{Timestamp}' = [r\text{Timestamp} \text{ EXCEPT } ![r] = m.\text{timestamp}] \\
& \quad \wedge \text{Reply}(m, [\text{src} \mapsto r, \\
& \quad \quad \quad \text{dest} \mapsto c, \\
& \quad \quad \quad \text{type} \mapsto M\text{ClientReply}, \\
& \quad \quad \quad \text{viewID} \mapsto r\text{ViewID}[r], \\
& \quad \quad \quad \text{seqNum} \mapsto m.\text{seqNum}, \\
& \quad \quad \quad \text{index} \mapsto \text{index}, \\
& \quad \quad \quad \text{timestamp} \mapsto m.\text{timestamp}, \\
& \quad \quad \quad \text{value} \mapsto m.\text{value}, \\
& \quad \quad \quad \text{succeeded} \mapsto \text{TRUE}])
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{UNCHANGED } \langle rStatus, rAbortPoint, rAbortReps \rangle \\
\vee & \wedge \vee \wedge \neg isSequential \\
& \wedge m.seqNum > rSeqNum[r][c] + 1 \\
& \vee \neg isLinear \\
& \wedge \vee \wedge IsPrimary(r) \\
& \wedge Abort(r, c, m) \\
& \vee \wedge \neg IsPrimary(r) \\
& \wedge Repair(r, c, m) \\
& \wedge \text{UNCHANGED } \langle rStatus, rAbortPoint, rAbortReps \rangle \\
& \wedge \text{UNCHANGED } \langle rLog, rIndex, rSeqNum, rTimestamp \rangle \\
\vee & \wedge m.viewID < rViewID[r] \\
& \wedge Reply(m, [src \mapsto r, \\
& \quad dest \mapsto c, \\
& \quad type \mapsto MClientReply, \\
& \quad viewID \mapsto rViewID[r], \\
& \quad seqNum \mapsto m.seqNum, \\
& \quad succeeded \mapsto FALSE]) \\
& \wedge \text{UNCHANGED } \langle rStatus, rLog, rIndex, rSeqNum, rTimestamp, rAbortPoint, rAbortReps \rangle \\
& \wedge \text{UNCHANGED } \langle globalVars, clientVars, rViewID, rLastViewID, rViewChanges \rangle
\end{aligned}$$

Replica 'r' handles replica 's' repair request 'm'

When a repair request is received, if the requested sequence number is in the session

log, the entry is returned. Otherwise, the primary aborts the request.

$$\begin{aligned}
HandleRepairRequest(r, s, m) & \triangleq \\
& \wedge m.viewID = rViewID[r] \\
& \wedge IsPrimary(r) \\
& \wedge rStatus[r] = SNormal \\
& \wedge LET \ offset \triangleq Len(rLog[r][m.client]) - (rSeqNum[r][m.client] - m.seqNum) \\
& IN \\
& \vee \wedge offset \leq Len(rLog[r][m.client]) \\
& \wedge Reply(m, [src \mapsto r, \\
& \quad dest \mapsto s, \\
& \quad type \mapsto MRepairReply, \\
& \quad viewID \mapsto rViewID[r], \\
& \quad client \mapsto m.client, \\
& \quad seqNum \mapsto m.seqNum, \\
& \quad value \mapsto rLog[r][m.client][offset].value, \\
& \quad timestamp \mapsto rLog[r][m.client][offset].timestamp]) \\
& \wedge \text{UNCHANGED } \langle rStatus, rAbortPoint, rAbortReps \rangle \\
& \vee \wedge offset = Len(rLog[r][m.client]) + 1 \\
& \wedge Abort(r, m.client, m) \\
& \wedge \text{UNCHANGED } \langle globalVars, clientVars, rLog, rIndex, rSeqNum, rTimestamp, rViewID, rLastViewID, rViewChanges \rangle
\end{aligned}$$

Replica 'r' handles replica 's' repair response 'm'

Repair responses are handled like client requests.

