# Tiga TLA+ Specification

─────────────── MODULE *Tiga* ───────────────

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

─────────────────────────────────────────────

## Bounds for Model Check [Configurable]

Time Range [Configurable]
$MaxTime \triangleq 3$

In *Tiga*, we assume client and coordinator are co-located
In this spec, we use "coordinator" to represent them
Each coordinator 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]
$MaxViews \triangleq 3$

The set of replicas and an ordering of them [Can be configured in TLA+ *Toolbox*]
$Replicas \triangleq 0 \,.\, .\, 2$
$ReplicaOrder \triangleq \langle 0,\, 1,\, 2 \rangle$
$Shards \triangleq 0 \,.\, .\, 2$
$Coords \triangleq 0 \,.\, .\, 1$
$LatencyBounds \triangleq [c \in Coords \mapsto 1]$

ASSUME  *IsFiniteSet*(*Replicas*)
ASSUME  *IsFiniteSet*(*Shards*)
ASSUME *ReplicaOrder* $\in Seq(Replicas)$
$Servers \triangleq \{$
$\quad [$
$\qquad replicaId \mapsto e[1],$
$\qquad shardId \mapsto e[2]$
$\quad ] : e \in Replicas \times Shards$
$\}$

─────────────────────────────────────────────

These variables are used to implment at-most-once primitives

## Constants
$F \triangleq (Cardinality(Replicas) - 1) \div 2$

$ceilHalfF \triangleq \text{IF } (F \div 2) * 2 = F \text{ THEN } F \div 2 \text{ ELSE } (F + 1) \div 2$

$floorHalfF \quad \triangleq \ F \div 2$

$QuorumSize \triangleq \ F + 1$

$FastQuorumSize \triangleq \ F + ceilHalfF + 1$

$RecoveryQuorumSize \triangleq \ ceilHalfF + 1$

$FastQuorums \triangleq \ \{R \in \text{SUBSET } (Replicas) : \\ \qquad\qquad\quad Cardinality(R) \geq FastQuorumSize\}$

$Quorums \triangleq \ \{R \in \text{SUBSET } (Replicas) : \\ \qquad\qquad Cardinality(R) * 2 > Cardinality(Replicas)\}$

## Server Status

$StNormal \triangleq 1$

$StViewChange \triangleq 2$

$StCrossShardSyncing \triangleq 3$

$StRecovering \triangleq 4$

$StFailing \triangleq 5$

## Message Types

$MTxn \triangleq 1$

$MLogEntry \triangleq 2$    *Log* entry, different from index, it includes command field, which can be large in practice

$MDeadlineNotification \triangleq 3$ Leaders send the message to other leaders for deadline agreement

$MInterReplicaSync \triangleq 4$    Synchronize within shard group (across replicas) to ensure strict serializability

$MFastReply \triangleq 5$ Fast Reply Message

$MSlowReply \triangleq 6$ Slow Reply Message

The following messages are mainly for view change within each sharding group

$MViewChangeReq \triangleq 7$      Sent by config manager when leader/sequencer failure detected

$MViewChange \triangleq 8$      Sent to $ACK$ view change

$MStartView \triangleq 9$      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

$MLocalSyncStatus \triangleq 10$      Sent by the leader to ensure $log$ durability

$MLocalCommit \triangleq 11$      Sent by followers as $ACK$

The following messages are used for periodic sync across sharding groups

This is an optional optimization to enable fast recovery

$MPeerShardCommitStatus \triangleq 12$

The following messages are mainly used for server recovery

$MCrashVectorReq \triangleq 13$

$MCrashVectorRep \triangleq 14$

$MRecoveryReq \triangleq 15$

$MRecoveryRep \triangleq 16$

$MStartViewReq \triangleq 17$

$MCrossShardConfirm \triangleq 19$

Config Manager ($CM$)'s Operations. Since $CM$ is supported by typical viewstamped replication ($VR$), in this spec, we do not repeat the $VR$'s failure recovery spec for $CM$

