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.

Summary:

* View-based protocol

* Views are identified by a monotonically increasing, globally unique identifier

* Each view assigns a primary plus a set of replicas that form the quorum

* Clients timestamp each request and send in parallel to all replicas in the quorum

* Each replica appends requests to a log in chronological order, and the primary executes requests

* If a request is received out of chronological order it is rejected

* Replies to clients include a checksum of the log on each replica

* If the client receives a reply indicating the request was received out of chronological order or if a checksum does not match the primary's checksum, the client initiates a reconciliation protocol

* To reconcile inconsistencies in the log, replicas pull logs from the primary, and once logs have been reconciled the original request is acknowledged

* Once the client receives matching acknowledgements from all the replicas in the quorum a request is committed

* View changes select the most recent log from a majority of the replicas to ensure the initial view log contains all committed requests from prior views


*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 to 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 reconcilitation 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.

$ReqIDs \triangleq [c \in Clients \mapsto i \in (1 \mathinner{\ldotp\ldotp})]$

$ViewIDs \triangleq [r \in Replicas \mapsto i \in (1 \mathinner{\ldotp\ldotp})]$

$Logs \triangleq [r \in Replicas \mapsto [i \in (1 \mathinner{\ldotp\ldotp}) \mapsto Value]]$

$Indexes \triangleq [r \in Replicas \mapsto i \in (1 \mathinner{\ldotp\ldotp})]$

$Timestamps \triangleq [r \in Replicas \mapsto i \in (1 \mathinner{\ldotp\ldotp})]$

$Checksums \triangleq [r \in Replicas \mapsto [i \in (1..) \mapsto t \in \text{Timestamps}]]$

$ClientRequest$
$[\ src \qquad \mapsto c \in Clients,$
$\quad dest \qquad \mapsto r \in Replicas,$
$\quad type \qquad \mapsto MClientRequest,$
$\quad viewID \mapsto i \in ViewIDs,$
$\quad reqID \quad \mapsto i \in ReqIDs,$
$\quad value \quad \mapsto v \in Values,$
$\quad timestamp \mapsto t \in \text{Timestamps}\ ]$

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

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

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

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

$RepairReply$
$[\ src \mapsto r \in Replicas,$

$$dest \mapsto r \in Replicas,$$
$$req \mapsto (ClientRequest),$$
$$type \mapsto MRepairReply,$$
$$viewID \mapsto i \in ViewIDs,$$
$$index \mapsto i \in Indexes,$$
$$log \mapsto l \in Logs]$$

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

$ViewChangeReply$
$$[src \quad \mapsto r \in Replicas,$$
$$dest \quad \mapsto r \in Replicas,$$
$$type \quad \mapsto MViewChangeReply,$$
$$viewID \mapsto i \in ViewIDs,$$
$$logViewID \mapsto i \in ViewIDs,$$
$$log \quad \mapsto l \in Logs]$$

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

---

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 $cReqID$

4

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 $x \triangleq$ CHOOSE $x \in S :$ TRUE
       IN    $\langle x \rangle \circ SeqFromSet(S \setminus \{x\})$

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

$$\text{IN} \quad SetReduce(Op, \ S \setminus \{s\}, \ Op(s, \ value))$$

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'
$$Sum(S) \triangleq \text{LET } \_op(a, \ b) \triangleq a + b$$
$$\text{IN} \quad SetReduce(\_op, \ S, \ 0)$$

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
$$Quorum(v) \triangleq$$
$$\text{LET}$$
$$\quad quorumSize \triangleq Len(replicas) \div 2 \quad\quad + 1$$
$$\quad index(i) \quad\triangleq PrimaryIndex(v) + (i - 1)$$
$$\quad member(i) \quad\triangleq \text{IF } index(i) > Len(replicas) \text{ THEN } replicas[index(i)\%Len(replicas)] \text{ ELSE } replicas[inde$$
$$\text{IN}$$
$$\quad \{member(i) : i \in 1 \,..\, quorumSize\}$$

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.

Messages between processes are unordered and can be dropped by the network at any time.

