

Defines the Just-In-Time *Paxos* (*JITPaxos*) protocol. *JITPaxos* is a variant of the *Paxos* consensus protocol designed for environments where process clocks are synchronized with high precision. The protocol relies on synchronized clocks to establish a global total ordering of events, avoiding coordination between replicas when requests arrive in the expected order, and reconciling requests only when they arrive out of order. This allows *JITPaxos* to reach consensus within a single round trip in the normal case, falling back to traditional replication strategies only when required.

*JITPaxos* uses a view-based approach to elect a primary and reconcile logs across views. Views are identified by a monotonically increasing, globally unique view *ID*. Each view deterministically assigns a quorum, and within the quorum a primary replica responsible for executing client requests and reconciling inconsistencies in the logs of the remaining replicas. *JITPaxos* replicas do not coordinate with each other in the normal case. Clients send timestamped requests in parallel to every replica in the view's quorum. When a replica receives a client request, if the request is received in chronological order, it's appended to the replica's *log*. If a request is received out of order (*i.e.* the request timestamp is less than the last timestamp in the replica's *log*), the request is rejected by the replica. Clients are responsible for identifying inconsistencies in the quorum's logs and initiating the reconciliation protocol. To help clients identify inconsistencies, replicas return a checksum representing the contents of the *log* up to the request point with each client reply. If a client's request is received out of chronological order, or if the checksums provided by the quorum do not match, the client must initiate the reconciliation protocol to reconcile the inconsistencies in the quorum's logs.

When requests are received out-of-order, the reconciliation protocol works to re-order requests using the view's primary as a reference. When a client initiates the reconciliation protocol for an inconsistent replica, the replica stops accepting client requests and sends a repair request to the primary. The primary responds with the subset of the *log* not yet reconciled on the replica, and the replica replaces the out-of-order entries in its *log*. Once the replica's *log* has been reconciled with the primary, it can acknowledge the reconciled request and begin accepting requests again. Once a client has reconciled all the divergent replicas and has received acknowledgement from each of the replicas in the quorum, the request can be committed.

View primaries and quorums are evenly distributed amongst view *IDs*. View changes can be initiated to change the primary or the set of replicas in the quorum. When a view change is initiated, each replica sends its *log* to the primary for the initiated view. Once the primary has received logs from a majority of replicas, it initializes the view with the *log* from the most recent in-sync replica, broadcasting the *log* to its peers. The use of quorums to determine both the commitment of a request and the initialization of new views ensures that each view *log* contains all prior committed requests.

EXTENDS *Naturals*, *Reals*, *Sequences*, *FiniteSets*, *TLC*

The set of *JITPaxos* replicas

CONSTANT *Replicas*

The set of *JITPaxos* clients

CONSTANT *Clients*

The set of possible values

CONSTANT *Values*

An empty value

CONSTANT *Nil*

Request/response types

CONSTANTS

*MClientRequest*,  
*MClientReply*,  
*MReconcileRequest*,  
*MReconcileReply*,  
*MRepairRequest*,  
*MRepairReply*,  
*MViewChangeRequest*,  
*MViewChangeReply*,  
*MStartViewRequest*

Replica statuses

CONSTANTS

*SInSync*,  
*SRepair*,  
*SViewChange*

This section specifies the message types and schemas used in this spec.

$ReqIDs \triangleq [c \in Clients \mapsto i \in (1 \dots)]$

$ViewIDs \triangleq [r \in Replicas \mapsto i \in (1 \dots)]$

$Logs \triangleq [r \in Replicas \mapsto [i \in (1 \dots) \mapsto Value]]$

$Indexes \triangleq [r \in Replicas \mapsto i \in (1 \dots)]$

$Timestamps \triangleq [r \in Replicas \mapsto i \in (1 \dots)]$

$Checksums \triangleq [r \in Replicas \mapsto [i \in (1 \dots) \mapsto t \in Timestamps]]$

*ClientRequest*

[ *src*  $\mapsto c \in Clients$ ,  
*dest*  $\mapsto r \in Replicas$ ,  
*type*  $\mapsto MClientRequest$ ,  
*viewID*  $\mapsto i \in ViewIDs$ ,  
*reqID*  $\mapsto i \in ReqIDs$ ,  
*value*  $\mapsto v \in Values$ ,  
*timestamp*  $\mapsto t \in Timestamps$  ]

*ClientReply*

[ *src*  $\mapsto r \in Replicas$ ,  
*dest*  $\mapsto c \in Clients$ ,  
*req*  $\mapsto (ClientRequest)$ ,  
*type*  $\mapsto MClientReply$ ,  
*viewID*  $\mapsto i \in ViewIDs$ ,  
*index*  $\mapsto i \in Indexes$ ,  
*checksum*  $\mapsto c \in Checksums$ ,  
*value*  $\mapsto v \in Values$ ,  
*timestamp*  $\mapsto t \in Timestamps$ ,  
*succeeded*  $\mapsto TRUE \vee FALSE$  ]