$MCMPrepare \triangleq 20$
$MCMPrepareReply \triangleq 21$
$MCMCommit \triangleq 22$

## Message Schemas

Each server is identified by a combination of $< replicaId, shardId > TxnID$ uniquely identifies one request on one server But across replicas, the same $TxnID$ may have different deadlines (the leader may modify the deadline to make the request eligible to enter the early-buffer) so $< deadline, txnId >$ uniquely identifes one request across replicas

$TxnID = [$
  $coordId \mapsto i$ in $(1 .. ),$
  $rId \quad \mapsto i$ in $(1 .. )$
$]$

$Txn = [$
  $mtype \mapsto MTxn$
  $txnId \mapsto TxnID,$
  $shards \mapsto Shards,$
  $command \mapsto command,$
  $st \quad \mapsto sendTime,$
  $bound \mapsto latencyBound$
$]$

$LogEntry = [$
  $mtype \mapsto MLogEntry$
  $txnId \mapsto TxnID,$
  $shards \mapsto Shards,$
  $command \mapsto command,$
  $deadline \mapsto deadline$
$]$

After the request arrives at the *shards* and is placed into its early buffer (either with deadline modified or not), the server will broadcast *DeadlineNotification* to all the other servers in the same replica group to tell them the deadline of the request on its own server

$DeadlineNotification = [$
  $mtype \mapsto MDeadlineNotification,$
  $gView \mapsto 0 \ldots x$
  $lView \mapsto 0 \ldots y$
  $sender \mapsto src \in Servers,$
  $dest \mapsto dst \in Servers,$
  $entry \mapsto LogEntry$
$]$

After leader has released the *txn*, it synchornizes the *log* with its followers. If followers are inconsistent, they will rectify their logs to keey consistent with leader

$InterReplicaSync = [$
   $mtype \quad \mapsto MInterReplicaSync,$
   $lView \quad \mapsto 0 \ldots y$
   $sender \quad \mapsto src \in Servers,$
   $dest \quad \mapsto dst \in Servers,$
   $entries \quad \mapsto [LogEntry \ldots]$
$]$

*logId* (*i.e.*, the position index of the *log* entry in the *log* list) 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

$FastReply = [$
   $mtype \quad \mapsto MFastReply,$
   $sender \quad \mapsto src \in Servers,$
   $dest \quad \mapsto dst \in Coords,$
   $gView \quad \mapsto 0 \ldots x$
   $lView \quad \mapsto 0 \ldots x$
   $txnId \quad \mapsto txnId$

   In real implementation, we use $SHA1 +$ Incremental Hash

   $hash \quad \mapsto [\ entries \mapsto log \text{ entries so far } cv \mapsto crashVector\ ]$
   $deadline \mapsto i \in (1 \mathrel{..} MaxTime + MaxBound),$
   $logId \quad \mapsto n \in (1 \mathrel{..})$
$]$

$SlowReply = [$
   $mtype \quad \mapsto MSlowReply,$
   $sender \quad \mapsto src \in Servers,$
   $dest \quad \mapsto c \in Coords,$
   $gView \quad \mapsto 0 \ldots x$
   $lView \quad \mapsto 0 \ldots x$
   $txnId \quad \mapsto txnId$
   $logId \quad \mapsto n \in (1 \mathrel{..})$
$]$

$ViewChangeReq = [$
   $mtype \mapsto MViewChangeReq,$
   $sender \mapsto src \in Replicas, \text{ (by } configManager)$
   $dest \mapsto dst \in Servers,$
   $gView \mapsto 0 \mathrel{..} x$
   $gVec \mapsto$ the *lViews* for each shard
$]$