Send a set of messages
$$Sends(ms) \triangleq$$
$$\quad \wedge\ messages' \quad\quad = messages \cup ms$$
$$\quad \wedge\ messageCount' = messageCount + Cardinality(ms)$$
$$\quad \wedge\ stepCount' \quad\ = stepCount + 1$$

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

*Ack* a message
$Ack(m) \triangleq$

$\quad \wedge messages' \quad = messages \setminus \{m\}$
$\quad \wedge messageCount' = messageCount + 1$
$\quad \wedge stepCount' \quad = stepCount + 1$

*Ack* a message and send a set of messages
$AckAndSends(m, ms) \triangleq$

$\quad \wedge messages' \quad = (messages \cup ms) \setminus \{m\}$
$\quad \wedge messageCount' = messageCount + Cardinality(ms)$
$\quad \wedge stepCount' \quad = stepCount + 1$

*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 *JITPaxos* clients.

Client 'c' sends value 'v' to all replicas
$ClientRequest(c, v) \triangleq$

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

Client 'c' handles a response 'm' from replica 'r'
$HandleClientReply(c, r, m) \triangleq$

If the reply view *ID* does not match the request view *ID*, update the client's view.
$\quad \wedge \quad \vee \quad \wedge m.viewID \neq m.req.viewID$
$\qquad\qquad \wedge \quad \vee \quad \wedge cViewID[c] < m.viewID$
$\qquad\qquad\qquad\qquad \wedge cViewID' = [cViewID \text{ EXCEPT } ![c] = m.viewID]$
$\qquad\qquad\qquad \vee \quad \wedge cViewID[c] \geq m.viewID$

7

$\land$ UNCHANGED $\langle cViewID \rangle$

$\quad \land Ack(m)$

$\quad \land$ UNCHANGED $\langle cReps,\ cCommits \rangle$

If the request and reply views match and the reply view matches the client's view, aggregate the replies for the associated client request.

$\lor\ \land m.viewID = m.req.viewID$

$\quad \land m.viewID = cViewID[c]$

$\quad \land\ \lor\ \land m.succeeded$

$\qquad\quad \land cReps' = [cReps$ EXCEPT $![c] =$

$\qquad\qquad\qquad\qquad (cReps[c] \setminus \{n\quad \in cReps[c] : \land n.src \qquad\quad = m.src$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \land n.viewID \quad = cViewID[c]$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \land n.req.reqID = m.req.reqID$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \land \neg n.succeeded\}) \cup \{m\}]$

$\quad\quad \lor\ \land \neg m.succeeded$

$\qquad\quad \land \neg \exists\, n \in cReps[c] : \land n.src \qquad\quad = m.src$

$\qquad\qquad\qquad\qquad\qquad\qquad \land n.viewID \quad = cViewID[c]$

$\qquad\qquad\qquad\qquad\qquad\qquad \land n.req.reqID = m.req.reqID$

$\qquad\qquad\qquad\qquad\qquad\qquad \land n.succeeded$

$\qquad\quad \land cReps' = [cReps$ EXCEPT $![c] = cReps[c] \cup \{m\}]$

$\quad \land$ LET $reps \qquad\qquad \triangleq\ \{n \in cReps'[c] : \land n.viewID \quad\ = cViewID[c]$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \land n.req.reqID\ = m.req.reqID\}$

$\qquad\quad isQuorum \quad\ \triangleq\ \{n.src : n \in \{n \in reps : n.succeeded\}\} = Quorum(cViewID[c])$

$\qquad\quad isCommitted\ \triangleq\ \land \forall\, n \in reps : n.succeeded$

$\qquad\qquad\qquad\qquad\qquad\quad \land Cardinality(\{n.checksum : n \in reps\}) = 1$

$\qquad\quad hasPrimary\ \triangleq\ \exists\, n \in reps : n.src = Primary(cViewID[c]) \land n.succeeded$

$\quad$ IN

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

$\qquad \lor\ \land isQuorum$

$\qquad\quad \land isCommitted$

$\qquad\quad \land$ LET $commit \triangleq$ CHOOSE $n \in reps : n.src = Primary(cViewID[c])$

$\qquad\qquad\quad$ IN $\quad cCommits' = [cCommits$ EXCEPT $![c] = cCommits[c] \cup \{commit\}]$

$\qquad\quad \land Ack(m)$

$\qquad$ 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*.