$HandleRepairReply(r, s, m) \triangleq$
 $HandleClientRequest(r, m.client, [m \text{ EXCEPT } !.src = m.client])$

Replica 'r' handles replica 's' abort request 'm'

If the aborted sequence number is in the session *log*, the entry is replaced with a no-op entry. If the sequence number can be appended to the *log*, it is.

$HandleAbortRequest(r, s, m) \triangleq$
 $\wedge m.viewID = rViewID[r]$
 $\wedge rStatus[r] \in \{SNormal, SAborting\}$
 $\wedge \text{LET}$
 $\quad offset \triangleq Len(rLog[r][m.client]) - (rSeqNum[r][m.client] - m.seqNum)$
 $\quad entry \triangleq [type \mapsto TNoOp, value \mapsto Nil, timestamp \mapsto 0]$
 IN
 $\wedge \vee \wedge offset \leq Len(rLog[r][m.client])$
 $\quad \wedge rLog' = [rLog \text{ EXCEPT } ![r] = [$
 $\quad \quad rLog[r] \text{ EXCEPT } ![m.client] = [$
 $\quad \quad rLog[r][m.client] \text{ EXCEPT } ![offset] = [$
 $\quad \quad rLog[r][m.client][offset] \text{ EXCEPT } !.type = TNoOp]]]$
 $\quad \wedge \text{UNCHANGED } \langle rTimestamp, rSeqNum \rangle$
 $\vee \wedge offset = Len(rLog[r][m.client]) + 1$
 $\quad \wedge rLog' = [rLog \text{ EXCEPT } ![r] = [$
 $\quad \quad rLog[r] \text{ EXCEPT } ![m.client] =$
 $\quad \quad Append(rLog[r][m.client], entry)]]$
 $\quad \wedge rTimestamp' = [rTimestamp \text{ EXCEPT } ![r] = Max(\{rTimestamp[r], m.timestamp\})]$
 $\quad \wedge rSeqNum' = [rSeqNum \text{ EXCEPT } ![r] = [$
 $\quad \quad rSeqNum[r] \text{ EXCEPT } ![m.client] = m.seqNum]]$
 $\wedge Replies(m, \{[src \mapsto r,$
 $\quad dest \mapsto Primary(rViewID[r]),$
 $\quad type \mapsto MAbortReply,$
 $\quad viewID \mapsto rViewID[r],$
 $\quad client \mapsto m.client,$
 $\quad seqNum \mapsto m.seqNum],$
 $\quad [src \mapsto r,$
 $\quad dest \mapsto m.client,$
 $\quad type \mapsto MClientReply,$
 $\quad viewID \mapsto rViewID[r],$
 $\quad seqNum \mapsto m.seqNum,$
 $\quad succeeded \mapsto FALSE]\})$
 $\wedge \text{UNCHANGED } \langle globalVars, clientVars, rStatus, rAbortPoint,$
 $\quad rAbortReps, rViewID, rLastViewID, rViewChanges \rangle$

Replica 'r' handles replica 's' repair response 'm'

$HandleAbortReply(r, s, m) \triangleq$
 $\wedge rStatus[r] = SAborting$
 $\wedge m.viewID = rViewID[r]$

$$\begin{aligned}
& \wedge \text{IsPrimary}(r) \\
& \wedge m.\text{seqNum} = r\text{AbortPoint}[r].\text{seqNum} \\
& \wedge r\text{AbortReps}' = [r\text{AbortReps} \text{ EXCEPT } ![r] = r\text{AbortReps}[r] \cup \{m\}] \\
& \wedge \text{LET } \text{reps} \triangleq \{ \text{res}.\text{src} : \text{res} \in \{ \text{resp} \in r\text{AbortReps}'[r] : \\
& \quad \wedge \text{resp}.\text{viewID} = r\text{ViewID}[r] \\
& \quad \wedge \text{resp}.\text{client} = r\text{AbortPoint}[r].\text{client} \\
& \quad \wedge \text{resp}.\text{seqNum} = r\text{AbortPoint}[r].\text{seqNum} \} \} \\
& \quad \text{isQuorum} \triangleq r \in \text{reps} \wedge \text{reps} \in \text{Quorums} \\
& \text{IN} \\
& \quad \vee \wedge \text{isQuorum} \\
& \quad \wedge r\text{Status}' = [r\text{Status} \text{ EXCEPT } ![r] = S\text{Normal}] \\
& \quad \vee \wedge \neg \text{isQuorum} \\
& \quad \wedge \text{UNCHANGED } \langle r\text{Status} \rangle \\
& \wedge \text{UNCHANGED } \langle \text{globalVars}, \text{messageVars}, \text{clientVars}, r\text{Log}, r\text{Index}, r\text{SeqNum}, \\
& \quad r\text{Timestamp}, r\text{AbortPoint}, r\text{ViewID}, r\text{ViewChanges}, r\text{LastViewID} \rangle
\end{aligned}$$