#### *ReconcileRequest*

[ *src*  $\mapsto c \in \text{Clients}$ ,  
  *dest*  $\mapsto r \in \text{Replicas}$ ,  
  *type*  $\mapsto \text{MReconcileRequest}$ ,  
  *viewID*  $\mapsto i \in \text{ViewIDs}$ ,  
  *reqID*  $\mapsto i \in \text{ReqIDs}$ ,  
  *index*  $\mapsto i \in \text{Indexes}$  ]

#### *ReconcileReply*

[ *src*  $\mapsto r \in \text{Replicas}$ ,  
  *dest*  $\mapsto c \in \text{Clients}$ ,  
  *req*  $\mapsto (\text{ClientRequest})$ ,  
  *type*  $\mapsto \text{MReconcileReply}$ ,  
  *viewID*  $\mapsto i \in \text{ViewIDs}$ ,  
  *index*  $\mapsto i \in \text{Indexes}$ ,  
  *checksum*  $\mapsto c \in \text{Checksums}$ ,  
  *value*  $\mapsto v \in \text{Values}$ ,  
  *timestamp*  $\mapsto t \in \text{Timestamps}$ ,  
  *succeeded*  $\mapsto \text{TRUE} \vee \text{FALSE}$  ]

#### *RepairRequest*

[ *src*  $\mapsto r \in \text{Replicas}$ ,  
  *dest*  $\mapsto r \in \text{Replicas}$ ,  
  *req*  $\mapsto (\text{ClientRequest})$ ,  
  *type*  $\mapsto \text{MRepairRequest}$ ,  
  *viewID*  $\mapsto i \in \text{ViewIDs}$ ,  
  *index*  $\mapsto i \in \text{Indexes}$  ]

#### *RepairReply*

[ *src*  $\mapsto r \in \text{Replicas}$ ,  
  *dest*  $\mapsto r \in \text{Replicas}$ ,  
  *req*  $\mapsto (\text{ClientRequest})$ ,  
  *type*  $\mapsto \text{MRepairReply}$ ,  
  *viewID*  $\mapsto i \in \text{ViewIDs}$ ,  
  *index*  $\mapsto i \in \text{Indexes}$ ,  
  *log*  $\mapsto l \in \text{Logs}$  ]

#### *ViewChangeRequest*

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

#### *ViewChangeReply*

[ *src*  $\mapsto r \in \text{Replicas}$ ,  
  *dest*  $\mapsto r \in \text{Replicas}$ ,  
  *type*  $\mapsto \text{MViewChangeReply}$ ,  
  *viewID*  $\mapsto i \in \text{ViewIDs}$ ,  
  *logViewID*  $\mapsto i \in \text{ViewIDs}$ ,  
  *log*  $\mapsto l \in \text{Logs}$  ]

#### *StartViewRequest*

[ *src*  $\mapsto r \in \text{Replicas}$ ,

$dest \mapsto r \in Replicas,$   
 $type \mapsto MStartViewRequest,$   
 $viewID \mapsto i \in ViewIDs,$   
 $log \mapsto l \in Logs ]$

---

The set of all messages on the network  
 VARIABLE  $messages$

The total number of messages sent  
 VARIABLE  $messageCount$

The total number of steps executed  
 VARIABLE  $stepCount$

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

Local client state

Strictly increasing representation of synchronized time  
 VARIABLE  $cTime$

The highest known view  $ID$  for a client  
 VARIABLE  $cViewID$

Client request  $IDs$   
 VARIABLE  $cReqID$

A client response buffer  
 VARIABLE  $cReps$

A set of all commits - used for model checking  
 VARIABLE  $cCommits$

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

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  $rLog$

A replica's sync index  
 VARIABLE  $rSyncIndex$

The view  $ID$  for the  $log$   
VARIABLE  $rLogViewID$

The set of view change replies  
VARIABLE  $rViewChangeReps$

$replicaVars \triangleq \langle rStatus, rViewID, rLog, rSyncIndex, rLogViewID, rViewChangeReps \rangle$   
 $vars \triangleq \langle messageVars, clientVars, replicaVars \rangle$

This section provides utilities for implementing 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  $\langle x \rangle \circ SeqFromSet(S \setminus \{x\})$

RECURSIVE  $SetReduce(-, -, -)$   
 $SetReduce(Op(-, -), S, value) \triangleq$   
IF  $S = \{\}$  THEN  
 $value$   
ELSE  
LET  $s \triangleq \text{CHOOSE } s \in S : \text{TRUE}$   
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  $SetReduce(_op, S, 0)$

The values of a sequence  
 $Range(s) \triangleq \{s[i] : i \in \text{DOMAIN } s\}$

This section provides helpers for the protocol.