$\qquad \lor\ \land \neg isCommitted$

$\qquad\quad \land\ \lor\ \land hasPrimary$

$\qquad\qquad\quad \land$ LET $primaryRep \triangleq$ CHOOSE $n \in reps : \land n.src = Primary(cViewID[c])$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \land n.succeeded$

$\qquad\qquad\qquad\quad retryReps \quad\ \triangleq\ \{n \in reps :$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \land n.src \qquad\quad \neq Primary(cViewID[c])$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \land n.checksum \neq primaryRep.checksum\}$

$\qquad\qquad\qquad\quad$ IN $\quad AckAndSends(m, \{[src \qquad \mapsto c,$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad dest \qquad \mapsto r,$

8

$$
\begin{array}{ll}
type & \mapsto MReconcileRequest,\\
viewID & \mapsto cViewID[c],\\
reqID & \mapsto m.req.reqID,\\
index & \mapsto primaryRep.index] : n \in retryReps\})
\end{array}
$$

$$
\begin{array}{l}
\quad\quad\quad\quad\quad \vee \; \wedge \neg hasPrimary\\
\quad\quad\quad\quad\quad\quad\quad \wedge Ack(m)\\
\quad\quad\quad\quad \wedge \textsc{unchanged}\; \langle cCommits \rangle
\end{array}
$$

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

$$
\begin{array}{l}
\quad\quad\quad \vee \; \wedge \neg isQuorum\\
\quad\quad\quad\quad\quad \wedge isCommitted\\
\quad\quad\quad\quad\quad \wedge Ack(m)\\
\quad\quad\quad\quad\quad \wedge \textsc{unchanged}\; \langle cCommits \rangle\\
\quad\quad \wedge \textsc{unchanged}\; \langle cViewID \rangle\\
\quad \wedge \textsc{unchanged}\; \langle replicaVars,\; cTime,\; cReqID \rangle
\end{array}
$$

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

---

This section models *JITPaxos* replicas.

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.
$$
\begin{array}{l}
\wedge \;\vee\; \wedge m.viewID = rViewID[r]\\
\quad\quad\quad \wedge \textsc{let}\; lastTimestamp \;\triangleq\; Max(\{rLog[r][i].timestamp : i \in \textsc{domain}\; rLog[r]\} \cup \{0\})\\
\quad\quad\quad \textsc{in}
\end{array}
$$

     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.

$$
\begin{array}{l}
\quad\quad\quad\quad \wedge \;\vee\; \wedge m.timestamp > lastTimestamp\\
\quad\quad\quad\quad\quad\quad\quad \wedge rLog' = [rLog\; \textsc{except}\; ![r] =\\
\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad Append(rLog[r], [value \;\;\;\;\;\; \mapsto m.value,\\
\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\; timestamp \mapsto m.timestamp])]
\end{array}
$$

$$
\begin{array}{ll}
\quad\quad\quad\quad\quad\quad \wedge Reply(m, [src & \mapsto r,\\
& dest & \mapsto c,\\
& req & \mapsto m,\\
& type & \mapsto MClientReply,\\
& viewID & \mapsto rViewID[r],\\
& index & \mapsto Len(rLog'[r]),\\
& checksum & \mapsto rLog'[r],\\
& value & \mapsto m.value,\\
& timestamp & \mapsto m.timestamp,\\
& succeeded & \mapsto \textsc{true}])
\end{array}
$$

    If the request timestamp matches the highest *log* timestamp, treat the

request as a duplicate. Return a successful response indicating the
entry was appended.

