# N TLA+ Specification

———————————————— MODULE $Nezha$ ————————————————

EXTENDS $Naturals$, $TLC$, $FiniteSets$, $Sequences$

## Bounds for Model Check [Configurable]

Time Range [Configurable]
$MaxTime \triangleq 3$

Each client is only allowed to submit $MaxReqNum$ requests [Configurable]
In the specification, we will only consider two roles, client and replicas
($i.e.$ it can be considered as co-locating one proxy with one client)
For the proxy-based design, we just need to replace client with proxy,
and then the specification describes the interaction between proxy and replicas
$MaxReqNum \triangleq 1$

The leader is only allowed to crash when the view $< MaxViews$ [Configurable]
$MaxViews \triangleq 3$

These variables are used to implment at-most-once primitives
$i.e.$ The variables record the messages processed by $Replicas/Clients$, so
that the $Replicas/Clients$ will not process twice
VARIABLE $vReplicaProcessed$,    Messages that have been processed by replicas
             $vClientProcessed$   Messages that have been processed by clients

VARIABLE $DebugAction$

## Constants

The set of replicas and an ordering of them
CONSTANTS $Replicas$, $ReplicaOrder$, $Clients$, $LatencyBounds$
ASSUME $IsFiniteSet(Replicas)$
ASSUME $ReplicaOrder \in Seq(Replicas)$

$F \triangleq (Cardinality(Replicas) - 1) \div 2$
$ceilHalfF \triangleq$ IF $(F \div 2) * 2 = F$ THEN $F \div 2$ ELSE $(F + 1) \div 2$
$floorHalfF \triangleq F \div 2$
$QuorumSize \triangleq F + 1$
$FastQuorumSize \triangleq F + ceilHalfF + 1$
$RecoveryQuorumSize \triangleq ceilHalfF + 1$
$FastQuorums \triangleq \{R \in$ SUBSET $(Replicas) : Cardinality(R) \geq FastQuorumSize\}$
$Quorums \triangleq \{R \in$ SUBSET $(Replicas) : Cardinality(R) * 2 > Cardinality(Replicas)\}$

Replica Statuses
$StNormal \triangleq 1$

$StViewChange \triangleq 2$
$StRecovering \triangleq 3$

Message Types
$MClientRequest \triangleq 1$  Sent by client to replicas
$MFastReply \triangleq 2$  Fast Reply Message
$MSlowReply \triangleq 3$  Slow Reply Message
$MLogIndex \triangleq 4$   $LogIndex$
$MLogEntry \triangleq 5$   Log entry, different from index, it includes command field, which can be large in practice
$MIndexSync \triangleq 6$  Sync message during the index sync process
$MMissEntryRequest \triangleq 7$  Sent by followers once they fail to find the entry on itself
$MMissEntryReply \triangleq 8$   Response to $MMissEntryRequest$, providing the missing entries

$MViewChangeReq \triangleq 9$       Sent when leader/sequencer failure detected
$MViewChange \triangleq 10$       Sent to $ACK$ view change
$MStartView \triangleq 11$       Sent by new leader to start view

The following messages are mainly used for periodic sync
Just as described in $NOPaxos$, it is an optional optimization to enable fast recovery after failure
$MSyncPrepare \triangleq 12$       Sent by the leader to ensure $log$ durability
$MSyncRep \triangleq 13$       Sent by followers as $ACK$
$MSyncCommit \triangleq 14$        Sent by leaders to indicate stable $log$

The following messages are mainly used for replica recovery
$MCrashVectorReq \triangleq 15$
$MCrashVectorRep \triangleq 16$
$MRecoveryReq \triangleq 17$
$MRecoveryRep \triangleq 18$
$MStateTransferReq \triangleq 19$
$MStateTransferRep \triangleq 20$

**Message Schemas**

$ViewIDs \triangleq [\, leaderNum \mapsto n \in (1 \,..\,)\,]$

\ * $< clientID, requestID >$  uniquely identifies one request on one replica
\ * But across replicas, the same $< clientID, requestID >$  may have different deadlines
\ * (the leader may modify the deadline to make the request eligible to enter the early-buffer)
\ * so $< deadline, clientID, reqID >$  uniquely identifes one request across replicas

$ClientRequest$
  $[\, mtype \quad\; \mapsto MClientRequest,$
   $sender \quad\; \mapsto c \in Clients,$
   $dest \quad\;\; \mapsto r \in Replicas,$
   $requestID \mapsto i \in (1 \,..\,),$
   $command \quad \mapsto$ "",
   $s \qquad\;\; \mapsto t \in (1 \,..\, MaxTime),$
   $l \qquad\;\; \mapsto l \in (1 \,..\, MaxBound)$
  $]$

\ ∗ *logSlotNum* is not necessary and it is not described in the paper
\ ∗  Here we include *logSlotNum* in *FastReply* and *SlowReply messages*
\ ∗  to facilitate the check of *Linearizability* invariant
*FastReply*
  [ *mtype*        ↦ *MFastReply*,
    *sender*      ↦ *r* ∈ *Replicas*,
    *dest*         ↦ *c* ∈ *Clients*,
    *viewID*      ↦ *v* ∈ *ViewIDs*,
    *requestID* ↦ *i* ∈  (1 . . *vClientReqNum*)
    *hash*        ↦  [ *log* ↦ *vLogs*[1 . . *n*],
                    *cv* ↦ *crashVector*
                  ]
    *deadline* ↦ *i* ∈  (1 . . *MaxTime* + *MaxBound*),
    *logSlotNum* ↦ *n* ∈  (1 . . )
  ]

*SlowReply*
  [ *mtype*        ↦ *MSlowReply*,
    *sender*       ↦ *r* ∈ *Replicas*,
    *dest*          ↦ *c* ∈ *Clients*,
    *viewID*      ↦ *v* ∈ *ViewIDs*,
    *requestID* ↦ *i* ∈  (1 . . *vClientReqNum*)
    *logSlotNum* ↦ *n* ∈  (1 . . )
  ]

*LogIndex*
  [ *mtype*        ↦ *MLogIndex*,
    *clientID* ↦ *c* ∈ *Clients*,
    *requestID* ↦ *i* ∈  (1 . . *vClientReqNum*),
    *deadline* ↦ *i* ∈  (1 . . *MaxTime* + *MaxBound*),
  ]

*LogEntry*
  [ *mtype*         ↦ *MLogEntry*,
    *clientID* ↦ *c* ∈ *Clients*,
    *requestID* ↦ *i* ∈  (1 . . *vClientReqNum*),
    *deadline* ↦ *i* ∈  (1 . . *MaxTime* + *MaxBound*),
    *command* ↦ ""
  ]

*IndexSync*
  [ *mtype*        ↦ *MIndexSync*,
    *sender*       ↦ *r* ∈ *Replicas*,
    *dest*          ↦ *c* ∈ *Clients*,
    *viewID*       ↦ *v* ∈ *ViewIDs*,
    *logindcies* ↦ *index* ∈ *vLogs*[*leaderIdx*]
  ]

 *MMissEntryRequest*
  [ *mtype*          ↦ *MMissEntryRequest*,
    *sender*        ↦ *r* ∈ *Replicas*,
    *dest*          ↦ *d* ∈ *Replicas*,

$$
\begin{array}{ll}
viewID & \mapsto v \in \mathit{ViewIDs}, \\
miss & \mapsto \{log\ indices\} \\
\end{array}
$$
$\ ]$

$\mathit{MMissEntryRequest}$
$$
\begin{array}{ll}
[\ mtype & \mapsto \mathit{MMissEntryReply}, \\
sender & \mapsto r \in \mathit{Replicas}, \\
dest & \mapsto d \in \mathit{Replicas}, \\
viewID & \mapsto v \in \mathit{ViewIDs}, \\
entries & \mapsto \{log\ entries\} \\
\end{array}
$$
$\ ]$

$\mathit{ViewChangeReq}$
$$
\begin{array}{l}
[\ mtype \mapsto \mathit{MViewChangeReq}, \\
sender \mapsto r \in \mathit{Replicas}, \\
dest \mapsto r \in \mathit{Replicas}, \\
viewID \mapsto v \in \mathit{ViewIDs}, \\
cv \quad \mapsto \text{crash vector} \\
\end{array}
$$
$\ ]$

$\mathit{ViewChange}$
$$
\begin{array}{ll}
[\ mtype & \mapsto \mathit{MViewChange}, \\
sender & \mapsto r \in \mathit{Replicas}, \\
dest & \mapsto r \in \mathit{Replicas}, \\
viewID & \mapsto v \in \mathit{ViewIDs}, \\
lastNormal & \mapsto v \in \mathit{ViewIDs}, \\
log & \mapsto l \in vLogs[1 \mathinner{.\,.} n], \\
cv & \mapsto \text{crash vector} \\
\end{array}
$$
$\ ]$

$\mathit{StartView}$
$$
\begin{array}{ll}
[\ mtype & \mapsto \mathit{MStartView}, \\
dest & \mapsto r \in \mathit{Replicas}, \\
viewID & \mapsto v \in \mathit{ViewIDs}, \\
log & \mapsto l \in vLogs[1 \mathinner{.\,.} n], \\
cv & \mapsto \text{crash vector} \\
\end{array}
$$
$\ ]$

$\mathit{SyncPrepare}$
$$
\begin{array}{ll}
[\ mtype & \mapsto \mathit{MSyncPrepare}, \\
dest & \mapsto r \in \mathit{Replicas}, \\
sender & \mapsto r \in \mathit{Replicas}, \\
viewID & \mapsto v \in \mathit{ViewIDs}, \\
log & \mapsto l \in vLogs[1 \mathinner{.\,.} n]\ ] \\
\end{array}
$$

$\mathit{SyncRep}$
$$
\begin{array}{ll}
[\ mtype & \mapsto \mathit{MSyncRep}, \\
dest & \mapsto r \in \mathit{Replicas}, \\
sender & \mapsto r \in \mathit{Replicas}, \\
viewID & \mapsto v \in \mathit{ViewIDs}, \\
logSlotNumber & \mapsto n \in\ (1 \mathinner{.\,.})\ ] \\
\end{array}
$$

$\mathit{SyncCommit}$

$[\ mtype \quad\quad \mapsto MSyncCommit,$
$\quad dest \quad\quad\quad \mapsto r \in Replicas,$
$\quad sender \quad\quad \mapsto r \in Replicas,$
$\quad viewID \quad\quad \mapsto v \in ViewIDs,$
$\quad log \quad\quad\quad \mapsto l \in vLogs[1\ ..\ n]\ ]$

$CrashVectorReq$
$[\ mtype \quad\quad \mapsto MCrashVectorReq,$
$\quad sender \quad\quad \mapsto r \in Replicas,$
$\quad dest \quad\quad\quad \mapsto r \in Replicas,$
$\quad nonce \quad\quad \mapsto nonce$
$]$

$CrashVectorRep$
$[\ mtype \quad\quad \mapsto MCrashVectorRep,$
$\quad sender \quad\quad \mapsto r \in Replicas,$
$\quad dest \quad\quad\quad \mapsto r \in Replicas,$
$\quad nonce \quad\quad \mapsto nonce,$
$\quad cv \quad\quad\quad \mapsto$ vector of counters
$]$

$RecoveryReq$
$[\ mtype \quad\quad \mapsto MRecoveryReq,$
$\quad sender \quad\quad \mapsto r \in Replicas,$
$\quad dest \quad\quad\quad \mapsto r \in Replicas,$
$\quad cv \quad\quad\quad \mapsto$ vector of counters
$]$