A sorted sequence of replicas  
 $replicas \triangleq SeqFromSet(Replicas)$

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

The primary for view 'v'

$Primary(v) \triangleq replicas[PrimaryIndex(v)]$

Quorum is the quorum for a given view

$Quorum(v) \triangleq$

LET

$quorumSize \triangleq Len(replicas) \div 2 + 1$

$index(i) \triangleq PrimaryIndex(v) + (i - 1)$

$member(i) \triangleq \text{IF } index(i) > Len(replicas) \text{ THEN } replicas[index(i) \% Len(replicas)] \text{ ELSE } replicas[index(i)]$

IN

$\{member(i) : i \in 1 \dots quorumSize\}$

A boolean indicating whether the given set is a quorum

$IsQuorum(S) \triangleq Cardinality(S) * 2 \geq Cardinality(Replicas)$

A boolean indicating whether the given set is a quorum that includes the given replica

$IsLocalQuorum(r, S) \triangleq IsQuorum(S) \wedge r \in S$

This section models the network.

Send a set of messages

$Sends(ms) \triangleq$

$\wedge messages' = messages \cup ms$

$\wedge messageCount' = messageCount + Cardinality(ms)$

$\wedge stepCount' = stepCount + 1$

Send a message

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

Ack a message

$Ack(m) \triangleq$

$\wedge messages' = messages \setminus \{m\}$

$\wedge messageCount' = messageCount + 1$

$\wedge stepCount' = stepCount + 1$

Ack a message and send a set of messages

$AckAndSends(m, ms) \triangleq$

$\wedge messages' = (messages \cup ms) \setminus \{m\}$

$\wedge messageCount' = messageCount + Cardinality(ms)$

$\wedge stepCount' = stepCount + 1$

Ack and send a message

$AckAndSend(m, n) \triangleq AckAndSends(m, \{n\})$

Reply to a message with a set of responses

$Replies(req, reps) \triangleq AckAndSends(req, reps)$

Reply to a message

$$Reply(req, resp) \triangleq AckAndSend(req, resp)$$

---

This section models client requests.

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

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

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

$$\begin{aligned}
HandleClientReply(c, r, m) &\triangleq \\
&\text{If the reply view } ID \text{ does not match the request view } ID, \text{ update the client's view.} \\
&\wedge \vee \wedge m.viewID \neq m.req.viewID \\
&\quad \wedge \vee \wedge cViewID[c] < m.viewID \\
&\quad \quad \wedge cViewID' = [cViewID \text{ EXCEPT } ![c] = m.viewID] \\
&\quad \vee \wedge cViewID[c] \geq m.viewID \\
&\quad \quad \wedge UNCHANGED \langle cViewID \rangle \\
&\quad \wedge Ack(m) \\
&\quad \wedge UNCHANGED \langle cReps, cCommits \rangle \\
&\text{If the request and reply views match and the reply view matches the client's view,} \\
&\text{aggregate the replies for the associated client request.} \\
&\vee \wedge m.viewID = m.req.viewID \\
&\quad \wedge m.viewID = cViewID[c] \\
&\quad \wedge \vee \wedge m.succeeded \\
&\quad \quad \wedge cReps' = [cReps \text{ EXCEPT } ![c] = \\
&\quad \quad \quad (cReps[c] \setminus \{n \in cReps[c] : \wedge n.src = m.src \\
&\quad \quad \quad \quad \wedge n.viewID = cViewID[c] \\
&\quad \quad \quad \quad \wedge n.req.reqID = m.req.reqID \\
&\quad \quad \quad \quad \wedge \neg n.succeeded\}) \cup \{m\}] \\
&\quad \vee \wedge \neg m.succeeded \\
&\quad \quad \wedge \neg \exists n \in cReps[c] : \wedge n.src = m.src \\
&\quad \quad \quad \wedge n.viewID = cViewID[c] \\
&\quad \quad \quad \wedge n.req.reqID = m.req.reqID \\
&\quad \quad \quad \wedge n.succeeded \\
&\quad \quad \wedge cReps' = [cReps \text{ EXCEPT } ![c] = cReps[c] \cup \{m\}] \\
&\quad \wedge LET reps \triangleq \{n \in cReps'[c] : \wedge n.viewID = cViewID[c] \\
&\quad \quad \quad \wedge n.req.reqID = m.req.reqID\}
\end{aligned}$$

$$\begin{aligned}
isQuorum &\triangleq \{n.src : n \in \{n \in reps : n.succeeded\}\} = Quorum(cViewID[c]) \\
isCommitted &\triangleq \wedge \forall n \in reps : n.succeeded \\
&\quad \wedge Cardinality(\{n.checksum : n \in reps\}) = 1 \\
hasPrimary &\triangleq \exists n \in reps : n.src = Primary(cViewID[c]) \wedge n.succeeded
\end{aligned}$$