Replica 'r' requests a view change

The view change is requested by sending a *ViewChangeRequest* to each replica.

$$\begin{aligned}
\text{ChangeView}(r) & \triangleq \\
& \wedge \text{Sends}(\{ \text{src} \mapsto r, \\
& \quad \text{dest} \mapsto d, \\
& \quad \text{type} \mapsto M\text{ViewChangeRequest}, \\
& \quad \text{viewID} \mapsto r\text{ViewID}[r] + 1 : d \in \text{Replicas} \}) \\
& \wedge \text{UNCHANGED } \langle \text{globalVars}, \text{clientVars}, \text{replicaVars} \rangle
\end{aligned}$$

Replica 'r' handles replica 's' view change request 'm'

Replicas respond to *ViewChangeRequests* with the contents of their logs

for reconciliation. When a new view change is requested, the replica updates its *ViewID* and transitions to the *ViewChange* status to block writes during the transition.

$$\begin{aligned}
\text{HandleViewChangeRequest}(r, s, m) & \triangleq \\
& \wedge r\text{ViewID}[r] < m.\text{viewID} \\
& \wedge r\text{ViewID}' = [r\text{ViewID} \text{ EXCEPT } ![r] = m.\text{viewID}] \\
& \wedge r\text{Status}' = [r\text{Status} \text{ EXCEPT } ![r] = S\text{ViewChange}] \\
& \wedge r\text{ViewChanges}' = [r\text{ViewChanges} \text{ EXCEPT } ![r] = \{\}] \\
& \wedge \text{Reply}(m, [\text{src} \mapsto r, \\
& \quad \text{dest} \mapsto \text{Primary}(m.\text{viewID}), \\
& \quad \text{type} \mapsto M\text{ViewChangeReply}, \\
& \quad \text{viewID} \mapsto m.\text{viewID}, \\
& \quad \text{lastViewID} \mapsto r\text{LastViewID}[r], \\
& \quad \text{log} \mapsto r\text{Log}[r]]) \\
& \wedge \text{UNCHANGED } \langle \text{globalVars}, \text{clientVars}, r\text{Log}, r\text{Index}, r\text{SeqNum}, \\
& \quad r\text{Timestamp}, r\text{AbortPoint}, r\text{AbortReps}, r\text{LastViewID} \rangle
\end{aligned}$$

Replica 'r' handles replica 's' view change response 'm'