$\lor\ \land\ m.timestamp = lastTimestamp$

$\qquad \land\ Reply(m, [src \qquad\quad \mapsto r,$

$\qquad\qquad\qquad\quad dest \qquad\quad \mapsto c,$

$\qquad\qquad\qquad\quad req \qquad\qquad \mapsto m,$

$\qquad\qquad\qquad\quad type \qquad\quad \mapsto MClientReply,$

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

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

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

$\qquad\qquad\qquad\quad value \qquad\ \ \mapsto m.value,$

$\qquad\qquad\qquad\quad timestamp \mapsto m.timestamp,$

$\qquad\qquad\qquad\quad succeeded\ \ \mapsto \text{TRUE}])$

$\qquad \land\ \text{UNCHANGED}\ \langle rLog\rangle$

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

$\lor\ \land\ m.timestamp < lastTimestamp$

$\qquad \land\ Reply(m, [src \qquad\quad \mapsto r,$

$\qquad\qquad\qquad\quad dest \qquad\quad \mapsto c,$

$\qquad\qquad\qquad\quad req \qquad\qquad \mapsto m,$

$\qquad\qquad\qquad\quad type \qquad\quad \mapsto MClientReply,$

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

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

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

$\qquad\qquad\qquad\quad value \qquad\ \ \mapsto m.value,$

$\qquad\qquad\qquad\quad timestamp \mapsto m.timestamp,$

$\qquad\qquad\qquad\quad succeeded\ \ \mapsto \text{FALSE}])$

$\qquad \land\ \text{UNCHANGED}\ \langle rLog\rangle$

$\land\ \text{UNCHANGED}\ \langle rViewID,\ rStatus,\ rViewChangeReps\rangle$

If the client's view is greater than the replica's view, reject the client's
request with the outdated view *ID* and enter the view change protocol.

$\lor\ \land\ m.viewID > rViewID[r]$

$\quad \land\ rViewID' \qquad\qquad = [rViewID \qquad\qquad \text{EXCEPT} \ ![r] \ = m.viewID]$

$\quad \land\ rStatus' \qquad\qquad\ = [rStatus \qquad\qquad\ \text{EXCEPT} \ ![r] \ = SViewChange]$

$\quad \land\ rViewChangeReps' = [rViewChangeReps\ \text{EXCEPT}\ ![r] = \{\}]$

$\quad \land\ Replies(m, \{[src \qquad\quad \mapsto r,$

$\qquad\qquad\qquad\quad dest \qquad\quad \mapsto c,$

$\qquad\qquad\qquad\quad req \qquad\qquad \mapsto m,$

$\qquad\qquad\qquad\quad type \qquad\quad \mapsto MClientReply,$

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

$\qquad\qquad\qquad\quad succeeded\ \mapsto \text{FALSE}],$

$\qquad\qquad\qquad [src \qquad\qquad \mapsto r,$

$\qquad\qquad\qquad\quad dest \qquad\quad \mapsto Primary(m.viewID),$

$\qquad\qquad\qquad\quad type \qquad\quad \mapsto MViewChangeReply,$

$\qquad\qquad\qquad\quad viewID \qquad \mapsto m.viewID,$

$$
\begin{aligned}
&\qquad\qquad\qquad logViewID \mapsto rLogViewID[r],\\
&\qquad\qquad\qquad log \qquad\quad \mapsto rLog[r]]\})
\end{aligned}
$$

$\quad \land \text{UNCHANGED } \langle rLog \rangle$

If the client's view is less than the replica's view, reject the client's request with the updated view $ID$ to force the client to retry.

$\quad \lor\; \land m.viewID < rViewID[r]$

$\qquad \land Reply(m, [src \qquad\;\; \mapsto r,$
$\qquad\qquad\qquad\quad dest \qquad\;\; \mapsto c,$
$\qquad\qquad\qquad\quad req \qquad\;\; \mapsto m,$
$\qquad\qquad\qquad\quad type \qquad\;\; \mapsto MClientReply,$
$\qquad\qquad\qquad\quad viewID \quad\;\; \mapsto rViewID[r],$
$\qquad\qquad\qquad\quad succeeded \mapsto \text{FALSE}])$

$\qquad \land \text{UNCHANGED } \langle rViewID, rStatus, rLog, rViewChangeReps \rangle$

$\quad \land \text{UNCHANGED } \langle clientVars, rLogViewID, rSyncIndex \rangle$

$HandleReconcileRequest(r, c, m) \;\triangleq$

$\quad \land rStatus[r] \quad = SInSync$
$\quad \land rViewID[r] = m.viewID$
$\quad \land\; \lor\; \land rSyncIndex[r] \geq m.index$