IN

If a quorum of successful replies have been received and the checksums match, add the primary reply to commits.

$$\begin{aligned}
&\vee \wedge isQuorum \\
&\quad \wedge isCommitted \\
&\quad \wedge LET commit \triangleq CHOOSE n \in reps : n.src = Primary(cViewID[c]) \\
&\quad \quad IN cCommits' = [cCommits EXCEPT ![c] = cCommits[c] \cup \{commit\}] \\
&\quad \wedge Ack(m)
\end{aligned}$$

If some reply failed or was returned with an incorrect checksum, send a *ReconcileRequest* to the inconsistent node to force it to reconcile its *log* with the primary's *log*.

$$\begin{aligned}
&\vee \wedge \neg isCommitted \\
&\quad \wedge \vee \wedge hasPrimary \\
&\quad \quad \wedge LET primaryRep \triangleq CHOOSE n \in reps : \wedge n.src = Primary(cViewID[c]) \\
&\quad \quad \quad \wedge n.succeeded \\
&\quad \quad \quad retryReps \triangleq \{n \in reps : \\
&\quad \quad \quad \quad \wedge n.src \neq Primary(cViewID[c]) \\
&\quad \quad \quad \quad \wedge n.checksum \neq primaryRep.checksum\} \\
&\quad \quad \quad IN AckAndSends(m, \{[src \mapsto c, \\
&\quad \quad \quad \quad dest \mapsto r, \\
&\quad \quad \quad \quad type \mapsto MReconcileRequest, \\
&\quad \quad \quad \quad viewID \mapsto cViewID[c], \\
&\quad \quad \quad \quad reqID \mapsto m.req.reqID, \\
&\quad \quad \quad \quad index \mapsto primaryRep.index] : n \in retryReps\})
\end{aligned}$$

$$\begin{aligned}
&\quad \vee \wedge \neg hasPrimary \\
&\quad \quad \wedge Ack(m) \\
&\quad \quad \wedge UNCHANGED \langle cCommits \rangle
\end{aligned}$$

If a quorum has not yet been reached, wait for more replies.

$$\begin{aligned}
&\vee \wedge \neg isQuorum \\
&\quad \wedge isCommitted \\
&\quad \wedge Ack(m) \\
&\quad \wedge UNCHANGED \langle cCommits \rangle \\
&\quad \wedge UNCHANGED \langle cViewID \rangle \\
&\quad \wedge UNCHANGED \langle replicaVars, cTime, cReqID \rangle
\end{aligned}$$

$$HandleReconcileReply(c, r, m) \triangleq HandleClientReply(c, r, m)$$


---

This section models the replica protocol.



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

$$HandleClientRequest(r, c, m) \triangleq$$

Client requests can only be handled if the replica is in-sync.

$$\wedge rStatus[r] = SInSync$$

If the client's view matches the replica's view, process the client's request.

$$\wedge \vee \wedge m.viewID = rViewID[r]$$

$$\wedge LET lastTimestamp \triangleq Max(\{rLog[r][i].timestamp : i \in DOMAIN\ rLog[r]\} \cup \{0\})$$

IN

If the request timestamp is greater than the highest *log* timestamp, append the entry to the *log* and return a successful response with the appended entry index.

$$\wedge \vee \wedge m.timestamp > lastTimestamp$$

$$\wedge rLog' = [rLog\ EXCEPT\ ![r] =$$

$$Append(rLog[r], [value \mapsto m.value,$$

$$timestamp \mapsto m.timestamp])]$$

$$\wedge Reply(m, [src \mapsto r,$$

$$dest \mapsto c,$$

$$req \mapsto m,$$

$$type \mapsto MClientReply,$$

$$viewID \mapsto rViewID[r],$$

$$index \mapsto Len(rLog'[r]),$$

$$checksum \mapsto rLog'[r],$$

$$value \mapsto m.value,$$

$$timestamp \mapsto m.timestamp,$$

$$succeeded \mapsto TRUE])$$

If the request timestamp matches the highest *log* timestamp, treat the request as a duplicate. Return a successful response indicating the entry was appended.