ViewChangeReplies are handled by the primary for the new view. Once

responses are received from a majority of the replicas including the new primary, the logs received from each replica are merged together to form the *log* for the new view. For each known session, the logs from each replica are merged by comparing each entry and keeping all non-empty sequential entries in the quorum. An updated timestamp is calculated from the reconciled *log*, and a *StartViewRequest* containing the new logs is sent to each replica.

```

HandleViewChangeReply(r, s, m)  $\triangleq$ 
   $\wedge$  IsPrimary(r)
   $\wedge$  rViewID[r] = m.viewID
   $\wedge$  rStatus[r] = SViewChange
   $\wedge$  rViewChanges' = [rViewChanges EXCEPT ![r] = rViewChanges[r]  $\cup$  {m}]
   $\wedge$  LET viewChanges  $\triangleq$  {v  $\in$  rViewChanges'[r] : v.viewID = rViewID[r]}
      viewSources  $\triangleq$  {v.src : v  $\in$  viewChanges}
      isQuorum  $\triangleq$  r  $\in$  viewSources  $\wedge$  viewSources  $\in$  Quorums
      lastViewIDs  $\triangleq$  {v.lastViewID : v  $\in$  viewChanges}
      lastViewID  $\triangleq$  (CHOOSE v1  $\in$  lastViewIDs :  $\forall$  v2  $\in$  lastViewIDs : v2  $\leq$  v1)
      lastViewChanges  $\triangleq$  {v2  $\in$  viewChanges : v2.lastViewID = lastViewID}
      viewLogs  $\triangleq$  [c  $\in$  Clients  $\mapsto$  {v1.log[c] : v1  $\in$  lastViewChanges}]
      mergeEnts(es)  $\triangleq$ 
        IF es = {}  $\vee \exists e \in es$  : e.type = TNoOp THEN
          CHOOSE e  $\in$  es : e.type = TNoOp
        ELSE
          CHOOSE e  $\in$  es : e.type  $\neq$  TNoOp  $\wedge$  e.index = Max({f.index : f  $\in$  es})
          range(ls)  $\triangleq$  Max({Len(l) : l  $\in$  ls})
          entries(ls, i)  $\triangleq$  {l[i] : l  $\in$  {k  $\in$  ls : i  $\leq$  Len(k)}}
          mergeLogs(ls)  $\triangleq$  [i  $\in$  1 .. range(ls)  $\mapsto$  mergeEnts(entries(ls, i))]
          viewLog  $\triangleq$  [c  $\in$  Clients  $\mapsto$  mergeLogs(viewLogs[c])]
          viewRange  $\triangleq$  Max({Len(viewLog[c]) : c  $\in$  Clients})
          liveEntries(l)  $\triangleq$  {l[i] : i  $\in$  {i  $\in$  DOMAIN l : l[i].type  $\neq$  TNoOp}}
          viewIndex  $\triangleq$  Sum({Cardinality(liveEntries(viewLog[c])) : c  $\in$  Clients})
          viewTimestamp  $\triangleq$  IF viewRange > 0 THEN
            Max(UNION { {l[i].timestamp : i  $\in$  DOMAIN l } :
                      l  $\in$  {viewLog[c] : c  $\in$  Clients} })
            ELSE 0
  IN
     $\vee \wedge$  isQuorum
       $\wedge$  Replies(m, {[src  $\mapsto$  r,
                     [dest  $\mapsto$  d,
                     [type  $\mapsto$  MStartViewRequest,
                     [viewID  $\mapsto$  rViewID[r],
                     [index  $\mapsto$  viewIndex,
                     [timestamp  $\mapsto$  viewTimestamp,
                     [log  $\mapsto$  viewLog] : d  $\in$  Replicas}])
     $\vee \wedge \neg$  isQuorum
       $\wedge$  Discard(m)

```

$$\wedge \text{UNCHANGED } \langle \text{globalVars}, \text{clientVars}, \text{rStatus}, \text{rViewID}, \text{rLog}, \text{rIndex}, \\ \text{rSeqNum}, \text{rTimestamp}, \text{rAbortPoint}, \text{rAbortReps}, \text{rLastViewID} \rangle$$

Replica 'r' handles replica 's' start view request 'm'

If the view is new, the replica updates its logs and session state from the request.