$ViewChange = [$
   $mtype \quad \mapsto MViewChange,$
   $sender \quad \mapsto src \in Servers,$
   $dest \quad \mapsto dst \in Servers,$
   $gView \quad \mapsto 0 \mathrel{..} x$
   $gVec \quad \mapsto$ the *lViews* for each shard
   $lView \quad \mapsto 0 \ldots x$
   $lastNormal \mapsto v \in ViewIDs,$

```
    lSyncPoint ↦ 0 . .
    entries    ↦ l ∈ vLogs[1 . . n],
    cv         ↦  crash vector
]

CrossShardConfirm =  [
    mtype     ↦ MCrossShardConfirm,
    sender    ↦ src ∈ Servers,
    dest      ↦ dst ∈ Servers,
    lView     ↦ 0 . . . x
    gView     ↦ 0 . . .
    entries ↦ l ∈ vLogs[1 . . n]
]

StartView =  [
    mtype     ↦ MStartView,
    sender    ↦ src ∈ Servers,
    dest      ↦ dst ∈ Servers,
    lView     ↦ 0 . . . x
    gView     ↦ 0 . . . x
    gVec      ↦  the lViews for each shard
    entries ↦ l ∈ vLogs[1 . . n],
    cv        ↦  crash vector
]


CrashVectorReq =  [
    mtype     ↦ MCrashVectorReq,
    sender    ↦ src ∈ Servers,
    dest      ↦ dst ∈ Servers,
    nonce     ↦ nonce
]

CrashVectorRep =  [
    mtype      ↦ MCrashVectorRep,
    sender     ↦ src ∈ Servers,
    dest       ↦ dst ∈ Servers,
    nonce      ↦ nonce,
    cv         ↦  vector of counters
]

RecoveryReq =  [
    mtype       ↦ MRecoveryReq,
    sender      ↦ src ∈ Servers,
    dest        ↦ dst ∈ Servers,
    cv          ↦  vector of counters
]

RecoveryRep = [
    mtype ↦ MRecoveryRep,
    sender ↦ src ∈ Servers,
    dest ↦ dst ∈ Servers,
    gView ↦ 0 . . x
```

$$lView \mapsto 0 \,.\,.\, x$$
$$cv \qquad \mapsto vector\ of\ counters$$
]

$StartViewReq = [$
$\quad mtype \qquad \mapsto MStartViewReq,$
$\quad sender \qquad \mapsto src \in Servers,$
$\quad dest \qquad \mapsto dst \in Servers,$
$\quad lView \qquad \mapsto 0 \,.\,.\, x$
$\quad cv \qquad \mapsto vector\ of\ counters$
]


Follower reports to its leader

$LocalSyncStatus = \ [$
$\quad mtype \qquad \mapsto MLocalSyncStatus,$
$\quad sender \qquad \mapsto src \in Servers,$
$\quad dest \qquad \mapsto dst \in Servers,$
$\quad lView \qquad \mapsto 0 \ldots x$
$\quad lSyncPoint \mapsto n \in \ (1 \,.\,.\, )$
$\quad cv \qquad \mapsto \ vector\ of\ counters$
]

Leader notifies its followers

$LocalCommit = \ [$
$\quad mtype \qquad \mapsto MLocalCommit,$
$\quad sender \qquad \mapsto src \in Servers,$
$\quad dest \qquad \mapsto dst \in Servers,$
$\quad lView \qquad \mapsto 0 \ldots x$
$\quad entries \qquad \mapsto log\ entries$
$\quad lCommitPoint \mapsto n \in \ (1 \ldots )$
]

Each server tells its neighbors (the servers in the same region but belong to different *shards*) its local commit status. This is optional optimization (only for checkpoint and failure recovery acceleration)

$PeerShardCommitStatus = [$
$\quad mtype \qquad \mapsto MPeerShardCommitStatus,$
$\quad sender \qquad \mapsto src \in Servers,$
$\quad dest \qquad \mapsto dst \in Servers,$
$\quad gView \qquad \mapsto 0 \,.\,.\, x$
$\quad deadline \mapsto the\ largest\ committed\ deadline$
]


Configuration Manager ($CM$)'s message to prepare global information (including $gView$ and $gVec$)

In our implementation, $CM$ is co-located on Shard $- 0$, but from design perspective, $CM$ is completed standalone and decoupled from *Tiga Servers*