$RecoveryRep$
$[\ mtype \quad\quad \mapsto MRecoveryRep,$
$\quad sender \quad\quad \mapsto r \in Replicas,$
$\quad dest \quad\quad\quad \mapsto r \in Replicas,$
$\quad viewID \quad\quad \mapsto v \in ViewIDs,$
$\quad cv \quad\quad\quad \mapsto$ vector of counters
$]$

$StateTransferReq$
$[\ mtype \quad\quad \mapsto MStateTransferReq,$
$\quad sender \quad\quad \mapsto r \in Replicas,$
$\quad dest \quad\quad\quad \mapsto r \in Replicas,$
$\quad cv \quad\quad\quad \mapsto$ vector of counters
$]$

$StateTransferRep$
$[\ mtype \quad\quad \mapsto MStateTransferRep,$
$\quad sender \quad\quad \mapsto r \in Replicas,$
$\quad dest \quad\quad\quad \mapsto r \in Replicas,$
$\quad viewID \quad\quad \mapsto v \in ViewIDs,$
$\quad log \quad\quad\quad \mapsto l \in vLogs[1\ ..\ n]\ ],$
$\quad cv \quad\quad\quad \mapsto$ vector of counters
$]$

## Variables

**Network State**

VARIABLE *messages*   Set of all messages sent

$$networkVars \quad \triangleq \langle messages \rangle$$
$$InitNetworkState \triangleq messages = \{\}$$

Used as a dummy value
$$NULLLog \triangleq [ \quad deadline \quad \mapsto 0,$$
$$clientID \quad \mapsto 0,$$
$$requestID \quad \mapsto 0$$
$$]$$

**Replica State**

| VARIABLES | | |
|---|---|---|
| *vLog*, | Log of values |
| *vEarlyBuffer*, | The early buffer to hold request, |
| | and release it after clock passes its deadline $(s + l)$ |
| *vReplicaStatus*, | One of *StNormal*, *StViewChange*, *StRecovering* |
| *vViewID*, | Current *viewID* replicas recognize |
| *vReplicaClock*, | Current Time of the replica |
| *vLastNormView*, | Last views in which replicas had status *StNormal* |
| *vViewChanges*, | Used for logging view change votes |
| *vSyncPoint*, | Latest synchronization point, |
| | to which the replica state (*vLog*) is consistent with the leader. |
| *vLateBuffer*, | The late buffer Used to store the requests |
| | which are not eligible to enter *vEarlyBuffer* |
| *vTentativeSync*, | Used by leader to mark current *syncPrepare* point (during periodic sync process) |
| | (Actually, *vSyncPoint* and *vTentativeSync* can be merged into one *Var* |
| | However, we decouple them to make the spec easy to understand) |
| *vSyncReps*, | Used for logging sync reps at leader |
| *vCommitPoint*, | Different from *vSyncPoint*, |
| | *vCommitPoint* indicates that the logs before this point has been replicated to majority |
| | So followers can safely execute requests (*log* entries) up to *vCommitPoint* |
| | Refer to "Acceleration of Recovery" para in *Sec* 6 |
| *vUUIDCounter*, | Locally unique string (for *CrashVectorReq*) |
| *vCrashVector*, | *CrashVector*, initialized as all-zero vector |
| *vCrashVectorReps*, | *CrashVectorRep* Set |
| *vRecoveryReps* | *RecoveryRep* Set |

$$replicaVars \quad \triangleq \langle vLog, vEarlyBuffer,$$
$$vViewID, vReplicaClock,$$
$$vLastNormView, vViewChanges, vReplicaStatus,$$
$$vSyncPoint, vLateBuffer,$$
$$vTentativeSync, vSyncReps, vCommitPoint,$$
$$vUUIDCounter, vCrashVector,$$

$$\langle vCrashVectorReps, vRecoveryReps\rangle$$

$InitReplicaState \triangleq$
  $\land vLog \qquad\qquad = [r \in Replicas \mapsto \langle\rangle]$
  $\land vEarlyBuffer \qquad = [r \in Replicas \mapsto \{\}]$
  $\land vViewID \qquad\quad = [r \in Replicas \mapsto 1]$    0 should also be okay
  $\land vReplicaClock \qquad = [r \in Replicas \mapsto 1]$
  $\land vLastNormView \quad = [r \in Replicas \mapsto 1]$
  $\land vViewChanges \qquad = [r \in Replicas \mapsto \{\}]$
  $\land vReplicaStatus \qquad = [r \in Replicas \mapsto StNormal]$
  $\land vSyncPoint \qquad\quad = [r \in Replicas \mapsto 0]$
  $\land vLateBuffer \qquad\quad = [r \in Replicas \mapsto \{\}]$
  $\land vTentativeSync \qquad = [r \in Replicas \mapsto 0]$
  $\land vSyncReps \qquad\quad = [r \in Replicas \mapsto \{\}]$
  $\land vCommitPoint \qquad = [r \in Replicas \mapsto 0]$
  $\land vCrashVector \qquad\; = [r \in Replicas \mapsto [rr \in Replicas \mapsto 0]]$
  $\land vCrashVectorReps = [r \in Replicas \mapsto \{\}]$
  $\land vRecoveryReps \qquad = [r \in Replicas \mapsto \{\}]$
  $\land vUUIDCounter \qquad = [c \in Replicas \mapsto 0]$

**Client State**

VARIABLES     $vClientClock,$      Current Clock Time of the client
              $vClientReqNum$    The number of requests that have been sent by this client

$InitClientState \triangleq$
  $\land vClientClock \qquad = [c \in Clients \mapsto 1]$
  $\land vClientReqNum \quad = [c \in Clients \mapsto 0]$

$clientVars \qquad\qquad \triangleq \langle vClientClock, vClientReqNum\rangle$

**Set of all vars**
$vars \triangleq \langle networkVars, replicaVars, clientVars\rangle$

\ * **Initial state**
$Init \triangleq \land InitNetworkState$
  $\qquad \land InitReplicaState$
  $\qquad \land InitClientState$
  $\qquad \land vReplicaProcessed = [r \in Replicas \mapsto \{\}]$
  $\qquad \land vClientProcessed = [c \in Clients \mapsto \{\}]$
  $\qquad \land DebugAction = \langle \text{"Init"}, \text{""}\rangle$

---

## Helpers

$NumofReplicas(status) \triangleq Cardinality(\{r \in Replicas \quad : vReplicaStatus[r] = status\})$

$DuplicateRep(ReplySet, m) \triangleq m.sender \in \{mm.sender : mm \in ReplySet\}$

$Pick(S) \triangleq$ CHOOSE $s \in S :$ TRUE

Convert a Set to Sequence
RECURSIVE $Set2Seq(\_)$
$Set2Seq(S) \triangleq$ IF $Cardinality(S) = 0$ THEN $\langle\rangle$
        ELSE
        LET
          $x \triangleq$ CHOOSE $x \in S :$ TRUE
        IN
          $\langle x \rangle \circ Set2Seq(S \setminus \{x\})$

Convert a Sequence to Set
$Seq2Set(seq) \triangleq \{seq[i] : i \in$ DOMAIN $seq\}$

$Max(S) \triangleq$ CHOOSE $x \in S : \forall\, y \in S : x \geq y$

$Min(S) \triangleq$ CHOOSE $x \in S : \forall\, y \in S : x \leq y$

**View ID Helpers**
$LeaderID(viewID) \triangleq (viewID \% Len(ReplicaOrder)) + ($IF $viewID \geq Len(ReplicaOrder)$ THEN 1 ELSE 0$)$

$Leader(viewID) \triangleq ReplicaOrder[LeaderID(viewID)]$   remember $\langle\rangle$ are 1-indexed

**Log Manipulation Helpers**

The order of 2 *log* entries are decided by the tuple $< deadline,\ clientID,\ requestID >$
Usually, deadline makes the two entries comparable
When 2 different entries have the same deadline, the tie is broken with *clientID*
Further, the tie is broken is *requestID*
(unnecessary if we only allow client to submit one request at one tick)
$EntryLeq(l1,\ l2)$        $\triangleq\ \wedge\ l1.deadline \leq l2.deadline$
                              $\wedge\ l1.clientID \leq l2.clientID$
                              $\wedge\ l1.requestID \leq l2.requestID$

$EntryEq(l1,\ l2)$          $\triangleq\ \wedge\ l1.deadline = l2.deadline$
                              $\wedge\ l1.clientID = l2.clientID$
                              $\wedge\ l1.requestID = l2.requestID$

$EntryLessThan(l1,\ l2) \triangleq\ \wedge\ EntryLeq(l1,\ l2)$
                                  $\wedge\ \neg(EntryEq(l1,\ l2))$

Find entry in one replica's *log* ( $< clientID,\ reqID >$ can uniquely identify the *log* entry)
We do not check deadline, because the leader may have modified the request's deadline
Return 0 when we fail to find it (remember Sequence is 1-indexed in TLA+, so 0 can serve as a dummy value)
$FindEntry(clientID,\ reqID,\ log) \triangleq$
                LET
                    $entryIndexSet \triangleq \{i \in 1\,..\,Len(log) : \wedge\ log[i].clientID = clientID$
                                                  $\wedge\ log[i].reqID = reqID\}$

IN
　IF $Cardinality(entryIndexSet) = 0$ THEN
　　$0$
　ELSE
　　$Pick(entryIndexSet)$

$SortLogSeq(seq) \triangleq SortSeq(seq, \text{LAMBDA } x, y : EntryLessThan(x, y))$

Given a set of logs, return the sorted *log* list
$GetSortLogSeq(S) \triangleq$ LET
　　　$seq \triangleq Set2Seq(S)$
　IN
　　$SortLogSeq(seq)$

Merge logs, first put all *log* items together, deduplicated (*i.e.* UNION them into a set). Then, do filtering and only keep those that have appeared in at least $\lceil f/2 \rceil + 1 replicas$.
$CountVotes(logll, x) \triangleq Cardinality(\{logSet \in logll : x \in logSet\})$

$MergeUnSyncLogs(unSyncedLogs, lastSyncedLog) \triangleq$
　　LET
　　　$unSyncedLogSet \triangleq \text{UNION } unSyncedLogs$
　　　$votedLogSet \triangleq \{x \in unSyncedLogSet :$
　　　　　　$\land EntryLessThan(lastSyncedLog, x)$
　　　　　　$\land CountVotes(unSyncedLogs, x) \geq RecoveryQuorumSize\}$
　　IN
　　　$GetSortLogSeq(votedLogSet)$

**Network Helpers**
Add a message to the network
$Send(ms) \triangleq messages' = messages \cup ms$

Convert the request format to a *log* format (by summing up $s$ and $l$ to get deadline)
$Req2Log(req) \triangleq [\quad mtype \quad\quad \mapsto MLogEntry,$
　　　　　$deadline \quad \mapsto req.s + req.l,$
　　　　　$clientID \quad \mapsto req.sender,$
　　　　　$requestID \quad \mapsto req.requestID,$
　　　　　$command \quad \mapsto req.command$
　　　　$]$

Index does not need to include command field, which is the body of the request/*log*, and can be very large
$GetLogIndex(entry) \triangleq [$
　　　　　$mtype \quad\quad \mapsto MLogIndex,$
　　　　　$deadline \quad \mapsto entry.deadline,$
　　　　　$clientID \quad \mapsto entry.clientID,$
　　　　　$requestID \quad \mapsto entry.requestID$
　　　　$]$