$$\vee \wedge m.timestamp = lastTimestamp$$

$$\wedge Reply(m, [src \mapsto r,$$

$$dest \mapsto c,$$

$$req \mapsto m,$$

$$type \mapsto MClientReply,$$

$$viewID \mapsto rViewID[r],$$

$$index \mapsto Len(rLog[r]),$$

$$checksum \mapsto rLog[r],$$

$$value \mapsto m.value,$$

$$timestamp \mapsto m.timestamp,$$

$$succeeded \mapsto TRUE])$$

$$\wedge UNCHANGED\ \langle rLog \rangle$$

If the request timestamp is less than the highest *log* timestamp, reject the request.

$$\vee \wedge m.timestamp < lastTimestamp$$

$$\wedge Reply(m, [src \mapsto r,$$

$$dest \mapsto c,$$

$$\begin{aligned}
& \begin{array}{lcl}
req & \mapsto & m, \\
type & \mapsto & MClientReply, \\
viewID & \mapsto & rViewID[r], \\
index & \mapsto & Len(rLog[r]), \\
checksum & \mapsto & rLog[r], \\
value & \mapsto & m.value, \\
timestamp & \mapsto & m.timestamp, \\
succeeded & \mapsto & FALSE]) \\
& \wedge \text{UNCHANGED } \langle rLog \rangle \\
& \wedge \text{UNCHANGED } \langle rViewID, rStatus, rViewChangeReps \rangle \\
& \text{If the client's view is greater than the replica's view, reject the client's} \\
& \text{request with the outdated view } ID \text{ and enter the view change protocol.} \\
& \vee \wedge m.viewID > rViewID[r] \\
& \wedge rViewID' = [rViewID \quad \text{EXCEPT } ![r] = m.viewID] \\
& \wedge rStatus' = [rStatus \quad \text{EXCEPT } ![r] = SViewChange] \\
& \wedge rViewChangeReps' = [rViewChangeReps \quad \text{EXCEPT } ![r] = \{\}] \\
& \wedge Replies(m, \{[src \mapsto r, \\
& \quad \quad \quad dest \mapsto c, \\
& \quad \quad \quad req \mapsto m, \\
& \quad \quad \quad type \mapsto MClientReply, \\
& \quad \quad \quad viewID \mapsto rViewID[r], \\
& \quad \quad \quad succeeded \mapsto FALSE], \\
& \quad [src \mapsto r, \\
& \quad \quad \quad dest \mapsto Primary(m.viewID), \\
& \quad \quad \quad type \mapsto MViewChangeReply, \\
& \quad \quad \quad viewID \mapsto m.viewID, \\
& \quad \quad \quad logViewID \mapsto rLogViewID[r], \\
& \quad \quad \quad log \mapsto rLog[r]]\}) \\
& \wedge \text{UNCHANGED } \langle rLog \rangle \\
& \text{If the client's view is less than the replica's view, reject the client's request} \\
& \text{with the updated view } ID \text{ to force the client to retry.} \\
& \vee \wedge m.viewID < rViewID[r] \\
& \wedge Reply(m, [src \mapsto r, \\
& \quad \quad \quad dest \mapsto c, \\
& \quad \quad \quad req \mapsto m, \\
& \quad \quad \quad type \mapsto MClientReply, \\
& \quad \quad \quad viewID \mapsto rViewID[r], \\
& \quad \quad \quad succeeded \mapsto FALSE]) \\
& \wedge \text{UNCHANGED } \langle rViewID, rStatus, rLog, rViewChangeReps \rangle \\
& \wedge \text{UNCHANGED } \langle clientVars, rLogViewID, rSyncIndex \rangle \\
& HandleReconcileRequest(r, c, m) \triangleq \\
& \wedge rStatus[r] = SInSync \\
& \wedge rViewID[r] = m.viewID \\
& \wedge \vee \wedge rSyncIndex[r] \geq m.index
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{Reply}(m, [src \mapsto r, \\
& \quad dest \mapsto c, \\
& \quad req \mapsto m, \\
& \quad type \mapsto MReconcileReply, \\
& \quad viewID \mapsto rViewID[r], \\
& \quad index \mapsto m.index, \\
& \quad checksum \mapsto [i \in 1 \dots m.index \mapsto rLog[r][i]], \\
& \quad value \mapsto rLog[r][m.index].value, \\
& \quad timestamp \mapsto rLog[r][m.index].timestamp, \\
& \quad succeeded \mapsto \text{TRUE}]) \\
& \wedge \text{UNCHANGED } \langle rStatus \rangle \\
\vee & \wedge rSyncIndex[r] < m.index \\
& \wedge \text{Primary}(rViewID[r]) \neq r \\
& \wedge rStatus' = [rStatus \text{ EXCEPT } ![r] = SRepair] \\
& \wedge \text{AckAndSend}(m, [src \mapsto r, \\
& \quad dest \mapsto \text{Primary}(rViewID[r]), \\
& \quad req \mapsto m, \\
& \quad type \mapsto MRepairRequest, \\
& \quad viewID \mapsto rViewID[r], \\
& \quad index \mapsto m.index]) \\
& \wedge \text{UNCHANGED } \langle clientVars, rViewID, rLog, rLogViewID, rSyncIndex, rViewChangeReps \rangle \\
\\
\text{HandleRepairRequest}(r, s, m) & \triangleq \\
& \wedge rStatus[r] = SInSync \\
& \wedge rViewID[r] = m.viewID \\
& \wedge \text{Primary}(rViewID[r]) = r \\
& \wedge \text{Reply}(m, [src \mapsto r, \\
& \quad dest \mapsto s, \\
& \quad req \mapsto m.req, \\
& \quad type \mapsto MRepairReply, \\
& \quad viewID \mapsto rViewID[r], \\
& \quad index \mapsto m.index, \\
& \quad log \mapsto [i \in 1 \dots m.index \mapsto rLog[r][i]]]) \\
& \wedge \text{UNCHANGED } \langle clientVars, replicaVars \rangle \\
\\
\text{HandleRepairReply}(r, s, m) & \triangleq \\
& \wedge rStatus[r] = SRepair \\
& \wedge rViewID[r] = m.viewID \\
& \wedge rStatus' = [rStatus \text{ EXCEPT } ![r] = SInSync] \\
& \wedge rLog' = [rLog \text{ EXCEPT } ![r] = m.log \circ SubSeq(rLog[r], Len(m.log), Len(rLog[r]))] \\
& \wedge rSyncIndex' = [rSyncIndex \text{ EXCEPT } ![r] = Len(rLog'[r])] \\
& \wedge \text{Reply}(m, [src \mapsto r, \\
& \quad dest \mapsto m.req.src, \\
& \quad req \mapsto m.req, \\
& \quad type \mapsto MReconcileReply,
\end{aligned}$$