$CMPrepare = [$

$$mtype \mapsto MCMPrepare,$$
$$sender \mapsto src \in Servers,$$
$$dest \mapsto dst \in Servers,$$
$$cView \mapsto 0 .. x$$
$$gView \mapsto 0 .. x$$
$$gVec \mapsto [shardId \mapsto lView]$$
]

$CMPrepareReply = [$
$$mtype \mapsto MCMPrepareReply,$$
$$sender \mapsto src \in Servers,$$
$$dest \mapsto dst \in Servers,$$
$$cView \mapsto 0 .. x$$
$$gView \mapsto 0 .. x$$
]

$CMCommit = [$
$$mtype \mapsto MCMPrepareReply,$$
$$sender \mapsto src \in Servers,$$
$$dest \mapsto dst \in Servers,$$
$$cView \mapsto 0 .. x$$
$$gView \mapsto 0 .. x$$
]

---

**Network State**

VARIABLES *messages*  Set of all messages sent

**Server State**

VARIABLES

Messages that have been processed by servers

$vServerProcessed,$

*Log* list of entries

$vLog,$

The sequencer to hold txns and release it after clock passes its deadline $(s + l)$

$vEarlyBuffer,$

The buffer to hold txns on followers because these txns come too late and cannot enter early-buffer

$vLateBuffer,$

Each leader server has a data structure of *DeadlineQuroum* to collect the deadlines from other servers for agreement

$vDeadlineQuorum,$

After servers have recovered their logs from the signle shard, they need confirmation from the other *shards* to ensure the recovered logs satisfy strict serializability

$vCrossShardConfirmQuorum,$

One of *StNormal*, *StViewChange*, *StFailing*, *StRecovering*

7

$vServerStatus,$

Global views of each server

$vGView,$

The g-vecs of each server

$vGVec,$

Local views of each server

$vLView,$

Current Time of the server

$vServerClock,$

Last *lView* in which this server had *StNormal* status

$vLastNormView,$

Used for collecting view change votes

$vViewChange,$

*vLSyncPoint* indicates to which the server state (*vLog*) is consistent with the leader.

$vLSyncPoint,$

*vLCommitPoint* indicates that the *log* entries before this point has been locally committed, *i.e.*, replicated to majority in this sharding groups. So followers can safely execute the logged txns

$vLCommitPoint,$

*vPeerCommitDeadline* records the peer's largest deadline that has been locally committed. This can be used to save data transfer during cross-shard confirmation

$vPeerCommitDeadline,$

*vLSyncQuorum* is used by each leader to collect the *LocalSyncStatus* messages from servers in the same sharding group

$vLSyncQuorum,$

Locally unique string (for *CrashVectorReq*)

$vUUIDCounter,$

*CrashVector*, initialized as all-zero vector

$vCrashVector,$
$vCrashVectorReps,$
$vRecoveryReps$


**Coordinator State**

VARIABLES      Current Clock Time of the coordinator

$vCoordClock,$

The txns that have been sent by this coordinator. This variable makes it easy to derive the Invariants

$vCoordTxns,$

Messages that have been processed by coordinators

$vCoordProcessed$

**Configuration Manager (CM) State**

VARIABLES

Since *CM* is supported by traditional *VR*, here we do not want to repeat *VR*'s failure recovery in this spec, so we make *CMStatus* always *StNormal*

$vCMStatus$,
$vCMView$,

Config Manager: the latest global info the manager maintains ($gView$ and $gVec$)

$vCMGInfo$,
$vCMPrepareGInfo$,

Config Manager: quorum of *CMPrepareReplies*

$vCMPrepareReps$,

$vCMProcessed$

VARIABLES   *ActionName*

$networkVars \triangleq \langle messages \rangle$

$serverStateVars \triangleq$
$\langle vLog, vEarlyBuffer, vLateBuffer,$
$vDeadlineQuorum, vCrossShardConfirmQuorum, vServerStatus,$
$vGView, vGVec, vLView, vServerClock, vLastNormView,$
$vViewChange, vLSyncPoint, vLCommitPoint,$
$vPeerCommitDeadline, vLSyncQuorum,$
$vUUIDCounter, vCrashVector, vCrashVectorReps,$
$vRecoveryReps, vServerProcessed \rangle$