$$\begin{aligned} \text{HandleStartViewRequest}(r, s, m) &\triangleq \\ &\wedge \vee \text{rViewID}[r] < m.\text{viewID} \\ &\quad \vee \wedge \text{rViewID}[r] = m.\text{viewID} \\ &\quad \wedge \text{rStatus}[r] = \text{SViewChange} \\ &\wedge \text{rLog}' = [\text{rLog} \quad \text{EXCEPT } ![r] = m.\text{log}] \\ &\wedge \text{rSeqNum}' = [\text{rSeqNum} \quad \text{EXCEPT } ![r] = [c \in \text{Clients} \mapsto 0]] \\ &\wedge \text{rIndex}' = [\text{rIndex} \quad \text{EXCEPT } ![r] = m.\text{index}] \\ &\wedge \text{rTimestamp}' = [\text{rTimestamp} \quad \text{EXCEPT } ![r] = m.\text{timestamp}] \\ &\wedge \text{rStatus}' = [\text{rStatus} \quad \text{EXCEPT } ![r] = \text{SNormal}] \\ &\wedge \text{rViewID}' = [\text{rViewID} \quad \text{EXCEPT } ![r] = m.\text{viewID}] \\ &\wedge \text{rLastViewID}' = [\text{rLastViewID} \quad \text{EXCEPT } ![r] = m.\text{viewID}] \\ &\wedge \text{Discard}(m) \\ &\wedge \text{UNCHANGED } \langle \text{globalVars}, \text{clientVars}, \text{rAbortPoint}, \text{rAbortReps}, \text{rViewChanges} \rangle \end{aligned}$$

$$\begin{aligned} \text{InitMessageVars} &\triangleq \\ &\wedge \text{messages} = \{\} \\ &\wedge \text{messageCount} = 0 \end{aligned}$$

$$\begin{aligned} \text{InitClientVars} &\triangleq \\ &\wedge \text{cTime} = 0 \\ &\wedge \text{cViewID} = [c \in \text{Clients} \mapsto 1] \\ &\wedge \text{cSeqNum} = [c \in \text{Clients} \mapsto 0] \\ &\wedge \text{cReps} = [c \in \text{Clients} \mapsto \{\}] \\ &\wedge \text{cCommits} = [c \in \text{Clients} \mapsto \{\}] \end{aligned}$$

$$\begin{aligned} \text{InitReplicaVars} &\triangleq \\ &\wedge \text{replicas} = \text{SeqFromSet}(\text{Replicas}) \\ &\wedge \text{rStatus} = [r \in \text{Replicas} \mapsto \text{SNormal}] \\ &\wedge \text{rViewID} = [r \in \text{Replicas} \mapsto 1] \\ &\wedge \text{rLog} = [r \in \text{Replicas} \mapsto [c \in \text{Clients} \mapsto \langle \rangle]] \\ &\wedge \text{rIndex} = [r \in \text{Replicas} \mapsto 0] \\ &\wedge \text{rSeqNum} = [r \in \text{Replicas} \mapsto [c \in \text{Clients} \mapsto 0]] \\ &\wedge \text{rTimestamp} = [r \in \text{Replicas} \mapsto 0] \\ &\wedge \text{rAbortPoint} = [r \in \text{Replicas} \mapsto [\text{client} \mapsto \text{Nil}, \text{seqNum} \mapsto 0]] \\ &\wedge \text{rAbortReps} = [r \in \text{Replicas} \mapsto \{\}] \\ &\wedge \text{rLastViewID} = [r \in \text{Replicas} \mapsto 1] \\ &\wedge \text{rViewChanges} = [r \in \text{Replicas} \mapsto \{\}] \end{aligned}$$

$$\text{Init} \triangleq$$

$\wedge \text{InitMessageVars}$
 $\wedge \text{InitClientVars}$
 $\wedge \text{InitReplicaVars}$

This section specifies the invariants for the protocol.

Merge client logs together
 $\text{MergeLogs}(l) \triangleq$
 LET $\text{entries} \triangleq \text{UNION } \{\{l[i][j] : j \in \text{DOMAIN } l[i]\} : i \in \text{DOMAIN } l\}$
 IN $[i \in 1 \dots \text{Cardinality}(\text{entries}) \mapsto \text{CHOOSE } e \in \text{entries} : e.\text{index} = i]$

The type invariant asserts that the leader's *log* will never contain a different value at the same index as a client commit.