$$\begin{aligned}
viewID &\mapsto rViewID[r], \\
index &\mapsto m.index, \\
checksum &\mapsto m.log, \\
value &\mapsto m.log[m.index].value, \\
timestamp &\mapsto m.log[m.index].timestamp, \\
succeeded &\mapsto \text{TRUE}) \\
&\wedge \text{UNCHANGED } \langle clientVars, rViewID, rLogViewID, rViewChangeReps \rangle
\end{aligned}$$

Replica 'r' requests a view change

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

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

$$\begin{aligned}
HandleViewChangeRequest(r, s, m) &\triangleq \\
&\wedge \vee \wedge rViewID[r] < m.viewID \\
&\quad \wedge rViewID' = [rViewID \text{ EXCEPT } ![r] = m.viewID] \\
&\quad \wedge rStatus' = [rStatus \text{ EXCEPT } ![r] = SViewChange] \\
&\quad \wedge rViewChangeReps' = [rViewChangeReps \text{ EXCEPT } ![r] = \{\}] \\
&\quad \wedge Reply(m, [src \mapsto r, \\
&\quad \quad dest \mapsto Primary(m.viewID), \\
&\quad \quad type \mapsto MViewChangeReply, \\
&\quad \quad viewID \mapsto m.viewID, \\
&\quad \quad logViewID \mapsto rLogViewID[r], \\
&\quad \quad log \mapsto rLog[r]]) \\
&\vee \wedge rViewID[r] \geq m.viewID \\
&\quad \wedge Ack(m) \\
&\quad \wedge \text{UNCHANGED } \langle rViewID, rStatus, rViewChangeReps \rangle \\
&\wedge \text{UNCHANGED } \langle clientVars, rLog, rLogViewID, rSyncIndex \rangle
\end{aligned}$$

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

$$\begin{aligned}
HandleViewChangeReply(r, s, m) &\triangleq \\
&\text{The view change protocol is run by the primary for the view.} \\
&\wedge Primary(m.viewID) = r \\
&\wedge rViewID[r] = m.viewID \\
&\wedge rStatus[r] = SViewChange \\
&\wedge rViewChangeReps' = [rViewChangeReps \text{ EXCEPT } ![r] = rViewChangeReps[r] \cup \{m\}] \\
&\wedge \text{LET } viewChanges \triangleq \{v \in rViewChangeReps'[r] : v.viewID = rViewID[r]\} \\
&\text{IN}
\end{aligned}$$

In order to ensure the new view is initialized with the latest view,  
a quorum of view change replies must be received to guarantee the last  
activated view is present in the set of replies.  
If view change replies have been received from a majority of the replicas,

initialize the view using the *log* from the highest activated view.