$coordStateVars \triangleq \langle vCoordClock, vCoordTxns, vCoordProcessed \rangle$

$configManagerStateVars \triangleq \langle vCMStatus, vCMView, vCMGInfo,$
$\qquad\qquad\qquad\qquad\qquad vCMPrepareGInfo, vCMPrepareReps,$
$\qquad\qquad\qquad\qquad\qquad vCMProcessed \rangle$

$InitNetworkState \triangleq messages = \{\}$

$InitServerState \triangleq$
$\quad \wedge\ vServerProcessed = [serverId \in Servers \mapsto \{\}]$
$\quad \wedge\ vLog = [serverId \in Servers \mapsto \langle \rangle]$
$\quad \wedge\ vEarlyBuffer\ \ = [serverId \in Servers \mapsto \{\}]$
$\quad \wedge\ vLateBuffer\ \ \ = [serverId \in Servers \mapsto \{\}]$
$\quad \wedge\ vDeadlineQuorum\ \ = \ \ [serverId \in Servers \mapsto \{\}]$
$\quad \wedge\ vCrossShardConfirmQuorum = [serverId \in Servers \mapsto \{\}]$
$\quad \wedge\ vServerStatus\ \ \ \ = \ \ [serverId \in Servers \mapsto StNormal]$
$\quad \wedge\ vGView\ \ = \ \ [serverId \in Servers \mapsto 0]$
$\quad \wedge\ vGVec = [$

$$
\begin{aligned}
&\quad serverId \in Servers \mapsto [ \\
&\qquad shardId \in Shards \mapsto 0 \\
&\quad ] \\
&] \\
&\wedge\ vLView\ =\ \ [serverId \in Servers \mapsto 0] \\
&\wedge\ vServerClock =\ \ [serverId \in Servers \mapsto 1] \\
&\wedge\ vLastNormView\ \ = [serverId \in Servers \mapsto 0] \\
&\wedge\ vViewChange = [serverId \in Servers \mapsto \{\}] \\
&\wedge\ vLSyncPoint\ \ = [serverId \in Servers \mapsto 0] \\
&\wedge\ vLCommitPoint\ \ = [serverId \in Servers \mapsto 0] \\
&\wedge\ vPeerCommitDeadline\ \ = [serverId \in Servers \mapsto \\
&\qquad [shardId \in Shards \mapsto 0] \\
&\quad ] \\
&\wedge\ vLSyncQuorum\ \ \ = [serverId \in Servers \mapsto \{\}] \\
&\wedge\ vUUIDCounter = [serverId \in Servers \mapsto 0] \\
&\wedge\ vCrashVector\ \ = [ \\
&\qquad serverId\ \ \in Servers \mapsto [ \\
&\qquad\quad rr \in Replicas \mapsto 0 \\
&\qquad ] \\
&\quad ] \\
&\wedge\ vCrashVectorReps = [serverId \in Servers \mapsto \{\}] \\
&\wedge\ vRecoveryReps\ \ \ = [serverId \in Servers \mapsto \{\}]
\end{aligned}
$$

$$
\begin{aligned}
InitCoordState\ &\triangleq \\
&\wedge\ vCoordProcessed = [c \in Coords \mapsto \{\}] \\
&\wedge\ vCoordClock\ \ \ = [c \in Coords\ \mapsto 1] \\
&\wedge\ vCoordTxns\ \ = [c \in Coords\ \mapsto \{\}]
\end{aligned}
$$

$$
\begin{aligned}
InitConfigManagerState\ &\triangleq \\
&\wedge\ vCMStatus = [ \\
&\qquad replicaId \in Replicas \mapsto StNormal \\
&\quad ] \\
&\wedge\ vCMView = [ \\
&\qquad replicaId \in Replicas \mapsto 0 \\
&\quad ] \\
&\wedge\ vCMGInfo = [ \\
&\qquad replicaId\ \ \in Replicas \mapsto [ \\
&\qquad\quad gView\ \ \mapsto 0, \\
&\qquad\quad gVec\ \ \ \mapsto [shardId \in Shards \mapsto 0] \\
&\qquad ] \\
&\quad ] \\
&\wedge\ vCMPrepareGInfo = [ \\
&\qquad replicaId\ \ \in Replicas \mapsto [ \\
&\qquad\quad gView\ \ \ \mapsto 0,
\end{aligned}
$$