$\text{TypeOK} \triangleq$
 $\forall c \in \text{Clients} :$
 $\forall e \in c\text{Commits}[c] :$
 $\neg \exists r \in \text{Replicas} :$
 $\wedge \text{Primary}(r\text{ViewID}[r]) = r$
 $\wedge r\text{Status}[r] = \text{SNormal}$
 $\wedge \text{LET } \text{logs} \triangleq \text{MergeLogs}(r\text{Log}[r])$
 IN
 $\wedge \text{Len}(\text{logs}) \geq e.\text{index}$
 $\wedge \text{logs}[e.\text{index}].\text{value} \neq e.\text{value}$

$\text{NextClientRequest} \triangleq$
 $\exists c \in \text{Clients} :$
 $\exists v \in \text{Values} :$
 $\text{ClientRequest}(c, v)$

$\text{NextChangeView} \triangleq$
 $\exists r \in \text{Replicas} :$
 $\text{ChangeView}(r)$

$\text{NextHandleClientRequest} \triangleq$
 $\exists m \in \text{messages} :$
 $\wedge m.\text{type} = \text{MClientRequest}$
 $\wedge \text{HandleClientRequest}(m.\text{dest}, m.\text{src}, m)$

$\text{NextHandleClientReply} \triangleq$
 $\exists m \in \text{messages} :$
 $\wedge m.\text{type} = \text{MClientReply}$
 $\wedge \text{HandleClientReply}(m.\text{dest}, m.\text{src}, m)$

$\text{NextHandleRepairRequest} \triangleq$
 $\exists m \in \text{messages} :$

$$\begin{aligned}
& \wedge m.type = MRepairRequest \\
& \wedge HandleRepairRequest(m.dest, m.src, m) \\
NextHandleRepairReply & \triangleq \\
& \exists m \in messages : \\
& \quad \wedge m.type = MRepairReply \\
& \quad \wedge HandleRepairReply(m.dest, m.src, m) \\
NextHandleAbortRequest & \triangleq \\
& \exists m \in messages : \\
& \quad \wedge m.type = MAbortRequest \\
& \quad \wedge HandleAbortRequest(m.dest, m.src, m) \\
NextHandleAbortReply & \triangleq \\
& \exists m \in messages : \\
& \quad \wedge m.type = MAbortReply \\
& \quad \wedge HandleAbortReply(m.dest, m.src, m) \\
NextHandleViewChangeRequest & \triangleq \\
& \exists m \in messages : \\
& \quad \wedge m.type = MViewChangeRequest \\
& \quad \wedge HandleViewChangeRequest(m.dest, m.src, m) \\
NextHandleViewChangeReply & \triangleq \\
& \exists m \in messages : \\
& \quad \wedge m.type = MViewChangeReply \\
& \quad \wedge HandleViewChangeReply(m.dest, m.src, m) \\
NextHandleStartViewRequest & \triangleq \\
& \exists m \in messages : \\
& \quad \wedge m.type = MStartViewRequest \\
& \quad \wedge HandleStartViewRequest(m.dest, m.src, m) \\
NextDropMessage & \triangleq \\
& \exists m \in messages : \\
& \quad \wedge Discard(m) \\
& \quad \wedge UNCHANGED \langle globalVars, clientVars, replicaVars \rangle \\
Next & \triangleq \\
& \vee NextClientRequest \\
& \vee NextChangeView \\
& \vee NextHandleClientRequest \\
& \vee NextHandleClientReply \\
& \vee NextHandleRepairRequest \\
& \vee NextHandleRepairReply \\
& \vee NextHandleAbortRequest \\
& \vee NextHandleAbortReply
\end{aligned}$$

$\vee \textit{NextHandleViewChangeRequest}$
 $\vee \textit{NextHandleViewChangeReply}$
 $\vee \textit{NextHandleStartViewRequest}$
 $\vee \textit{NextDropMessage}$

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

\backslash * Modification History
 \backslash * Last modified *Mon Sep 28 17:46:34 PDT 2020* by *jordanhalterman*
 \backslash * Created *Fri Sep 18 22:45:21 PDT 2020* by *jordanhalterman*