$GetLogIndexFromReply(reply) \triangleq [$
$\qquad\qquad\qquad mtype \qquad \mapsto MLogIndex,$
$\qquad\qquad\qquad deadline \qquad \mapsto reply.deadline,$
$\qquad\qquad\qquad clientID \qquad \mapsto reply.dest,$
$\qquad\qquad\qquad requestID \quad \mapsto reply.requestID$
$\qquad\qquad ]$

$IndexEq(index, msg) \triangleq \land index.deadline = msg.deadline$
$\qquad\qquad\qquad\qquad\quad \land index.clientID = msg.clientID$
$\qquad\qquad\qquad\qquad\quad \land index.requestID = msg.requestID$

Add local time to the message (for easy debug)
$Msg2RLog(msg, r) \triangleq msg @@ [tl \mapsto vReplicaClock[r]]$

$LastLog(logList) \qquad \triangleq \text{IF } Len(logList) = 0 \text{ THEN } NULLLog \text{ ELSE } logList[Len(logList)]$

$MergeCrashVector(cv1, cv2) \triangleq [r \in Replicas \mapsto Max(\{cv1[r], cv2[r]\})]$

$CheckCrashVector(m, r) \triangleq$
$\quad \text{IF } m.cv[m.sender] < vCrashVector[r][m.sender] \text{ THEN}$
$\qquad \text{FALSE }$ Potential stray message
$\quad \text{ELSE}$
$\qquad vCrashVector' = [vCrashVector \text{ EXCEPT } ![r] = MergeCrashVector(m.cv, vCrashVector[r])]$

$FilterStrayMessage(MSet, cv) \triangleq \{m \in MSet : m.cv[m.sender] \geq cv[m.sender]\}$

---

## Message Handlers and Actions

**Client action**
Client $c$ sends a request
We assume client can only send one request in one tick of time
If time has reached the bound, this client cannot send request any more

$ClientSendRequest(c) \triangleq \qquad \land vClientClock[c] < MaxTime$
$\qquad\qquad\qquad\qquad\qquad\quad \land vClientReqNum[c] < MaxReqNum$
$\qquad\qquad\qquad\qquad\qquad\quad \land Send(\{[mtype \mapsto MClientRequest,$
$\qquad\qquad\qquad\qquad\qquad\qquad\quad sender \qquad\quad \mapsto c,$ clientID
$\qquad\qquad\qquad\qquad\qquad\qquad\quad requestID \quad\quad \mapsto vClientReqNum[c] + 1,$ requestID
$\qquad\qquad\qquad\qquad\qquad\qquad\quad command \qquad\; \mapsto \text{""},$
$\qquad\qquad\qquad\qquad\qquad\qquad\quad s \qquad\qquad\qquad \mapsto vClientClock[c],$ submission time
$\qquad\qquad\qquad\qquad\qquad\qquad\quad l \qquad\qquad\qquad \mapsto LatencyBounds[c],$ latency bound
$\qquad\qquad\qquad\qquad\qquad\qquad\quad dest \qquad\qquad \mapsto r$
$\qquad\qquad\qquad\qquad\qquad\qquad\quad ] : r \in Replicas\})$
$\qquad\qquad\qquad\qquad\qquad\quad \land vClientClock' = [vClientClock \text{ EXCEPT } ![c] = vClientClock[c] + 1]$

10

$$\land vClientReqNum' = [vClientReqNum \text{ EXCEPT } ![c] = vClientReqNum[c] + 1]$$
$$\land \text{UNCHANGED } \langle replicaVars \rangle$$

$Duplicate(entry,\ logSet) \triangleq$
  LET
      $findSet \triangleq \{x \in logSet : \land x.clientID = entry.clientID$
                                $\land x.requestID = entry.requestID\}$
  IN
      $Cardinality(findSet) > 0$

$HandleClientRequest(r,\ m) \triangleq$
  LET
    $mlog \triangleq Req2Log(m)$
  IN
    If the request is duplicate, it will no longer be appended to the $log$
    Replicas simply reply the previous execution result of this request
    (we do not model execution in this spec)
    $\land \neg Duplicate(mlog,\ Seq2Set(vLog[r]) \cup vEarlyBuffer[r])$
    $\land vReplicaStatus[r] = StNormal$
       The request can enter the early buffer
    $\land \lor \land EntryLessThan(LastLog(vLog[r]),\ mlog)$
         $\land vEarlyBuffer' = [$
             $vEarlyBuffer \text{ EXCEPT } ![r] = vEarlyBuffer[r] \cup \{mlog\}$
          $]$
         $\land \text{UNCHANGED } \langle networkVars,\ clientVars,$
                        $vLog,\ vViewID,\ vReplicaClock,$
                        $vLastNormView,\ vViewChanges,\ vReplicaStatus,$
                        $vSyncPoint,\ vLateBuffer,$
                        $vTentativeSync,\ vSyncReps,\ vCommitPoint,$
                        $vUUIDCounter,\ vCrashVector,$
                        $vCrashVectorReps,\ vRecoveryReps \rangle$
      (1) Followers' early buffers do not accept the request
         if its deadline is smaller than previously appended (last released) entry,
         so followers directly put the request into the late buffer
      (2) Leader modifies its deadline to be larger than the last released entry
         so as to make it eligible for entering the early buffer
      $\lor \land EntryLessThan(mlog,\ LastLog(vLog[r]))$
         $\land \text{IF } r = Leader(vViewID[r]) \text{ THEN }$ this replica is the leader in the current view
              $\land vEarlyBuffer' = [$
                  $vEarlyBuffer \text{ EXCEPT } ![r] = vEarlyBuffer[r] \cup \{[$
                    $mtype \quad\quad \mapsto MLogEntry,$
                    $clientID \quad \mapsto mlog.clientID,$
                    $requestID \quad \mapsto mlog.requestID,$

$$
\begin{aligned}
&\qquad\qquad\qquad\quad deadline \quad\ \mapsto LastLog(vLog[r]).deadline + 1,\\
&\qquad\qquad\qquad\quad command \ \mapsto mlog.command\\
&\qquad\qquad\quad ]\}\\
&\qquad\quad ]
\end{aligned}
$$

$\land$ UNCHANGED $\langle networkVars, clientVars,$
$\qquad\qquad\qquad\qquad vLog, vViewID, vReplicaClock,$
$\qquad\qquad\qquad\qquad vLastNormView, vViewChanges, vReplicaStatus,$
$\qquad\qquad\qquad\qquad vSyncPoint, vLateBuffer,$
$\qquad\qquad\qquad\qquad vTentativeSync, vSyncReps, vCommitPoint,$
$\qquad\qquad\qquad\qquad vUUIDCounter, vCrashVector,$
$\qquad\qquad\qquad\qquad vCrashVectorReps, vRecoveryReps\rangle$

ELSE     this replica is a follower in the current view

$\qquad\quad \land vLateBuffer' = [$
$\qquad\qquad\qquad vLateBuffer$ EXCEPT $![r] = vLateBuffer[r] \cup \{mlog\}$
$\qquad\quad ]$

$\qquad\quad \land$ UNCHANGED $\langle networkVars, clientVars,$
$\qquad\qquad\qquad\qquad\qquad vLog, vEarlyBuffer, vViewID, vReplicaClock,$
$\qquad\qquad\qquad\qquad\qquad vLastNormView, vViewChanges, vReplicaStatus,$
$\qquad\qquad\qquad\qquad\qquad vSyncPoint, vTentativeSync,$
$\qquad\qquad\qquad\qquad\qquad vSyncReps, vCommitPoint,$
$\qquad\qquad\qquad\qquad\qquad vUUIDCounter, vCrashVector,$
$\qquad\qquad\qquad\qquad\qquad vCrashVectorReps, vRecoveryReps\rangle$

Release relevant requests from *vEarlyBuffer* and append to *vLog*,
and then send a fast reply

$FlushEarlyBuffer(r) \triangleq$

LET
$\quad validLogSet \triangleq \{x \in vEarlyBuffer[r] :$
$\qquad\qquad\qquad\qquad\quad \land x.deadline < vReplicaClock[r]$   < rather than $\leq$
$\qquad\qquad\qquad\qquad\quad \land EntryLessThan(LastLog(vLog[r]), x)\}$
$\quad validLogs \triangleq GetSortLogSeq(validLogSet)$
$\quad newLogStart \triangleq Len(vLog[r]) + 1$

IN
$\land\ vLog' = [vLog$ EXCEPT $![r] = vLog[r] \circ validLogs]$
$\land\ vEarlyBuffer' = [vEarlyBuffer$ EXCEPT $![r]$
$\qquad\qquad\qquad = \{x \in vEarlyBuffer[r] : x.deadline \geq vReplicaClock[r]\}]$   $\geq$ rather than >
$\land\ Send(\{[mtype \qquad\ \mapsto MFastReply,$
$\qquad\quad sender \qquad\ \mapsto r,$
$\qquad\quad dest \qquad\qquad \mapsto vLog'[r][i].clientID,$
$\qquad\quad viewID \qquad\ \mapsto vViewID[r],$
$\qquad\quad requestID \quad\ \mapsto vLog'[r][i].requestID,$
$\qquad\quad hash \qquad\qquad \mapsto [$
$\qquad\qquad\qquad\qquad log \mapsto SubSeq(vLog'[r], 1, i),$
$\qquad\qquad\qquad\qquad cv \mapsto vCrashVector$

$$
\begin{aligned}
&\qquad\qquad\qquad ],\\
&\qquad\qquad deadline \quad \mapsto vLog'[r][i].deadline,\\
&\qquad\qquad logSlotNum \mapsto i\\
&\qquad\qquad ] : i \in newLogStart \,..\, Len(vLog'[r])\})\\
&\quad\wedge\ \text{IF } r = Leader(vViewID[r]) \quad \text{THEN}\\
&\qquad\qquad \wedge vSyncPoint' = [vSyncPoint \text{ EXCEPT } ![r] = Len(vLog'[r])]\\
&\qquad\qquad \wedge \text{UNCHANGED}\quad \langle\ clientVars,\ vViewID,\ vLastNormView,\ vViewChanges,\\
&\qquad\qquad\qquad\qquad\qquad\qquad vReplicaStatus,\ vReplicaClock,\ vLateBuffer,\\
&\qquad\qquad\qquad\qquad\qquad\qquad vTentativeSync,\ vSyncReps,\ vCommitPoint,\\
&\qquad\qquad\qquad\qquad\qquad\qquad vUUIDCounter,\ vCrashVector,\\
&\qquad\qquad\qquad\qquad\qquad\qquad vCrashVectorReps,\ vRecoveryReps\rangle\\
&\qquad \text{ELSE}\\
&\qquad\qquad \text{UNCHANGED}\quad \langle clientVars,\ vViewID,\ vLastNormView,\ vViewChanges,\\
&\qquad\qquad\qquad\qquad\qquad\quad vReplicaStatus,\ vReplicaClock,\\
&\qquad\qquad\qquad\qquad\qquad\quad vSyncPoint,\ vLateBuffer,\\
&\qquad\qquad\qquad\qquad\qquad\quad vTentativeSync,\ vSyncReps,\ vCommitPoint,\\
&\qquad\qquad\qquad\qquad\qquad\quad vUUIDCounter,\ vCrashVector,\\
&\qquad\qquad\qquad\qquad\qquad\quad vCrashVectorReps,\ vRecoveryReps\ \rangle
\end{aligned}
$$