$$gVec \quad \mapsto [shardId \in Shards \mapsto 0]$$
$$]$$
$$]$$
$$\land \ vCMPrepareReps = [$$
$$replicaId \in Replicas \mapsto \ \{\}$$
$$]$$
$$\land \ vCMProcessed = [$$
$$replicaId \in Replicas \mapsto \{\}$$
$$]$$

$PickMax(S) \ \triangleq \ \text{CHOOSE} \ \ x \in S : \forall\, y \in S : y \leq x$

$PickMin(S) \ \triangleq \ \text{CHOOSE} \ \ x \in S : \forall\, y \in S : y \geq x$

$Min(a,\, b) \ \triangleq \ \text{IF} \ a < b \ \text{THEN} \ a \ \text{ELSE} \ \ b$

$Max(a,\, b) \ \triangleq \ \text{IF} \ a < b \ \text{THEN} \ b \ \text{ELSE} \ \ a$

$Send(ms) \ \triangleq \ messages' = messages \cup ms$

$SeqToSet(s) \ \triangleq$
$\{s[i] : i \in \text{DOMAIN} \ s\}$

$IsInjective(s) \ \triangleq$

TRUE iff the sequence $s$ contains no duplicates where two elements a, $b$ of $s$ are defined to be duplicates iff $a = b$. In other words,
$Cardinality(ToSet(s)) = Len(s)$

This definition is overridden by *TLC* in the *Java* class *SequencesExt*. The operator is overridden by the *Java* method with the same name.

Also see Functions!Injective operator.

$\forall\, i,\, j \in \text{DOMAIN} \ s : (s[i] = s[j]) \Rightarrow (i = j)$

$SetToSeq(S) \ \triangleq$

Convert a set to some sequence that contains all the elements of the set exactly once, and contains no other elements.

$\text{CHOOSE} \ f \in [1 \mathinner{\ldotp\ldotp} Cardinality(S) \to S] : IsInjective(f)$

$Remove(s,\, e) \ \triangleq$

The sequence $s$ with $e$ removed or $s$ iff $e \notin Range(s)$

$SelectSeq(s, \text{LAMBDA} \ t : t \neq e)$

$SetToSortSeq(S,\, op(\_,\, \_)) \ \triangleq$

Convert a set to a sorted sequence that contains all the elements of the set exactly once, and contains no other elements. Not defined via CHOOSE like *SetToSeq* but with an additional conjunct, because this variant works efficiently without a dedicated *TLC* override.

$SortSeq(SetToSeq(S),\ op)$

**View ID Helpers**

$LeaderID(viewId)\ \triangleq\ ReplicaOrder[(viewId\%Len(ReplicaOrder))+1]$ remember $\langle\rangle$ are 1-indexed

$isLeader(replicaId,\ viewId)\ \triangleq\ (replicaId = LeaderID(viewId))$

$PrintVal(id,\ exp)\ \triangleq\ Print(\langle id,\ exp\rangle,\ \text{TRUE})$

$ViewGreater(gv1,\ lv1,\ gv2,\ lv2)\ \triangleq$
   IF $gv1 > gv2$ THEN TRUE
    ELSE
      IF  $\wedge\ gv1 = gv2$
          $\wedge\ lv1 > lv2$
     THEN    TRUE
     ELSE    FALSE

Coordinator $c$ submits a *txn*. We assume Coordinator can only send one *txn* in one tick of time.
If time has reached the bound, this client cannot send request any more

$LastAppendedDeadline(Log)\ \triangleq$ IF $Len(Log) = 0$ THEN $0$
                                ELSE  $Tail(Log).deadline$