$\qquad\qquad \land Reply(m, [src \qquad\;\; \mapsto r,$
$\qquad\qquad\qquad\qquad\;\; dest \qquad\;\; \mapsto c,$
$\qquad\qquad\qquad\qquad\;\; req \qquad\;\; \mapsto m,$
$\qquad\qquad\qquad\qquad\;\; type \qquad\;\;\; \mapsto MReconcileReply,$
$\qquad\qquad\qquad\qquad\;\; viewID \quad\;\; \mapsto rViewID[r],$
$\qquad\qquad\qquad\qquad\;\; index \qquad\; \mapsto m.index,$
$\qquad\qquad\qquad\qquad\;\; checksum \;\; \mapsto [i \in 1 \,..\, m.index \mapsto rLog[r][i]],$
$\qquad\qquad\qquad\qquad\;\; value \qquad\;\; \mapsto rLog[r][m.index].value,$
$\qquad\qquad\qquad\qquad\;\; timestamp \mapsto rLog[r][m.index].timestamp,$
$\qquad\qquad\qquad\qquad\;\; succeeded \;\; \mapsto \text{TRUE}])$

$\qquad\quad \land \text{UNCHANGED } \langle rStatus \rangle$

$\qquad \lor\; \land rSyncIndex[r] < m.index$
$\qquad\qquad \land Primary(rViewID[r]) \neq r$
$\qquad\qquad \land rStatus' = [rStatus \text{ EXCEPT } ![r] = SRepair]$
$\qquad\qquad \land AckAndSend(m, [src \qquad\; \mapsto r,$
$\qquad\qquad\qquad\qquad\qquad\; dest \qquad\; \mapsto Primary(rViewID[r]),$
$\qquad\qquad\qquad\qquad\qquad\; req \qquad\;\; \mapsto m,$
$\qquad\qquad\qquad\qquad\qquad\; type \qquad\; \mapsto MRepairRequest,$
$\qquad\qquad\qquad\qquad\qquad\; viewID \mapsto rViewID[r],$
$\qquad\qquad\qquad\qquad\qquad\; index \quad\; \mapsto m.index])$

$\quad \land \text{UNCHANGED } \langle clientVars, rViewID, rLog, rLogViewID, rSyncIndex, rViewChangeReps \rangle$

$HandleRepairRequest(r, s, m) \;\triangleq$

$\quad \land rStatus[r] \quad = SInSync$
$\quad \land rViewID[r] = m.viewID$
$\quad \land Primary(rViewID[r]) = r$

$\qquad \wedge\; Reply(m, [src \quad\; \mapsto r,$
$\qquad\qquad\qquad\quad\; dest \quad\;\; \mapsto s,$
$\qquad\qquad\qquad\quad\; req \quad\;\; \mapsto m.req,$
$\qquad\qquad\qquad\quad\; type \quad\;\; \mapsto MRepairReply,$
$\qquad\qquad\qquad\quad\; viewID \mapsto rViewID[r],$
$\qquad\qquad\qquad\quad\; index \quad \mapsto m.index,$
$\qquad\qquad\qquad\quad\; log \quad\;\;\; \mapsto [i \in 1 \mathinner{\ldotp\ldotp} m.index \mapsto rLog[r][i]]])$
$\qquad \wedge\; \text{UNCHANGED}\; \langle clientVars,\, replicaVars \rangle$