$$\begin{aligned}
& \vee \wedge \text{IsLocalQuorum}(r, \{v.\text{src} : v \in \text{viewChanges}\}) \\
& \wedge \text{LET } \text{latestViewID} \triangleq \text{Max}(\{v.\text{logViewID} : v \in \text{viewChanges}\}) \\
& \quad \text{latestChange} \triangleq \text{CHOOSE } v \in \text{viewChanges} : \\
& \quad \quad \wedge v.\text{logViewID} = \text{latestViewID} \\
& \quad \quad \wedge v.\text{src} \in \text{Quorum}(\text{latestViewID}) \\
& \text{IN } \text{AckAndSends}(m, \{[\text{src} \mapsto r, \\
& \quad \quad \text{dest} \mapsto d, \\
& \quad \quad \text{type} \mapsto \text{MStartViewRequest}, \\
& \quad \quad \text{viewID} \mapsto r\text{ViewID}[r], \\
& \quad \quad \text{log} \mapsto \text{latestChange.log}] : d \in \text{Replicas}\}) \\
& \text{If view change replies have not yet been received from a quorum, record} \\
& \text{the view change reply and discard the message.} \\
& \vee \wedge \neg \text{IsLocalQuorum}(r, \{v.\text{src} : v \in \text{viewChanges}\}) \\
& \quad \wedge \text{Ack}(m) \\
& \wedge \text{UNCHANGED } \langle \text{clientVars}, r\text{Status}, r\text{ViewID}, r\text{Log}, r\text{LogViewID}, r\text{SyncIndex} \rangle
\end{aligned}$$

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

$$\begin{aligned}
& \text{HandleStartViewRequest}(r, s, m) \triangleq \\
& \quad \text{To activate a view, the replica must either not know of the view or already} \\
& \quad \text{be participating in the view change protocol for the view.} \\
& \wedge \vee r\text{ViewID}[r] < m.\text{viewID} \\
& \quad \vee \wedge r\text{ViewID}[r] = m.\text{viewID} \\
& \quad \wedge r\text{Status}[r] = S\text{ViewChange} \\
& \quad \text{If the replica is part of the quorum for the activated view, update the } \text{log} \\
& \quad \text{and record the activated view for use in the view change protocol.} \\
& \wedge \vee \wedge r \in \text{Quorum}(m.\text{viewID}) \\
& \quad \wedge r\text{Log}' = [r\text{Log} \text{ EXCEPT } ![r] = m.\text{log}] \\
& \quad \wedge r\text{LogViewID}' = [r\text{LogViewID} \text{ EXCEPT } ![r] = m.\text{viewID}] \\
& \quad \wedge r\text{SyncIndex}' = [r\text{SyncIndex} \text{ EXCEPT } ![r] = \text{Len}(m.\text{log})] \\
& \quad \vee \wedge r \notin \text{Quorum}(m.\text{viewID}) \\
& \quad \wedge \text{UNCHANGED } \langle r\text{Log}, r\text{LogViewID}, r\text{SyncIndex} \rangle \\
& \quad \text{Update the replica's view ID and status and clean up view change state.} \\
& \wedge r\text{ViewID}' = [r\text{ViewID} \text{ EXCEPT } ![r] = m.\text{viewID}] \\
& \wedge r\text{Status}' = [r\text{Status} \text{ EXCEPT } ![r] = S\text{InSync}] \\
& \wedge \text{LET } \text{viewChanges} \triangleq \{v \in r\text{ViewChangeReps}[r] : v.\text{viewID} = r\text{ViewID}[r]\} \\
& \quad \text{IN } r\text{ViewChangeReps}' = [r\text{ViewChangeReps} \text{ EXCEPT } ![r] = r\text{ViewChangeReps}[r] \setminus \text{viewChanges}] \\
& \wedge \text{Ack}(m) \\
& \wedge \text{UNCHANGED } \langle \text{clientVars} \rangle
\end{aligned}$$


---


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

$$\begin{aligned}
InitClientVars &\triangleq \\
&\wedge cTime = 0 \\
&\wedge cViewID = [c \in Clients \mapsto 1] \\
&\wedge cReqID = [c \in Clients \mapsto 0] \\
&\wedge cReps = [c \in Clients \mapsto \{\}] \\
&\wedge cCommits = [c \in Clients \mapsto \{\}] \\
InitReplicaVars &\triangleq \\
&\wedge rStatus = [r \in Replicas \mapsto SInSync] \\
&\wedge rViewID = [r \in Replicas \mapsto 1] \\
&\wedge rLog = [r \in Replicas \mapsto \langle \rangle] \\
&\wedge rSyncIndex = [r \in Replicas \mapsto 0] \\
&\wedge rLogViewID = [r \in Replicas \mapsto 1] \\
&\wedge rViewChangeReps = [r \in Replicas \mapsto \{\}] \\
Init &\triangleq \\
&\wedge InitMessageVars \\
&\wedge InitClientVars \\
&\wedge InitReplicaVars
\end{aligned}$$

---

This section specifies the invariants for the protocol.