$CoordSubmitTxn(c)\ \ \ \triangleq$
   $\wedge\ vCoordClock[c] < MaxTime$
   $\wedge\ Cardinality(vCoordTxns[c]) < MaxReqNum$
   $\wedge$ LET
         $txnId\ \triangleq\ [$
           $coordId \mapsto c,$
           $rId\ \ \ \ \ \mapsto Cardinality(vCoordTxns[c]) + 1$
         $]$
      IN
     $\wedge\ Send(\{[mtype\ \ \ \mapsto MTxn,$
            $txnId\ \ \ \ \ \ \mapsto txnId,$
            $command \mapsto\ \text{""},$
             Here we assume involves all *shards*
            $shards\ \ \mapsto Shards,$
            $st\ \ \ \ \ \ \ \mapsto vCoordClock[c],$
            $bound\ \ \ \mapsto LatencyBounds[c],$
            $sender\ \ \mapsto c,$
            $dest\ \ \ \ \ \mapsto serverId$
        $] : serverId\ \ \in Servers\})$
     $\wedge\ vCoordClock' = [vCoordClock \text{ EXCEPT } ![c] = vCoordClock[c] + 1]$
     $\wedge\ vCoordTxns' = [vCoordTxns \text{ EXCEPT } ![c] = vCoordTxns[c] \cup \{txnId\}]$

$HandleTxn(m) \triangleq$
    LET
        $myServerId \triangleq m.dest$
        $newLog \triangleq [$
            $mtype \quad \mapsto MLogEntry,$
            $txnId \quad \mapsto m.txnId,$
            $command \mapsto m.command,$
            $shards \quad \mapsto m.shards,$
            $deadline \quad \mapsto Max(LastAppendedDeadline(vLog[myServerId]), m.st + m.bound)$
        $]$
        $serversInOneReplica \triangleq \{s \in Servers : s.replicaId = myServerId.replicaId\}$
    IN
    $\lor$  $\land isLeader(myServerId.replicaId, vLView[myServerId])$
        $\land vEarlyBuffer' = [$
        $vEarlyBuffer$ EXCEPT $![myServerId]$
            $= vEarlyBuffer[myServerId] \cup \{newLog\}]$
         Broadcast deadline notifications to other *shards*
        $\land Send(\{[$
        $mtype \quad \mapsto MDeadlineNotification,$
        $gView \quad \mapsto vGView[myServerId],$
        $lView \quad \mapsto vLView[myServerId],$
        $sender \quad \mapsto myServerId,$
        $dest \quad \mapsto dstServerId,$
        $entry \quad \mapsto newLog$
        $] : dstServerId \in serversInOneReplica\})$
        $\land$ UNCHANGED $\langle vLateBuffer \rangle$
    $\lor$  $\land \neg isLeader(myServerId.replicaId, vLView[myServerId])$
        $\land$  $\lor$  $\land newLog.deadline = (m.st + m.bound)$
              $\land vEarlyBuffer' = [$
                $vEarlyBuffer$ EXCEPT $![myServerId]$
                    $= vEarlyBuffer[myServerId] \cup \{newLog\}$
                $]$
              $\land$ UNCHANGED $\langle vLateBuffer \rangle$
            $\lor$  $\land \neg(newLog.deadline = (m.st + m.bound))$
              $\land$  $vLateBuffer' = [$
                $vLateBuffer$ EXCEPT $![myServerId]$
                    $= vLateBuffer[myServerId] \cup \{newLog\}$
                $]$
              $\land$ UNCHANGED $\langle vEarlyBuffer \rangle$
        $\land$ UNCHANGED $\langle networkVars \rangle$

$HandleDeadlineNotification(m) \triangleq$
    LET
        $myServerId \triangleq m.dest$

13

$$quorum \stackrel{\Delta}{=} \{$$
$$\quad msg \in vDeadlineQuorum[myServerId]$$
$$\quad\quad : \quad \wedge msg.entry.txnId = m.entry.txnId$$
$$\quad\quad\quad \wedge msg.gView = m.gView$$
$$\quad\quad\quad \wedge m.gView = vGView[myServerId]$$
$$\} \cup \{m\}$$