$HandleRepairReply(r,\, s,\, m)\; \triangleq$
$\qquad \wedge\; rStatus[r]\;\;\; = SRepair$
$\qquad \wedge\; rViewID[r] = m.viewID$
$\qquad \wedge\; rStatus' \qquad = [rStatus \qquad \text{EXCEPT}\; ![r]\;\; = SInSync]$
$\qquad \wedge\; rLog' \qquad\;\; = [rLog \qquad\;\; \text{EXCEPT}\; ![r]\;\; = m.log \circ SubSeq(rLog[r],\, Len(m.log),\, Len(rLog[r]))]$
$\qquad \wedge\; rSyncIndex' = [rSyncIndex\; \text{EXCEPT}\; ![r] = Len(rLog'[r])]$
$\qquad \wedge\; Reply(m, [src \qquad\;\; \mapsto r,$
$\qquad\qquad\qquad\quad\; dest \qquad\;\;\; \mapsto m.req.src,$
$\qquad\qquad\qquad\quad\; req \qquad\;\;\; \mapsto m.req,$
$\qquad\qquad\qquad\quad\; type \qquad\;\;\; \mapsto MReconcileReply,$
$\qquad\qquad\qquad\quad\; viewID \qquad \mapsto rViewID[r],$
$\qquad\qquad\qquad\quad\; index \qquad\; \mapsto m.index,$
$\qquad\qquad\qquad\quad\; checksum \;\; \mapsto m.log,$
$\qquad\qquad\qquad\quad\; value \qquad\; \mapsto m.log[m.index].value,$
$\qquad\qquad\qquad\quad\; timestamp \mapsto m.log[m.index].timestamp,$
$\qquad\qquad\qquad\quad\; succeeded \;\; \mapsto \text{TRUE}])$
$\qquad \wedge\; \text{UNCHANGED}\; \langle clientVars,\, rViewID,\, rLogViewID,\, rViewChangeReps \rangle$

Replica 'r' requests a view change
$ChangeView(r)\; \triangleq$
$\qquad \wedge\; Sends(\{[src \qquad \mapsto r,$
$\qquad\qquad\qquad\; dest \quad\;\; \mapsto d,$
$\qquad\qquad\qquad\; type \quad\;\; \mapsto MViewChangeRequest,$
$\qquad\qquad\qquad\; viewID \mapsto rViewID[r] + 1] : d \in Replicas\})$
$\qquad \wedge\; \text{UNCHANGED}\; \langle clientVars,\, replicaVars \rangle$

Replica 'r' handles replica 's' view change request 'm'
$HandleViewChangeRequest(r,\, s,\, m)\; \triangleq$
$\qquad \wedge\; \vee\; \wedge\; rViewID[r] < m.viewID$
$\qquad\qquad\quad \wedge\; rViewID' \qquad\qquad = [rViewID \qquad\qquad \text{EXCEPT}\; ![r]\;\; = m.viewID]$
$\qquad\qquad\quad \wedge\; rStatus' \qquad\qquad = [rStatus \qquad\qquad \text{EXCEPT}\; ![r]\;\;\; = SViewChange]$
$\qquad\qquad\quad \wedge\; rViewChangeReps' = [rViewChangeReps\; \text{EXCEPT}\; ![r] = \{\}]$
$\qquad\qquad\quad \wedge\; Reply(m, [src \qquad\;\; \mapsto r,$
$\qquad\qquad\qquad\qquad\;\; dest \qquad\;\;\; \mapsto Primary(m.viewID),$
$\qquad\qquad\qquad\qquad\;\; type \qquad\;\;\; \mapsto MViewChangeReply,$
$\qquad\qquad\qquad\qquad\;\; viewID \qquad \mapsto m.viewID,$
$\qquad\qquad\qquad\qquad\;\; logViewID \mapsto rLogViewID[r],$

$$\qquad\qquad log \qquad\quad \mapsto rLog[r]])$$

$$\lor \;\land\; rViewID[r] \geq m.viewID$$
$$\quad\; \land\; Ack(m)$$
$$\quad\; \land\; \text{UNCHANGED } \langle rViewID,\; rStatus,\; rViewChangeReps \rangle$$
$$\land\; \text{UNCHANGED } \langle clientVars,\; rLog,\; rLogViewID,\; rSyncIndex \rangle$$

Replica 'r' handles replica 's' view change reply 'm'
$$HandleViewChangeReply(r,\; s,\; m) \;\triangleq$$

       The view change protocol is run by the primary for the view.

$$\land\; Primary(m.viewID) = r$$
$$\land\; rViewID[r] = m.viewID$$
$$\land\; rStatus[r] \;\; = SViewChange$$
$$\land\; rViewChangeReps' = [rViewChangeReps \text{ EXCEPT } ![r] = rViewChangeReps[r] \cup \{m\}]$$
$$\land\; \text{LET } viewChanges \;\triangleq\; \{v \in rViewChangeReps'[r] : v.viewID = rViewID[r]\}$$
$$\quad\; \text{IN}$$