Clock can be random value ($RandomElement(1 \,..\, MaxTime)$),

because clock sync algorithm can give negative offset, or even fails

But *Nezha* depend on clock for performance but not for correctness

If the replica clock goes beyond $MaxTime$, it will stop processing

Since Clock is moved, then replicas can release relevant requests and append to logs

$$
\begin{aligned}
ReplicaClockMove(r) \ \triangleq\ &\wedge \text{IF } vReplicaClock[r] < MaxTime \text{ THEN}\\
&\qquad vReplicaClock' = [\\
&\qquad\qquad vReplicaClock \text{ EXCEPT } ![r] = RandomElement(1 \,..\, MaxTime)\\
&\qquad ]\\
&\quad \text{ELSE}\\
&\qquad \text{UNCHANGED } vReplicaClock\\
&\wedge \text{UNCHANGED } \langle networkVars,\ clientVars,\\
&\qquad\qquad\qquad\qquad vLog,\ vEarlyBuffer,\ vViewID,\\
&\qquad\qquad\qquad\qquad vLastNormView,\ vViewChanges,\ vReplicaStatus,\\
&\qquad\qquad\qquad\qquad vSyncPoint,\ vLateBuffer,\ vTentativeSync,\\
&\qquad\qquad\qquad\qquad vSyncReps,\ vCommitPoint,\\
&\qquad\qquad\qquad\qquad vUUIDCounter,\ vCrashVector,\\
&\qquad\qquad\qquad\qquad vCrashVectorReps,\ vRecoveryReps\rangle
\end{aligned}
$$

Client clock move does not change any other things

$$
\begin{aligned}
ClientClockMove(c) \ \triangleq\ &\wedge \text{ IF } vClientClock[c] < MaxTime \text{ THEN}\\
&\qquad vClientClock' = [\\
&\qquad\qquad vClientClock \text{ EXCEPT } ![c] = RandomElement(1 \,..\, MaxTime)\\
&\qquad ]\\
&\quad \text{ELSE}\\
&\qquad \text{UNCHANGED } vClientClock\\
&\wedge \text{ UNCHANGED } \langle networkVars,\ replicaVars,\ vClientReqNum\rangle
\end{aligned}
$$

## Index Synchronization to Fix Set Inequality

$StartIndexSync(r) \triangleq$
  LET
    $indices \triangleq \{GetLogIndex(vLog[r][i]) : i \in 1 .. Len(vLog[r])\}$
  IN
  $\wedge\ r = Leader(vViewID[r])$
  $\wedge\ vReplicaStatus[r] = StNormal$
  $\wedge\ Cardinality(indices) > 0$   leader has *log* entries to sync
  $\wedge\ Send(\{[mtype \quad \mapsto MIndexSync,$
           $sender \quad \mapsto r,$
           $dest \quad \mapsto d,$
           $viewID \quad \mapsto vViewID[r],$
           $logindcies \mapsto indices] : d \in Replicas\})$
  $\wedge$ UNCHANGED $\langle clientVars, replicaVars \rangle$


$GetSyncLogs(logSeq, indices) \triangleq$
  LET
    $logSet \triangleq \{l \in Seq2Set(logSeq) : \exists\, index \in indices : EntryEq(index, l)\}$
  IN
    $GetSortLogSeq(logSet)$

$GetUnSyncLogs(logSeq, lastSyncedLog) \triangleq$
  LET
    $logSet \triangleq \{l \in Seq2Set(logSeq) : EntryLessThan(lastSyncedLog, l)\}$
  IN
    $GetSortLogSeq(logSet)$

$HandleIndexSync(r, m) \triangleq$
  $\wedge\ r \neq Leader(vViewID[r])$
  $\wedge\ vReplicaStatus[r] = StNormal$
  $\wedge\ m.viewID = vViewID[r]$
  $\wedge\ m.sender = Leader(vViewID[r])$
  $\wedge\ vSyncPoint[r] < Len(m.logindcies)$
  $\wedge$ LET
    $entries \triangleq \{vLog[r][i] : i \in 1 .. Len(vLog[r])\}$
    $indices \triangleq \{GetLogIndex(vLog[r][i]) : i \in 1 .. Len(vLog[r])\}$
    $missedEntries \triangleq m.indices \setminus indices$
   IN
     Missing some *log entries* $\rightarrow$ Send *MMissEntryRequest*
    IF $Cardinality(missedEntries) > 0$ THEN
      $\wedge\ Send(\{[\ mtype \quad \mapsto MMissEntryRequest,$

14

$$
\begin{aligned}
&sender \quad \mapsto r, \\
&dest \quad \mapsto d, \\
&viewID \quad \mapsto vViewID[r], \\
&miss \quad \mapsto missedEntries] : d \in (Replicas \setminus \{r\})\})
\end{aligned}
$$

$\land$ UNCHANGED $\langle vLog,\ vSyncPoint \rangle$

ELSE

   LET

$$
\begin{aligned}
syncLogs \quad &\triangleq\quad GetSyncLogs(vLog[r],\ indices) \\
unsyncLogs \quad &\triangleq\quad GetUnSyncLogs(vLog[r],\ LastLog(syncLogs))
\end{aligned}
$$

   IN

$\land\ vLog' = [vLog$ EXCEPT $![r] = syncLogs \circ unsyncLogs]$

$\land\ vSyncPoint' = [vSyncPoint$ EXCEPT $![r] = Len(syncLogs)]$

$\land\ Send(\{[\ mtype \quad \mapsto MSlowReply,$

$$
\begin{aligned}
&sender \quad \mapsto r, \\
&dest \quad \mapsto vLog'[r][i].clientID, \\
&viewID \quad \mapsto vViewID[r], \\
&requestID \quad \mapsto vLog'[r][i].requestID, \\
&logSlotNum \mapsto i] : i \in (1\ ..\ Len(syncLogs))\})
\end{aligned}
$$

$\land$ UNCHANGED $\langle clientVars,\ vEarlyBuffer,\ vViewID,\ vReplicaClock,$
$vLastNormView,\ vViewChanges,\ vReplicaStatus,$
$vLateBuffer,\ vTentativeSync,\ vSyncReps,\ vCommitPoint,$
$vUUIDCounter,\ vCrashVector,$
$vCrashVectorReps,\ vRecoveryReps \rangle$

$FindEntries(log,\ indices)\ \triangleq$
   $\{l \in Seq2Set(log)\ : \exists\, x \in indices : IndexEq(l,\ x)\}$

$HandleMissEntryRequest(r,\ m)\ \triangleq$

$\land\ m.viewID = vViewID[r]$

$\land$ LET

   $findentries\ \triangleq\ FindEntries(vLog[r],\ m.miss)$

  IN

$\land\ Cardinality(findentries) > 0$

$\land\ Send(\{[\ mtype \quad \mapsto MMissEntryReply,$

$$
\begin{aligned}
&sender \quad \mapsto r, \\
&dest \quad \mapsto m.sender, \\
&viewID \quad \mapsto vViewID[r], \\
&entries \quad \mapsto findentries]\})
\end{aligned}
$$

$\land$ UNCHANGED $\langle clientVars,\ replicaVars \rangle$

$HandleMissEntryReply(r, m) \triangleq$
 $\wedge\ m.viewID = vViewID[r]$
 $\wedge$ LET
  $mergedSet \triangleq Seq2Set(vLog[r]) \cup m.entries$
  IN
  $vLog' = [vLog$ EXCEPT $![r] = GetSortLogSeq(mergedSet)]$
 $\wedge$ UNCHANGED $\langle networkVars, clientVars,$
        $vEarlyBuffer, vViewID, vReplicaClock,$
        $vLastNormView, vViewChanges, vReplicaStatus,$
        $vSyncPoint, vLateBuffer,$
        $vTentativeSync, vSyncReps, vCommitPoint,$
        $vUUIDCounter, vCrashVector,$
        $vCrashVectorReps, vRecoveryReps\rangle$

---

## Replica Rejoin

Failed replica loses all states

$StartReplicaFail(r) \triangleq$
 $\wedge\ NumofReplicas(StRecovering) < F$   We assume at most $F$ replicas can fail at the same time
 $\wedge\ vReplicaStatus' = [vReplicaStatus$ EXCEPT $![r] = StRecovering]$
 $\wedge\ vLog' = [vLog$ EXCEPT $![r] = \langle\rangle]$
 $\wedge\ vEarlyBuffer' = [vEarlyBuffer$ EXCEPT $![r] = \{\}]$
 $\wedge\ vViewID' = [vViewID$ EXCEPT $![r] = 1]$
 $\wedge\ vLastNormView' = [vLastNormView$ EXCEPT $![r] = 1]$
 $\wedge\ vViewChanges' = [vViewChanges$ EXCEPT $![r] = \{\}]$
 $\wedge\ vSyncPoint' = [vSyncPoint$ EXCEPT $![r] = 0]$
 $\wedge\ vLateBuffer' = [vLateBuffer$ EXCEPT $![r] = \{\}]$
 $\wedge\ vTentativeSync' = [vTentativeSync$ EXCEPT $![r] = 0]$
 $\wedge\ vSyncReps' = [vSyncReps$ EXCEPT $![r] = \{\}]$
 $\wedge\ vCommitPoint' = [vCommitPoint$ EXCEPT $![r] = 0]$
 $\wedge\ vCrashVector' = [vCrashVector$ EXCEPT $![r] = [rr \in Replicas \mapsto 0]]$
 $\wedge\ vCrashVectorReps' = [vCrashVectorReps$ EXCEPT $![r] = \{\}]$
 $\wedge\ vRecoveryReps' = [vRecoveryReps$ EXCEPT $![r] = \{\}]$
 $\wedge$ UNCHANGED $\langle vReplicaClock, vUUIDCounter, clientVars, networkVars\rangle$

Recovering replica starts recovery (by first sending $CrashVectorReq$)

$StartReplicaRecovery(r) \triangleq$
 $\wedge\ vReplicaStatus[r] = StRecovering$
 $\wedge\ vUUIDCounter' = [vUUIDCounter$ EXCEPT $![r] = vUUIDCounter[r] + 1]$
 $\wedge\ Send(\{[mtype \mapsto MCrashVectorReq,$
     $sender \mapsto r,$
     $dest \mapsto d,$

$$nonce \ \mapsto vUUIDCounter'[r]] : d \in Replicas\})$$
$$\land \text{UNCHANGED} \ \langle vLog, \ vEarlyBuffer, \ vViewID, \ vReplicaClock,$$
$$vLastNormView, \ vViewChanges, \ vReplicaStatus,$$
$$vSyncPoint, \ vLateBuffer,$$
$$vTentativeSync, \ vSyncReps, \ vCommitPoint,$$
$$vCrashVector, \ vCrashVectorReps, \ vRecoveryReps,$$
$$clientVars \ \rangle$$

$HandleCrashVectorReq(r, \ m) \ \triangleq$
$\quad \land vReplicaStatus[r] = StNormal$
$\quad \land Send(\{[mtype \ \mapsto MCrashVectorRep,$
$\qquad\qquad sender \ \mapsto r,$
$\qquad\qquad dest \quad \mapsto m.sender,$
$\qquad\qquad nonce \ \mapsto m.nonce,$
$\qquad\qquad cv \qquad \mapsto vCrashVector[r]]\})$
$\quad \land \text{UNCHANGED} \ \langle replicaVars, \ clientVars\rangle$