The linearizability invariant verifies that once a client has received matching acks from a quorum and committed a value, thereafter the value is always present at the committed index on all in-sync replicas.

$$\begin{aligned}
Linearizability &\triangleq \\
&\forall c \in Clients : \\
&\quad \forall e \in cCommits[c] : \\
&\quad \quad \neg \exists r \in Replicas : \\
&\quad \quad \quad \wedge rStatus[r] = SInSync \\
&\quad \quad \quad \wedge rViewID[r] \geq e.viewID \\
&\quad \quad \quad \wedge r \in Quorum(rViewID[r]) \\
&\quad \quad \quad \wedge rLog[r][e.index].value \neq e.value
\end{aligned}$$

---


$$\begin{aligned}
NextClientRequest &\triangleq \\
&\exists c \in Clients : \\
&\quad \exists v \in Values : \\
&\quad \quad ClientRequest(c, v)
\end{aligned}$$

$$\begin{aligned}
NextChangeView &\triangleq \\
&\exists r \in Replicas : \\
&\quad ChangeView(r)
\end{aligned}$$

$$NextHandleClientRequest \triangleq$$

$$\begin{aligned}
& \exists m \in \text{messages} : \\
& \quad \wedge m.type = MClientRequest \\
& \quad \wedge HandleClientRequest(m.dest, m.src, m) \\
NextHandleClientReply & \triangleq \\
& \exists m \in \text{messages} : \\
& \quad \wedge m.type = MClientReply \\
& \quad \wedge HandleClientReply(m.dest, m.src, m) \\
NextHandleReconcileRequest & \triangleq \\
& \exists m \in \text{messages} : \\
& \quad \wedge m.type = MReconcileRequest \\
& \quad \wedge HandleReconcileRequest(m.dest, m.src, m) \\
NextHandleReconcileReply & \triangleq \\
& \exists m \in \text{messages} : \\
& \quad \wedge m.type = MReconcileReply \\
& \quad \wedge HandleReconcileReply(m.dest, m.src, m) \\
NextHandleRepairRequest & \triangleq \\
& \exists m \in \text{messages} : \\
& \quad \wedge m.type = MRepairRequest \\
& \quad \wedge HandleRepairRequest(m.dest, m.src, m) \\
NextHandleRepairReply & \triangleq \\
& \exists m \in \text{messages} : \\
& \quad \wedge m.type = MRepairReply \\
& \quad \wedge HandleRepairReply(m.dest, m.src, m) \\
NextHandleViewChangeRequest & \triangleq \\
& \exists m \in \text{messages} : \\
& \quad \wedge m.type = MViewChangeRequest \\
& \quad \wedge HandleViewChangeRequest(m.dest, m.src, m) \\
NextHandleViewChangeReply & \triangleq \\
& \exists m \in \text{messages} : \\
& \quad \wedge m.type = MViewChangeReply \\
& \quad \wedge HandleViewChangeReply(m.dest, m.src, m) \\
NextHandleStartViewRequest & \triangleq \\
& \exists m \in \text{messages} : \\
& \quad \wedge m.type = MStartViewRequest \\
& \quad \wedge HandleStartViewRequest(m.dest, m.src, m) \\
NextDropMessage & \triangleq \\
& \exists m \in \text{messages} : \\
& \quad \wedge Ack(m) \\
& \quad \wedge UNCHANGED \langle clientVars, replicaVars \rangle
\end{aligned}$$

$$Next \triangleq$$

- $\vee NextClientRequest$
- $\vee NextChangeView$
- $\vee NextHandleClientRequest$
- $\vee NextHandleClientReply$
- $\vee NextHandleReconcileRequest$
- $\vee NextHandleReconcileReply$
- $\vee NextHandleRepairRequest$
- $\vee NextHandleRepairReply$
- $\vee NextHandleViewChangeRequest$
- $\vee NextHandleViewChangeReply$
- $\vee NextHandleStartViewRequest$
- $\vee NextDropMessage$

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


---

\ \* Modification History  
 \ \* Last modified *Wed Sep 30 18:02:32 PDT 2020* by *jordanhalterman*  
 \ \* Created *Fri Sep 18 22:45:21 PDT 2020* by *jordanhalterman*