           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.

$$\lor \;\land\; IsLocalQuorum(r,\; \{v.src : v \in viewChanges\})$$
$$\quad\; \land\; \text{LET } latestViewID \;\triangleq\; Max(\{v.logViewID : v \in viewChanges\})$$
$$\qquad\qquad\quad latestChange \;\triangleq\; \text{CHOOSE } v \in viewChanges :$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \land\quad v.logViewID = latestViewID$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \land\quad v.src \in Quorum(latestViewID)$$
$$\quad\;\;\; \text{IN} \quad AckAndSends(m,\; \{[src \qquad \mapsto r,$$
$$\qquad\qquad\qquad\qquad\qquad\qquad dest \qquad \mapsto d,$$
$$\qquad\qquad\qquad\qquad\qquad\qquad type \qquad \mapsto MStartViewRequest,$$
$$\qquad\qquad\qquad\qquad\qquad\qquad viewID \mapsto rViewID[r],$$
$$\qquad\qquad\qquad\qquad\qquad\qquad log \qquad \mapsto latestChange.log] : d \in Replicas\})$$

           If view change replies have not yet been received from a quorum, record
the view change reply and discard the message.

$$\lor \;\land\; \neg IsLocalQuorum(r,\; \{v.src : v \in viewChanges\})$$
$$\quad\; \land\; Ack(m)$$
$$\land\; \text{UNCHANGED } \langle clientVars,\; rStatus,\; rViewID,\; rLog,\; rLogViewID,\; rSyncIndex \rangle$$

Replica 'r' handles replica 's' start view request 'm'
$$HandleStartViewRequest(r,\; s,\; m) \;\triangleq$$

       To activate a view, the replica must either not know of the view or already
be participating in the view change protocol for the view.

$$\land\; \lor\; rViewID[r] < m.viewID$$
$$\quad\; \lor\; \land\; rViewID[r] = m.viewID$$
$$\qquad\;\; \land\; rStatus[r] \;\; = SViewChange$$

       If the replica is part of the quorum for the activated view, update the *log*
and record the activated view for use in the view change protocol.

$\land\ \lor\ \land\ r \in Quorum(m.viewID)$
$\qquad\ \land\ rLog' \qquad\quad = [rLog \qquad\quad \text{EXCEPT }![r] \quad = m.log]$
$\qquad\ \land\ rLogViewID' = [rLogViewID \text{ EXCEPT }![r] = m.viewID]$
$\qquad\ \land\ rSyncIndex' \ = [rSyncIndex \text{ EXCEPT }![r] \ = Len(m.log)]$
$\quad\ \lor\ \land\ r \notin Quorum(m.viewID)$
$\qquad\ \land\ \text{UNCHANGED }\langle rLog,\ rLogViewID,\ rSyncIndex\rangle$

Update the replica's view *ID* and status and clean up view change state.

$\land\ rViewID' = [rViewID \text{ EXCEPT }![r] = m.viewID]$
$\land\ rStatus' \ = [rStatus \text{ EXCEPT }![r] \ = SInSync]$
$\land\ \text{LET } viewChanges \ \triangleq\ \{v \in rViewChangeReps[r] : v.viewID = rViewID[r]\}$
$\quad\ \text{IN}\quad rViewChangeReps' = [rViewChangeReps \text{ EXCEPT }![r] = rViewChangeReps[r] \setminus viewChanges]$
$\land\ Ack(m)$
$\land\ \text{UNCHANGED }\langle clientVars\rangle$

---

$InitMessageVars \ \triangleq$
$\quad \land\ messages \qquad\ = \{\}$
$\quad \land\ messageCount = 0$
$\quad \land\ stepCount \qquad = 0$

$InitClientVars \ \triangleq$
$\quad \land\ cTime \qquad = 0$
$\quad \land\ cViewID \ \ = [c \in Clients \mapsto 1]$
$\quad \land\ cReqID \ \ \ = [c \in Clients \mapsto 0]$
$\quad \land\ cReps \qquad = [c \in Clients \mapsto \{\}]$
$\quad \land\ cCommits = [c \in Clients \mapsto \{\}]$