$HandleCrashVectorRep(r, \ m) \ \triangleq$
$\quad \land vReplicaStatus[r] = StRecovering$
$\quad \land vUUIDCounter[r] = m.nonce$
$\quad \land Cardinality(vCrashVectorReps[r]) \leq F$
$\quad \land \neg DuplicateRep(vCrashVectorReps[r], \ m)$
$\quad \land vCrashVectorReps' = [vCrashVectorReps \ \text{EXCEPT} \ ![r] = vCrashVectorReps[r] \cup \{m\}]$
$\quad \land vCrashVector' = [vCrashVector \ \text{EXCEPT} \ ![r] = MergeCrashVector(vCrashVector[r], \ m.cv)]$
$\quad \land \text{IF} \ Cardinality(vCrashVectorReps') = F + 1 \ \text{THEN}$ <span style="background:#cccccc">got enough replies and can settle down $cv$</span>
$\qquad Send(\{[mtype \ \mapsto MRecoveryReq,$
$\qquad\qquad sender \ \mapsto r,$
$\qquad\qquad dest \quad \mapsto d,$
$\qquad\qquad nonce \ \mapsto m.nonce,$
$\qquad\qquad cv \qquad \mapsto vCrashVector'[r]] : d \in Replicas\})$
$\qquad \text{ELSE}$
$\qquad \text{UNCHANGED} \ \langle networkVars\rangle$

$\quad \land \text{UNCHANGED} \ \langle vLog, \ vEarlyBuffer, \ vViewID, \ vReplicaClock,$
$\qquad\qquad vLastNormView, \ vViewChanges, \ vReplicaStatus,$
$\qquad\qquad vSyncPoint, \ vLateBuffer,$
$\qquad\qquad vTentativeSync, \ vSyncReps, \ vCommitPoint,$
$\qquad\qquad vUUIDCounter, \ vRecoveryReps,$
$\qquad\qquad clientVars\rangle$

$HandleRecoveryReq(r, \ m) \ \triangleq$
$\quad \land vReplicaStatus[r] = StNormal$
$\quad \land vCrashVector' = [vCrashVector \ \text{EXCEPT} \ ![r] = MergeCrashVector(vCrashVector[r], \ m.cv)]$

$$\wedge \; Send(\{[\; mtype \;\mapsto MRecoveryRep,$$
$$sender \;\mapsto r,$$
$$dest \quad \mapsto m.sender,$$
$$viewID \mapsto vViewID[r],$$
$$cv \qquad \mapsto vCrashVector'[r]] : d \in Replicas\})$$

$\wedge$ UNCHANGED $\langle vLog,\; vEarlyBuffer,\; vViewID,\; vReplicaClock,$
$\qquad\qquad\qquad vLastNormView,\; vViewChanges,\; vReplicaStatus,$
$\qquad\qquad\qquad vSyncPoint,\; vLateBuffer,$
$\qquad\qquad\qquad vTentativeSync,\; vSyncReps,\; vCommitPoint,$
$\qquad\qquad\qquad vUUIDCounter,\; vCrashVectorReps,\; vRecoveryReps,$
$\qquad\qquad\qquad clientVars \quad \rangle$

$HandleRecoveryRep(r,\; m) \;\triangleq$
$\quad \wedge\; vReplicaStatus[r] = StRecovering$
$\quad \wedge\; Cardinality(vRecoveryReps[r]) \leq F$
$\quad \wedge\; \neg DuplicateRep(vRecoveryReps[r],\; m.sender)$
$\quad \wedge\; CheckCrashVector(m,\; r)$

Note: After crash vector is updated, those previously accepted messages may also become stray message. Those messages should also be filtered out.

$\quad \wedge\; vRecoveryReps' = [vRecoveryReps$ EXCEPT
$\qquad\qquad\qquad\qquad ![r] = FilterStrayMessage(vRecoveryReps[r] \cup \{m\},\; vCrashVector'[r])\; ]$

$\quad \wedge$ IF $Cardinality(vRecoveryReps') = F + 1$ THEN  got enough replies
$\qquad$ LET
$\qquad\qquad newView \;\triangleq\; Max(\{mm.viewID : mm \in vRecoveryReps'[r]\})$
$\qquad\qquad leaderId \;\triangleq\; newView\%Cardinality(Replicas)$
$\qquad$ IN
$\qquad\qquad Send(\{[mtype \;\mapsto MStateTransferReq,$
$\qquad\qquad\qquad sender \mapsto r,$
$\qquad\qquad\qquad dest \quad \mapsto leaderId,$
$\qquad\qquad\qquad cv \qquad \mapsto vCrashVector'[r]] : d \in Replicas\})$
$\qquad$ ELSE
$\qquad$ UNCHANGED $\langle networkVars \rangle$

$\quad \wedge$ UNCHANGED $\langle vLog,\; vEarlyBuffer,\; vViewID,\; vReplicaClock,$
$\qquad\qquad\qquad\qquad vLastNormView,\; vViewChanges,\; vReplicaStatus,$
$\qquad\qquad\qquad\qquad vSyncPoint,\; vLateBuffer,$
$\qquad\qquad\qquad\qquad vTentativeSync,\; vSyncReps,\; vCommitPoint,$
$\qquad\qquad\qquad\qquad vUUIDCounter,\; vCrashVectorReps,$
$\qquad\qquad\qquad\qquad clientVars \rangle$

$HandleStateTransferReq(r,\; m) \;\triangleq$
$\quad \wedge\; vReplicaStatus[r] = StNormal$

$\land$ *CheckCrashVector*$(m, r)$
$\land$ *Send*$(\{[$ *mtype* $\mapsto$ *MStateTransferRep*,
          *sender* $\mapsto r$,
          *dest* $\mapsto m.sender$,
          *log* $\mapsto vLog[r]$,
          *sp* $\mapsto vSyncPoint[r]$,
          *cp* $\mapsto vCommitPoint[r]$,
          *cv* $\mapsto vCrashVector'[r]]\})$
$\land$ UNCHANGED $\langle vLog, vEarlyBuffer, vViewID, vReplicaClock,$
            $vLastNormView, vViewChanges, vReplicaStatus,$
            $vSyncPoint, vLateBuffer,$
            $vTentativeSync, vSyncReps, vCommitPoint,$
            $vUUIDCounter, vCrashVectorReps, vRecoveryReps,$
            $clientVars\rangle$

$HandleStateTransferRep(r, m) \triangleq$
    $\land vReplicaStatus[r] = StRecovering$
    $\land CheckCrashVector(m, r)$
    $\land vLog' = [vLog$ EXCEPT $![r] = m.log]$
    $\land vSyncPoint' = [vSyncPoint$ EXCEPT $![r] = m.sp]$
    $\land vCommitPoint' = [vCommitPoint$ EXCEPT $![r] = m.cp]$
    $\land vViewID' = [vViewID$ EXCEPT $![r] = m.viewID]$
    $\land vEarlyBuffer' = [vEarlyBuffer$ EXCEPT $![r] = \{\}]$
    $\land vLastNormView' = [vLastNormView$ EXCEPT $![r] = m.viewID]$
    $\land vViewChanges' = [vViewChanges$ EXCEPT $![r] = \{\}]$
    $\land vReplicaStatus' = [vReplicaStatus$ EXCEPT $![r] = StNormal]$
    $\land vLateBuffer' = [vLateBuffer$ EXCEPT $![r] = \{\}]$
    $\land vTentativeSync' = [vTentativeSync$ EXCEPT $![r] = m.sp]$
    $\land vSyncReps' = [vSyncReps$ EXCEPT $![r] = \{\}]$
    $\land vCrashVectorReps' = [vCrashVectorReps$ EXCEPT $![r] = \{\}]$
    $\land vRecoveryReps' = [vRecoveryReps$ EXCEPT $![r] = \{\}]$
    $\land$ UNCHANGED $\langle vReplicaClock, vUUIDCounter, clientVars\rangle$

## Leader Change

Replica $r$ starts a *Leader* change
$StartLeaderChange(r) \triangleq$
  $\land Send(\{[mtype \mapsto MViewChangeReq,$
          *sender* $\mapsto r$,
          *dest* $\mapsto d$,
          *viewID* $\mapsto vViewID[r] + 1$,
          *cv* $\mapsto vCrashVector[r]] : d \in Replicas\})$
  $\land$ UNCHANGED $\langle replicaVars, clientVars\rangle$

**View Change Handlers**

$HandleViewChangeReq(r, m) \triangleq$

  LET

    $currentViewID \quad \triangleq \quad vViewID[r]$

    $newViewID \qquad \triangleq \quad Max(\{currentViewID, m.viewID\})$

    $newLeaderNum \quad \triangleq \quad LeaderID(newViewID)$

  IN

    Recovering replica does not participate in view change

    $\wedge\ vReplicaStatus[r] \neq StRecovering$

    $\wedge\ currentViewID \quad \neq newViewID$

    $\wedge\ CheckCrashVector(m, r)$

    $\wedge\ vReplicaStatus' = [vReplicaStatus \ \text{EXCEPT} \ ![r] = StViewChange]$

    $\wedge\ vViewID' \qquad = [vViewID \ \text{EXCEPT} \ ![r] = newViewID]$

    $\wedge\ vViewChanges' = [vViewChanges \ \text{EXCEPT} \ ![r] = \{\}]$

    $\wedge\ Send(\{[mtype \qquad\quad \mapsto MViewChange,$

                 $dest \qquad\qquad \mapsto Leader(newViewID),$

                 $sender \qquad\quad\ \mapsto r,$

                 $viewID \qquad\quad \mapsto newViewID,$

                 $lastNormal \quad\ \mapsto vLastNormView[r],$

                 $syncedLog \quad\ \mapsto SubSeq(vLog[r], 1, vSyncPoint[r]),$

                 $unsyncedLog \mapsto SubSeq(vLog[r], vSyncPoint[r] + 1, Len(vLog[r])),$

                 $cv \qquad\qquad\ \mapsto vCrashVector[r]]\} \ \cup$

              Send the *MViewChangeReqs* in case this is an entirely new view

            $\{[mtype \ \mapsto MViewChangeReq,$

              $sender \mapsto r,$

              $dest \quad\ \mapsto d,$

              $viewID \mapsto newViewID,$

              $cv \qquad \mapsto vCrashVector[r]] : d \in Replicas\})$

    $\wedge\ \text{UNCHANGED} \ \langle clientVars, vLog, vEarlyBuffer, vReplicaClock,$

                       $vLastNormView, vSyncPoint, vLateBuffer,$

                       $vTentativeSync, vSyncReps, vCommitPoint,$

                       $vUUIDCounter, vCrashVectorReps, vRecoveryReps\rangle$

 

  Replica $r$ receives *MViewChange*, $m$

$HandleViewChange(r, m) \triangleq$

    Recovering replica does not participate in view change

    $\wedge\ vReplicaStatus[r] \neq StRecovering$

    Add the message to the *log*

    $\wedge\ vViewID[r] \qquad\quad = m.viewID$

    $\wedge\ vReplicaStatus[r] \quad = StViewChange$

    This replica is the leader

    $\wedge\ Leader(vViewID[r]) = r$

    $\wedge\ CheckCrashVector(m, r)$

Note: Similar to *vRecoveryReps*, (potential) stray messages should be filtered out.

$\wedge\, vViewChanges' = [vViewChanges \text{ EXCEPT}$
$\qquad\qquad\qquad ![r] = FilterStrayMessage(vViewChanges[r] \cup \{m\}, vCrashVector'[r])]$

If there's enough replies, start the new view

$\wedge$ LET
$\quad isViewPromise(M) \;\triangleq\; \wedge\, \{n.sender : n \in M\} \in Quorums$
$\qquad\qquad\qquad\qquad\quad \wedge\, \exists\, n \in M \;\; : n.sender = r$
$\quad vCMs \qquad\qquad\;\; \triangleq\; \{n \in vViewChanges'[r] :$
$\qquad\qquad\qquad\qquad\quad \wedge\, n.mtype \;\; = MViewChange$
$\qquad\qquad\qquad\qquad\quad \wedge\, n.viewID = vViewID[r]\}$

Create the state for the new view
$\quad normalViews \;\;\triangleq\; \{n.lastNormal : n \in vCMs\}$

Choose the largest normal view (*i.e.* the newest)
$\quad lastNormal \qquad\;\; \triangleq\; (\text{CHOOSE } v \in normalViews : \forall\, v2 \in normalViews : v2 \leq v)$

For logs before *vSyncPoint* (*i.e.* *syncedLog*), we directly copy from the *bestCandiates*

For *unsyncedLog*, we do quorum check to decide which ones should be added to recovery *Log*
$\quad goodCandidates \;\triangleq\;\; \{o \in vCMs : o.lastNormal = lastNormal\}$

*bestCandidate* can only be picked from *goodCandidates*,

because previous views may include invalid logs
$\quad bestCandidate \quad\; \triangleq\; \text{CHOOSE } n \in goodCandidates :$
$\qquad\qquad\qquad\qquad\quad \forall\, y \in goodCandidates : Len(n.syncedLog) \geq Len(y.syncedLog)$
$\quad unSyncedLogs \;\;\triangleq\; \{Seq2Set(n.unsyncedLog) : n \in goodCandidates\}$

$\quad$ IN
$\qquad$ IF *isViewPromise*(vCMs) THEN
$\qquad\quad Send(\{[mtype \qquad\; \mapsto MStartView,$
$\qquad\qquad\quad dest \qquad\qquad \mapsto d,$
$\qquad\qquad\quad viewID \qquad\; \mapsto vViewID[r],$
$\qquad\qquad\quad log \qquad\qquad\;\; \mapsto bestCandidate.syncedLog$
$\qquad\qquad\qquad\qquad\qquad\qquad \circ MergeUnSyncLogs(unSyncedLogs, LastLog(bestCandidate.syncedLog))$
$\qquad\qquad\quad ] : d \in Replicas\})$
$\qquad$ ELSE
$\qquad\quad$ UNCHANGED *networkVars*
$\wedge$ UNCHANGED $\langle clientVars, \; vLog, \; vEarlyBuffer, vViewID, vReplicaClock,$
$\qquad\qquad\qquad\quad vLastNormView, vReplicaStatus, vSyncPoint, vLateBuffer,$
$\qquad\qquad\qquad\quad vTentativeSync, vSyncReps, vCommitPoint,$
$\qquad\qquad\qquad\quad vUUIDCounter, vCrashVectorReps, vRecoveryReps\rangle$


Replica $r$ receives a *MStartView*, $m$

$HandleStartView(r, m) \;\triangleq\;$
$\quad \wedge\, vReplicaStatus[r] \neq StRecovering$
$\quad \wedge\, \vee\, vViewID[r] \quad\; < m.viewID$
$\qquad \vee\, vViewID[r] \quad\; = m.viewID \wedge vReplicaStatus[r] = StViewChange$
$\quad \wedge\, CheckCrashVector(m, r)$

$\wedge\ vLog' \qquad\qquad = [vLog\ \text{EXCEPT}\ ![r] = m.log]$
$\wedge\ vReplicaStatus' \quad = [vReplicaStatus\ \text{EXCEPT}\ ![r] = StNormal]$
$\wedge\ vViewID' \qquad\quad = [vViewID\ \text{EXCEPT}\ ![r] = m.viewID]$
$\wedge\ vLastNormView' = [vLastNormView\ \text{EXCEPT}\ ![r] = m.viewID]$
$\wedge\ vEarlyBuffer' = [vEarlyBuffer\ \text{EXCEPT}\ ![r] = \{\}]$  clear Early Buffer for the new view
$\wedge\ vLateBuffer' = [vLateBuffer\ \text{EXCEPT}\ ![r] = \{\}]$  clear Late Buffer for the new view
$\wedge\ vSyncPoint' = [vSyncPoint\ \text{EXCEPT}\ ![r] = Len(m.log)]$
$\wedge\ vTentativeSync' = [vTentativeSync\ \text{EXCEPT}\ ![r] = Len(m.log)]$

Send replies (in the new view) for all *log* items

$\wedge\ \text{IF}\ r = Leader(m.viewID)\ \text{THEN}$      Leader only sends fast reply
$\qquad Send(\{[\ mtype \qquad \mapsto MFastReply,$
$\qquad\qquad\qquad sender \qquad \mapsto r,$
$\qquad\qquad\qquad dest \qquad\qquad \mapsto m.log[i].clientID,$
$\qquad\qquad\qquad viewID \qquad \mapsto m.viewID,$
$\qquad\qquad\qquad requestID \quad \mapsto m.log[i].requestID,$
$\qquad\qquad\qquad hash \qquad\qquad \mapsto [$
$\qquad\qquad\qquad\qquad\qquad\qquad log \mapsto SubSeq(m.log,\ 1,\ i),$
$\qquad\qquad\qquad\qquad\qquad\qquad cv\ \mapsto vCrashVector$
$\qquad\qquad\qquad\qquad\qquad ],$
$\qquad\qquad\qquad deadline \quad\ \mapsto m.log[i].deadline,$
$\qquad\qquad\qquad logSlotNum \mapsto i] : i \in (1\ ..\ Len(m.log))\})$
$\quad\ \text{ELSE}$    While staring view, followers knows the *log* is synced with the leader, so send slow-reply
$\qquad Send(\{[\ mtype \qquad \mapsto MSlowReply,$
$\qquad\qquad\qquad sender \qquad \mapsto r,$
$\qquad\qquad\qquad dest \qquad\qquad \mapsto m.log[i].clientID,$
$\qquad\qquad\qquad viewID \qquad \mapsto m.viewID,$
$\qquad\qquad\qquad requestID \quad\ \mapsto m.log[i].requestID,$
$\qquad\qquad\qquad logSlotNum \mapsto i] : i \in (1\ ..\ Len(m.log))\})$
$\wedge\ \text{UNCHANGED}\ \langle clientVars,\ vReplicaClock,\ vViewChanges,$
$\qquad\qquad\qquad\qquad vSyncReps,\ vCommitPoint,\ vCrashVector,$
$\qquad\qquad\qquad\qquad vUUIDCounter,\ vCrashVectorReps,\ vRecoveryReps\rangle$

---

## Periodic Synchronization

Leader replica $r$ conduct synchronization periodically

This periodic sync process is different from index sync process

It ensures that all replicas' logs are stable up to their *CommitPoint* (for fast recovery)

Our *CommitPoint* is essentially the *sync-point* defined in *NOPaxos* paper

Just as mentioned in *NOPaxos* paper, it is an optional optimization for fast recovery

*Nezha* still works even without this part

$StartSync(r) \triangleq$
$\quad \wedge\ Leader(vViewID[r]) = r$
$\quad \wedge\ vReplicaStatus[r] \qquad = StNormal$
$\quad \wedge\ vTentativeSync[r] < Len(vLog[r])$   If $\geq$ then no need to sync
$\quad \wedge\ vSyncReps' \qquad\qquad = [vSyncReps\ \text{EXCEPT}\ ![r] = \{\}]$

$$\wedge\ vTentativeSync' \quad = [vTentativeSync \text{ EXCEPT } ![r] = Len(vLog[r])]$$
$$\wedge\ Send(\{[mtype \qquad \mapsto MSyncPrepare,$$
$$\qquad\qquad sender \qquad \mapsto r,$$
$$\qquad\qquad dest \qquad \mapsto d,$$
$$\qquad\qquad viewID \qquad \mapsto vViewID[r],$$
$$\qquad\qquad log \qquad\quad \mapsto vLog[r]] : d \in Replicas\})$$
$$\wedge\ \text{UNCHANGED } \langle clientVars,\ vLog,\ vEarlyBuffer,\ vViewID,\ vReplicaClock,$$
$$\qquad\qquad\qquad vLastNormView,\ vViewChanges,\ vReplicaStatus,$$
$$\qquad\qquad\qquad vSyncPoint,\ vLateBuffer,\ vCommitPoint,$$
$$\qquad\qquad\qquad vUUIDCounter,\ vCrashVector,$$
$$\qquad\qquad\qquad vCrashVectorReps,\ vRecoveryReps\rangle$$

Replica $r$ receives $MSyncPrepare,\ m$

$HandleSyncPrepare(r,\ m) \triangleq$

LET
$$newLog \triangleq m.log \circ GetUnSyncLogs(vLog[r],\ LastLog(m.log))$$
IN
$$\wedge\ vReplicaStatus[r] = StNormal$$
$$\wedge\ m.viewID \qquad\quad = vViewID[r]$$
$$\wedge\ m.sender \qquad\quad = Leader(vViewID[r])$$
$$\wedge\ \text{IF} \quad vSyncPoint[r]\ < Len(m.log) \text{ THEN}$$
$$\qquad\quad \wedge\ vSyncPoint' = [vSyncPoint \text{ EXCEPT } ![r] = Len(m.log)]$$
$$\qquad\quad \wedge\ vLog' \qquad = [vLog \text{ EXCEPT } ![r] = newLog]$$
$$\qquad\quad \wedge\ Send(\{[\ mtype \qquad \mapsto MSlowReply,$$
$$\qquad\qquad\qquad sender \qquad \mapsto r,$$
$$\qquad\qquad\qquad dest \qquad\quad \mapsto m.log[i].clientID,$$
$$\qquad\qquad\qquad viewID \qquad \mapsto m.viewID,$$
$$\qquad\qquad\qquad requestID \quad \mapsto m.log[i].requestID,$$
$$\qquad\qquad\qquad logSlotNum \mapsto i] : i \in (1 .. Len(m.log))\})$$
$$\qquad \text{ELSE}$$
$$\qquad\qquad \text{UNCHANGED } \langle vLog,\ vSyncPoint\rangle$$
$$\wedge\ Send(\{[mtype \qquad\qquad \mapsto MSyncRep,$$
$$\qquad\quad sender \qquad\qquad \mapsto r,$$
$$\qquad\quad dest \qquad\qquad\quad \mapsto m.sender,$$
$$\qquad\quad viewID \qquad\qquad \mapsto vViewID[r],$$
$$\qquad\quad logSlotNumber \mapsto Len(m.log)]\}$$
$$\qquad )$$
$$\wedge\ \text{UNCHANGED } \langle clientVars,\ vEarlyBuffer,\ vViewID,\ \ vReplicaClock,$$
$$\qquad\qquad\qquad vLastNormView,\ vViewChanges,\ vReplicaStatus,$$
$$\qquad\qquad\qquad vLateBuffer,\ vTentativeSync,\ vSyncReps,\ vCommitPoint,$$
$$\qquad\qquad\qquad vUUIDCounter,\ vCrashVector,$$
$$\qquad\qquad\qquad vCrashVectorReps,\ vRecoveryReps\rangle$$