$InitReplicaVars \ \triangleq$
$\quad \land\ rStatus \qquad\qquad\ = [r \in Replicas \mapsto SInSync]$
$\quad \land\ rViewID \qquad\qquad = [r \in Replicas \mapsto 1]$
$\quad \land\ rLog \qquad\qquad\quad = [r \in Replicas \mapsto \langle\rangle]$
$\quad \land\ rSyncIndex \qquad\ = [r \in Replicas \mapsto 0]$
$\quad \land\ rLogViewID \qquad = [r \in Replicas \mapsto 1]$
$\quad \land\ rViewChangeReps = [r \in Replicas \mapsto \{\}]$

$Init \ \triangleq$
$\quad \land\ InitMessageVars$
$\quad \land\ InitClientVars$
$\quad \land\ InitReplicaVars$

---

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.

$Linearizability \triangleq$
  $\forall\, c \in Clients :$
    $\forall\, e \in cCommits[c] :$
      $\neg\exists\, r \in Replicas :$
          $\land\ rStatus[r]\quad = SInSync$
          $\land\ rViewID[r] \geq e.viewID$
          $\land\ r \in Quorum(rViewID[r])$
          $\land\ rLog[r][e.index].value \neq e.value$

---

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

$NextChangeView \triangleq$
  $\exists\, r \in Replicas :$
    $ChangeView(r)$

$NextHandleClientRequest \triangleq$
  $\exists\, m \in messages :$
    $\land\ m.type = MClientRequest$
    $\land\ HandleClientRequest(m.dest,\, m.src,\, m)$

$NextHandleClientReply \triangleq$
  $\exists\, m \in messages :$
    $\land\ m.type = MClientReply$
    $\land\ HandleClientReply(m.dest,\, m.src,\, m)$

$NextHandleReconcileRequest \triangleq$
  $\exists\, m \in messages :$
    $\land\ m.type = MReconcileRequest$
    $\land\ HandleReconcileRequest(m.dest,\, m.src,\, m)$

$NextHandleReconcileReply \triangleq$
  $\exists\, m \in messages :$
    $\land\ m.type = MReconcileReply$
    $\land\ HandleReconcileReply(m.dest,\, m.src,\, m)$

$NextHandleRepairRequest \triangleq$
  $\exists\, m \in messages :$
    $\land\ m.type = MRepairRequest$
    $\land\ HandleRepairRequest(m.dest,\, m.src,\, m)$

$NextHandleRepairReply \triangleq$
  $\exists\, m \in messages :$
    $\land\ m.type = MRepairReply$

$\land$ *HandleRepairReply*(*m.dest*, *m.src*, *m*)

*NextHandleViewChangeRequest* $\triangleq$
    $\exists\, m \in messages :$
      $\land\, m.type = MViewChangeRequest$
      $\land\, HandleViewChangeRequest(m.dest, m.src, m)$

*NextHandleViewChangeReply* $\triangleq$
    $\exists\, m \in messages :$
      $\land\, m.type = MViewChangeReply$
      $\land\, HandleViewChangeReply(m.dest, m.src, m)$

*NextHandleStartViewRequest* $\triangleq$
    $\exists\, m \in messages :$
      $\land\, m.type = MStartViewRequest$
      $\land\, HandleStartViewRequest(m.dest, m.src, m)$

*NextDropMessage* $\triangleq$
    $\exists\, m \in messages :$
      $\land\, Ack(m)$
      $\land\,$ UNCHANGED $\langle clientVars,\ replicaVars \rangle$

*Next* $\triangleq$
    $\lor\, NextClientRequest$
    $\lor\, NextChangeView$
    $\lor\, NextHandleClientRequest$
    $\lor\, NextHandleClientReply$
    $\lor\, NextHandleReconcileRequest$
    $\lor\, NextHandleReconcileReply$
    $\lor\, NextHandleRepairRequest$
    $\lor\, NextHandleRepairReply$
    $\lor\, NextHandleViewChangeRequest$
    $\lor\, NextHandleViewChangeReply$
    $\lor\, NextHandleStartViewRequest$
    $\lor\, NextDropMessage$

*Spec* $\triangleq$ *Init* $\land\, \Box[Next]_{vars}$