$HandleSyncRep(r, m) \triangleq$
$\quad \wedge m.viewID \qquad\quad = vViewID[r]$
$\quad \wedge vReplicaStatus[r] \;\; = StNormal$
$\quad \wedge vSyncReps' \qquad\quad = [vSyncReps \text{ EXCEPT } ![r] = vSyncReps[r] \cup \{m\}]$
$\quad \wedge \text{LET } isViewPromise(M) \;\triangleq\; \wedge \{n.sender : n \in M\} \in Quorums$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge \exists\, n \in M \; : n.sender = r$
$\qquad\quad sRMs \qquad\qquad\qquad \triangleq \{n \in vSyncReps'[r] :$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge n.mtype \qquad\quad = MSyncRep$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge n.viewID \qquad\quad = vViewID[r]$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge n.logSlotNumber = vTentativeSync[r]\}$
$\qquad\quad committedLog \qquad\quad \triangleq \text{ IF } vTentativeSync[r] \geq 1 \text{ THEN}$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad SubSeq(vLog[r], 1, vTentativeSync[r])$
$\qquad\qquad\qquad\qquad\qquad\qquad\quad \text{ELSE}$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \langle\rangle$
$\qquad \text{IN}$
$\qquad\quad \text{IF } isViewPromise(sRMs) \text{ THEN}$
$\qquad\qquad \wedge Send(\{[mtype \qquad\quad \mapsto MSyncCommit,$
$\qquad\qquad\qquad\qquad sender \qquad\quad \mapsto r,$
$\qquad\qquad\qquad\qquad dest \qquad\qquad \mapsto d,$
$\qquad\qquad\qquad\qquad viewID \qquad\quad \mapsto vViewID[r],$
$\qquad\qquad\qquad\qquad log \qquad\qquad\quad \mapsto committedLog] :$
$\qquad\qquad\qquad\qquad d \in Replicas\})$
$\qquad\qquad \wedge vCommitPoint' = [vCommitPoint \text{ EXCEPT } ![r] = vTentativeSync[r]]$
$\qquad\quad \text{ELSE}$
$\qquad\qquad \text{UNCHANGED } \langle networkVars, vCommitPoint\rangle$
$\quad \wedge \text{UNCHANGED } \langle clientVars, vLog, vEarlyBuffer, vViewID,$
$\qquad\qquad\qquad\qquad\qquad vReplicaClock, vLastNormView, vViewChanges,$
$\qquad\qquad\qquad\qquad\qquad vReplicaStatus, vSyncPoint, vLateBuffer,$
$\qquad\qquad\qquad\qquad\qquad vTentativeSync, vUUIDCounter, vCrashVector,$
$\qquad\qquad\qquad\qquad\qquad vCrashVectorReps, vRecoveryReps\rangle$

$HandleSyncCommit(r, m) \;\triangleq$
$\quad \text{LET}$
$\qquad newLog \;\triangleq\; m.log \circ GetUnSyncLogs(vLog[r], LastLog(m.log))$
$\quad \text{IN}$
$\quad \wedge vReplicaStatus[r] = StNormal$
$\quad \wedge m.viewID \qquad\quad = vViewID[r]$
$\quad \wedge m.sender \qquad\quad = Leader(vViewID[r])$
$\quad \wedge \text{IF } Len(m.log) \leq \; vCommitPoint[r] \text{ THEN}$
$\qquad\quad \text{UNCHANGED } \langle vCommitPoint, vLog\rangle$
$\qquad \text{ELSE}$
$\qquad\quad \wedge vLog' \qquad\quad = [vLog \text{ EXCEPT } ![r] = newLog]$

$\wedge\ vCommitPoint' = [vCommitPoint\ \text{EXCEPT}\ ![r] = Len(m.log)]$
$\wedge\ Send(\{[mtype \qquad \mapsto MSlowReply,$
$\qquad\qquad sender \qquad \mapsto r,$
$\qquad\qquad dest \qquad\quad \mapsto m.log[i].clientID,$
$\qquad\qquad viewID \qquad \mapsto m.viewID,$
$\qquad\qquad requestID \quad \mapsto m.log[i].requestID,$
$\qquad\qquad logSlotNum \mapsto i] : i \in (1\ ..\ Len(m.log))\})$
$\wedge\ \text{UNCHANGED}\ \langle networkVars,\ clientVars,\ vEarlyBuffer,$
$\qquad\qquad\qquad vViewID,\ vReplicaClock,\ vLastNormView, vViewChanges,$
$\qquad\qquad\qquad vReplicaStatus, vSyncPoint, vLateBuffer,$
$\qquad\qquad\qquad vTentativeSync, vSyncReps,$
$\qquad\qquad\qquad vUUIDCounter, vCrashVector,$
$\qquad\qquad\qquad vCrashVectorReps, vRecoveryReps\rangle$

---

## Invariants and Helper Functions

A request/*log* is committed in two possible cases:
(1) A fast quorum has sent either slow-reply messages, or fast-reply messages with consistent hashes [Fast Path]
(2) A simple quorum has sent slow-reply messages [Slow Path] Both quorums should include the leader

Check whether $log < clientID, requestID >$ is committed at position $logSlotNum$
$Committed(clientID, requestID, logSlotNum) \triangleq$

  Fast path
 $\vee\ \exists M \in \text{SUBSET}\ (\{m \in messages : \wedge\ \vee\ m.mtype = MFastReply$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \vee\ m.mtype = MSlowReply$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge\ m.logSlotNum = logSlotNum$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge\ m.dest = clientID$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge\ m.requestID = requestID\}) :$

  Sent from a fast quorum
  $\wedge\ \{m.sender : m \in M\} \in FastQuorums$
  Matching view-id
  $\wedge\ \exists m1 \in M : \forall m2 \in M : m1.viewID = m2.viewID$
  One from the leader
  $\wedge\ \exists m \in M : m.sender = Leader(m.viewID)$
  Hash values are consistent
  $\wedge\ \text{LET}$
$\qquad\qquad leaderReply\ \triangleq\ \text{CHOOSE}\ m \in M : m.sender = Leader(m.viewID)$
   IN
  $\forall m1 \in M :\ \text{IF}\ m1.mtype = MFastReply\ \text{THEN}$
$\qquad\qquad\qquad m1.hash = leaderReply.hash$
$\qquad\qquad\qquad\quad \text{ELSE}$
$\qquad\qquad\qquad\qquad \text{TRUE}$   *SlowReply* has consistent hash for sure
  Slow path
 $\vee\ \exists M \in \text{SUBSET}\ (\{m \in messages : \wedge\ \vee\ m.mtype = MSlowReply$

$$\lor \land m.mtype = MFastReply \quad \text{Leader only sends fast-reply}$$
$$\land m.sender = Leader(m.viewID)$$
$$\land m.logSlotNum = logSlotNum$$
$$\land m.dest = clientID$$
$$\land m.requestID = requestID\}):$$
$$\land \{m.sender : m \in M\} \in Quorums$$

Matching view-id
$$\land \exists m1 \in M : \forall m2 \in M : m1.viewID = m2.viewID$$

One from the leader
$$\land \exists m \in M : m.sender = Leader(m.viewID)$$

Check whether $log < clientID,\ requestID >$ is committed in view $viewID$
$$CommittedInView(clientID,\ requestID,\ viewID) \triangleq$$

Fast path
$$\lor \exists M \in \text{SUBSET } (\{m \in messages : \land \lor m.mtype = MFastReply$$
$$\lor m.mtype = MSlowReply$$
$$\land m.dest = clientID$$
$$\land m.requestID = requestID$$
$$\land m.viewID = viewID\}):$$

Sent from a fast quorum
$$\land \{m.sender : m \in M\} \in FastQuorums$$

One from the leader
$$\land \exists m \in M : m.sender = Leader(m.viewID)$$

Hash values are the same
$$\land \quad \text{LET}$$
$$leaderReply \triangleq \text{CHOOSE } m \in M : m.sender = Leader(m.viewID)$$
$$\quad \text{IN}$$
$$\forall m1 \in M : \text{IF } m1.mtype = MFastReply \text{ THEN}$$
$$m1.hash = leaderReply.hash$$
$$\quad\quad \text{ELSE}$$
$$\quad\quad\quad \text{TRUE} \quad SlowReply \text{ has consistent hash for sure}$$

Slow path
$$\lor \exists M \in \text{SUBSET } (\{m \in messages : \land \lor m.mtype = MSlowReply$$
$$\lor \land m.mtype = MFastReply \quad \text{Leader only sends fast-reply}$$
$$\land m.sender = Leader(m.viewID)$$
$$\land m.dest = clientID$$
$$\land m.requestID = requestID$$
$$\land m.viewID = viewID\}):$$
$$\land \{m.sender : m \in M\} \in Quorums$$

Hash values are the same
$$\land \exists m1 \in M : \forall m2 \in M : m1.hash = m2.hash$$

One from the leader
$$\land \exists m \in M : m.sender = Leader(m.viewID)$$

26

$SystemRecovered(viewID) \triangleq \quad \wedge \exists RM \in \text{SUBSET } (Replicas) :$
$\qquad\qquad\qquad\qquad\qquad\qquad \wedge Cardinality(RM) \geq QuorumSize$
$\qquad\qquad\qquad\qquad\qquad\qquad \wedge \forall r \in RM : vLastNormView[r] \geq viewID$
$\qquad\qquad\qquad\qquad\qquad\qquad \wedge \forall r \in RM : vReplicaStatus[r] = StNormal$ These replicas must be normal

$\qquad\qquad\qquad\qquad\qquad$ The leader of this view has also recovered or even goes beyond this view
$\qquad\qquad\qquad\qquad\qquad\qquad \wedge \quad vLastNormView[Leader(viewID)] \geq viewID$

**Invariants**

Durability: *Committed* Requests always survive failure

*i.e.* If a request is committed in one view, then it will remain committed in the higher views

One thing to note, the check of "committed" only happens when the system is still "normal"

While the system is under recovery (*i.e.* less than $f + 1$ replicas are normal),

the check of committed does not make sense

$Durability \triangleq \forall v1, v2 \in 1 .. MaxViews :$

$\qquad\qquad\qquad\qquad$ If a request is committed in lower view ($v1$,),

$\qquad\qquad\qquad\qquad$ it is impossible to make this request uncommited in higher view ($v2$)

$\qquad\qquad\qquad\qquad\quad \neg( \wedge v1 < v2$

$\qquad\qquad\qquad\qquad\qquad\qquad$ To check *Durability* of request in higher views,
$\qquad\qquad\qquad\qquad\qquad\qquad$ the system should have entered the higher views
$\qquad\qquad\qquad\qquad\qquad\qquad \wedge SystemRecovered(v2)$
$\qquad\qquad\qquad\qquad\qquad\qquad \wedge \exists c \in Clients :$
$\qquad\qquad\qquad\qquad\qquad\qquad\quad \exists r \in 1 .. MaxReqNum :$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge CommittedInView(c, r, v1)$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge \neg CommittedInView(c, r, v2))$

Consistency: *Committed* requests have the same history even after view changes

*i.e.* If a request is committed in a lower view ($v1$), then (based on *Durability* Property)

it remains committed in higher view ($v2$)

Consistency requires the history of the request (*i.e.* all the request before this request) remain the same

$Consistency \triangleq$

$\qquad \forall v1, v2 \in 1 .. MaxViews :$

$\qquad\qquad\qquad \neg( \wedge v1 < v2$

$\qquad\qquad\qquad\qquad\qquad$ To check *Consistency* of request in higher views,
$\qquad\qquad\qquad\qquad\qquad$ the system should have entered the higher views
$\qquad\qquad\qquad\qquad\qquad \wedge SystemRecovered(v2)$
$\qquad\qquad\qquad\qquad\qquad \wedge \exists c \in Clients :$
$\qquad\qquad\qquad\qquad\qquad\quad \exists r \in 1 .. MaxReqNum :$
$\qquad\qquad\qquad\qquad\qquad\quad \exists t \in 1 .. MaxTime :$
$\qquad\qquad\qquad\qquad\qquad\qquad$ Durability has been checked in another invariant
$\qquad\qquad\qquad\qquad\qquad\quad \wedge CommittedInView(c, r, v1)$
$\qquad\qquad\qquad\qquad\qquad\quad \wedge CommittedInView(c, r, v2)$
$\qquad\qquad\qquad\qquad\qquad\quad \wedge \text{LET}$
$\qquad\qquad\qquad\qquad\qquad\qquad v1LeaderReply \triangleq \text{CHOOSE } m \in messages :$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge m.mtype = MFastReply$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \wedge m.deadline = t$

27

$$
\begin{aligned}
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \land\ m.dest = c \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \land\ m.requestID = r \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \land\ m.viewID = v1 \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \land\ m.sender = Leader(v1) \\
&\qquad\qquad v2LeaderReply\ \triangleq\ \text{CHOOSE}\ m \in messages : \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \land\ m.mtype = MFastReply \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \land\ m.deadline = t \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \land\ m.dest = c \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \land\ m.requestID = r \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \land\ m.viewID = v2 \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \land\ m.sender = Leader(v2) \\
&\qquad\qquad \text{IN} \\
&\qquad\qquad\qquad v1LeaderReply.hash \neq v2LeaderReply.hash )
\end{aligned}
$$

*Linearizability*: Only one request can be committed for a given position

*i.e.* If one request has committed at position $i$, then no contrary observation can be made

*i.e.* there cannot be a second request committed at the same position

$Linearizability\ \triangleq$

  LET

    $maxLogPosition\ \triangleq\ Max(\{1\}\ \cup$

      $\{m.logSlotNum : m \in \{m \in messages :$

                    $\lor\ m.mtype = MFastReply$

                    $\lor\ m.mtype = MSlowReply\}\})$

  IN    $\neg(\exists\, c1,\ c2 \in Clients :$

        $\exists\, r1,\ r2 \in 1\,..\,MaxReqNum :$

          $\land\ \langle c1,\ r1 \rangle \neq \langle c2,\ r2 \rangle$

          $\land\ \exists\, i \in (1\,..\,maxLogPosition) :$

            $\land\ Committed(c1,\ r1,\ i)$

            $\land\ Committed(c2,\ r2,\ i)$

        )

---

## Main Transition Function

$Next\ \triangleq$  Handle Messages

    $\lor\ \exists\, m \in messages :$

               $\land\ m.mtype = MClientRequest$

               $\land\ m \notin vReplicaProcessed[m.dest]$

               $\land\ HandleClientRequest(m.dest,\ m)$

               $\land\ vReplicaProcessed' =$

                   $[vReplicaProcessed\ \text{EXCEPT}\ ![m.dest] =$

                   $vReplicaProcessed[m.dest] \cup \{Msg2RLog(m,\ m.dest)\}]$

               $\land\ \text{UNCHANGED}\ vClientProcessed$

               $\land\ DebugAction' = \langle\ \text{“HandleClientRequest”},\ m\ \rangle$

    $\lor\ \exists\, m \in messages :$

$\land\, m.mtype = MViewChangeReq$
$\land\, m \notin vReplicaProcessed[m.dest]$
$\land\, HandleViewChangeReq(m.dest,\, m)$
$\land\, vReplicaProcessed' =$
$\quad [vReplicaProcessed \text{ EXCEPT } ![m.dest] =$
$\quad vReplicaProcessed[m.dest] \cup \{Msg2RLog(m,\, m.dest)\}]$
$\land\, \text{UNCHANGED } vClientProcessed$
$\land\, DebugAction' = \langle \text{``HandleViewChangeReq''},\, m \rangle$

$\lor\, \exists\, m \in messages :$
$\qquad \land\, m.mtype = MViewChange$
$\qquad \land\, m \notin vReplicaProcessed[m.dest]$
$\qquad \land\, HandleViewChange(m.dest,\, m)$
$\qquad \land\, vReplicaProcessed' =$
$\qquad\quad [vReplicaProcessed \text{ EXCEPT } ![m.dest] =$
$\qquad\quad vReplicaProcessed[m.dest] \cup \{Msg2RLog(m,\, m.dest)\}]$
$\qquad \land\, \text{UNCHANGED } vClientProcessed$
$\qquad \land\, DebugAction' = \langle \text{``HandleViewChange''},\, m \rangle$

$\lor\, \exists\, m \in messages :$
$\qquad \land\, m.mtype = MStartView$
$\qquad \land\, m \notin vReplicaProcessed[m.dest]$
$\qquad \land\, HandleStartView(m.dest,\, m)$
$\qquad \land\, vReplicaProcessed' =$
$\qquad\quad [vReplicaProcessed \text{ EXCEPT } ![m.dest] =$
$\qquad\quad vReplicaProcessed[m.dest] \cup \{Msg2RLog(m,\, m.dest)\}]$
$\qquad \land\, \text{UNCHANGED } vClientProcessed$
$\qquad \land\, DebugAction' = \langle \text{``HandleStartView''},\, m \rangle$

$\lor\, \exists\, m \in messages :$
$\qquad \land\, m.mtype = MSyncPrepare$
$\qquad \land\, m \notin vReplicaProcessed[m.dest]$
$\qquad \land\, HandleSyncPrepare(m.dest,\, m)$
$\qquad \land\, vReplicaProcessed' =$
$\qquad\quad [vReplicaProcessed \text{ EXCEPT } ![m.dest] =$
$\qquad\qquad vReplicaProcessed[m.dest] \cup \{Msg2RLog(m,\, m.dest)\}]$
$\qquad \land\, \text{UNCHANGED } vClientProcessed$
$\qquad \land\, DebugAction' = \langle \text{``HandleSyncPrepare''},\, m \rangle$

$\lor\, \exists\, m \in messages :$
$\qquad \land\, m.mtype = MSyncRep$
$\qquad \land\, m \notin vReplicaProcessed[m.dest]$
$\qquad \land\, HandleSyncRep(m.dest,\, m)$
$\qquad \land\, vReplicaProcessed' =$
$\qquad\quad [vReplicaProcessed \text{ EXCEPT } ![m.dest] =$
$\qquad\quad vReplicaProcessed[m.dest] \cup \{Msg2RLog(m,\, m.dest)\}]$

$$\land \text{UNCHANGED } vClientProcessed$$
$$\land DebugAction' = \langle \text{``HandleSyncRep''}, m \rangle$$
$$\lor \exists\, m \in messages :$$
$$\land m.mtype = MSyncCommit$$
$$\land m \notin vReplicaProcessed[m.dest]$$
$$\land HandleSyncCommit(m.dest, m)$$
$$\land vReplicaProcessed' =$$
$$[vReplicaProcessed \text{ EXCEPT } ![m.dest] =$$
$$vReplicaProcessed[m.dest] \cup \{Msg2RLog(m, m.dest)\}]$$
$$\land \text{UNCHANGED } vClientProcessed$$
$$\land DebugAction' = \langle \text{``HandleSyncCommit''}, m \rangle$$

$$\lor \exists\, m \in messages :$$
$$\land m.mtype = MMissEntryRequest$$
$$\land m \notin vReplicaProcessed[m.dest]$$
$$\land HandleMissEntryRequest(m.dest, m)$$
$$\land vReplicaProcessed' =$$
$$[vReplicaProcessed \text{ EXCEPT } ![m.dest] =$$
$$vReplicaProcessed[m.dest] \cup \{Msg2RLog(m, m.dest)\}]$$
$$\land \text{UNCHANGED } vClientProcessed$$
$$\land DebugAction' = \langle \text{``HandleMissEntryRequest''}, m \rangle$$

$$\lor \exists\, m \in messages :$$
$$\land m.mtype = MMissEntryReply$$
$$\land m \notin vReplicaProcessed[m.dest]$$
$$\land HandleMissEntryReply(m.dest, m)$$
$$\land vReplicaProcessed' =$$
$$[vReplicaProcessed \text{ EXCEPT } ![m.dest] =$$
$$vReplicaProcessed[m.dest] \cup \{Msg2RLog(m, m.dest)\}]$$
$$\land \text{UNCHANGED } vClientProcessed$$
$$\land DebugAction' = \langle \text{``HandleMissEntryReply''}, m \rangle$$

Client Actions
$$\lor \exists\, c \in Clients :$$
$$\land vClientReqNum[c] < MaxReqNum$$
$$\land ClientSendRequest(c)$$
$$\land \text{UNCHANGED } \langle vReplicaProcessed, vClientProcessed \rangle$$
$$\land DebugAction' = \langle \text{``ClientSendRequest''}, \text{``''} \rangle$$

Start Synchronization
$$\lor \exists\, r \in Replicas :$$
$$\land StartSync(r)$$
$$\land \text{UNCHANGED } \langle vReplicaProcessed, vClientProcessed \rangle$$
$$\land DebugAction' = \langle \text{``StartSync''}, \text{``''} \rangle$$

Replica Fail

$\lor \exists\, r \in Replicas :$
$\qquad\qquad \land vReplicaStatus[r] = StNormal$
$\qquad\qquad \land StartReplicaFail(r)$
$\qquad\qquad \land \textsc{unchanged}\ \langle vReplicaProcessed,\ vClientProcessed \rangle$
$\qquad\qquad \land DebugAction' = \langle \text{``StartReplicaFail''},\ \text{``''} \rangle$

Leader Change
$\lor \exists\, r \in Replicas :$
$\qquad\qquad \land vViewID[r] < MaxViews$
$\qquad\qquad \land StartLeaderChange(r)$
$\qquad\qquad \land \textsc{unchanged}\ \langle vReplicaProcessed,\ vClientProcessed \rangle$
$\qquad\qquad \land DebugAction' = \langle \text{``StartLeaderChange''},\ \text{``''} \rangle$

Replica Rejoin
$\lor \exists\, r \in Replicas :$
$\qquad\qquad \land vReplicaStatus[r] = StRecovering$
$\qquad\qquad \land StartReplicaRecovery(r)$
$\qquad\qquad \land \textsc{unchanged}\ \langle vReplicaProcessed,\ vClientProcessed \rangle$
$\qquad\qquad \land DebugAction' = \langle \text{``StartReplicaRecovery''},\ \text{``''} \rangle$

Replica Actions:
$\lor \exists\, r \in Replicas :$
$\qquad\qquad \land StartIndexSync(r)$
$\qquad\qquad \land \textsc{unchanged}\ \langle vReplicaProcessed,\ vClientProcessed \rangle$
$\qquad\qquad \land DebugAction' = \langle \text{``StartIndexSync''},\ \text{``''} \rangle$

$\lor \exists\, r \in Replicas :$
$\qquad\qquad \land FlushEarlyBuffer(r)$
$\qquad\qquad \land \textsc{unchanged}\ \langle vReplicaProcessed,\ vClientProcessed \rangle$
$\qquad\qquad \land DebugAction' = \langle \text{``FlushReplicaBuffer''},\ \text{``''} \rangle$

Clock Move
$\lor \exists\, r \in Replicas :$
$\qquad\qquad \land ReplicaClockMove(r)$
$\qquad\qquad \land \textsc{unchanged}\ \langle vReplicaProcessed,\ vClientProcessed \rangle$
$\qquad\qquad \land DebugAction' = \langle \text{``ReplicaClockMove''},\ \text{``''} \rangle$

$\lor \exists\, c \in Clients\ :$
$\qquad\qquad \land ClientClockMove(c)$
$\qquad\qquad \land \textsc{unchanged}\ \langle vReplicaProcessed,\ vClientProcessed \rangle$
$\qquad\qquad \land DebugAction' = \langle \text{``ClientClockMove''},\ \text{``''} \rangle$