IN

Only leader does deadline agreement
$$\wedge \; vGView[myServerId] = m.gView$$
$$\wedge \; vGVec[myServerId][m.sender.shardId] = m.lView$$
$$\wedge \; isLeader(myServerId.replicaId, vLView[myServerId])$$
$$\wedge \; vDeadlineQuorum' = [$$
$$\quad vDeadlineQuorum \text{ EXCEPT } ![myServerId]$$
$$\quad\quad = vDeadlineQuorum[myServerId] \cup \{m\}$$
$$]$$
$$\wedge \; \text{IF} \quad Cardinality(quorum) = Cardinality(m.entry.shards)$$

THEN

Deadline quorum established : Update the deadline of the *txn* in Sequencer

LET
$$maxDeadlineTxn \stackrel{\Delta}{=}$$
$$\quad \text{CHOOSE } x \in quorum :$$
$$\quad\quad \forall y \quad \in quorum :$$
$$\quad\quad\quad y.entry.deadline \leq x.entry.deadline$$
$$sequencingTxn \stackrel{\Delta}{=}$$
$$\quad \text{CHOOSE } x \in vEarlyBuffer[myServerId] :$$
$$\quad\quad x.txnId = m.entry.txnId$$

IN

IF $maxDeadlineTxn.entry.deadline > sequencingTxn.deadline$

THEN
$$vEarlyBuffer' = [vEarlyBuffer \text{ EXCEPT } ![myServerId]$$
$$\quad = (vEarlyBuffer[myServerId] \setminus \{sequencingTxn\}) \cup \{maxDeadlineTxn.entry\}]$$
ELSE UNCHANGED $\langle vEarlyBuffer \rangle$

ELSE

Deadline quorum not sufficient so far: do not take further actions
UNCHANGED $\langle vEarlyBuffer \rangle$

$HandleInterReplicaSync(m) \stackrel{\Delta}{=}$
$$\quad \wedge m.lView = vLView[m.dest]$$

Even if $m$'s *crashVector* is newer (larger value), we do not accept *it. The* consistency of *crashVector* will finally be solved during viewchange

$$\wedge m.crashVector[m.sender] = vCrashVector[m.sender]$$
$$\wedge \neg isLeader(m.dest.replicaId, vLView[m.dest])$$
$$\wedge \text{LET}$$

$mySeverId \triangleq m.dest$

$syncedTxnIds \triangleq \{m.entries[i].txnId : i \in 1 .. Len(m.entries)\}$

$currentSyncPoint \triangleq Len(vLSyncPoint[mySeverId])$

IN

$\vee \quad \wedge \ currentSyncPoint < Len(m.entries)$

$\qquad \wedge \ vLog' = [vLog \text{ EXCEPT } ![mySeverId] = m.entries]$

Kick synced entries out of *earlyBuffer*

$\qquad \wedge \ vEarlyBuffer' = [$

$\qquad\qquad vEarlyBuffer \text{ EXCEPT } ![mySeverId]$

$\qquad\qquad\quad = \{msg \in vEarlyBuffer[mySeverId] :$

$\qquad\qquad\qquad msg.txnId \notin syncedTxnIds\}$

$\qquad\qquad ]$

Kick synced entries out of late buffer. In actual implementation, *InterReplicaSync* only carries *log* indices, and the entries are fetched from Late Buffer first, if still missing, then it will go to ask leader. Such a design can save much unncessary transmission in practice.

$\qquad \wedge \ vLateBuffer' = [$

$\qquad\qquad vLateBuffer \text{ EXCEPT } ![mySeverId]$

$\qquad\qquad\quad = \{msg \in vLateBuffer[mySeverId] :$

$\qquad\qquad\qquad msg.txnId \notin syncedTxnIds\}$

$\qquad\qquad ]$

Kick synced entries out of deadline quorum. These txns have been synced, no need to record in *DeadlineQuorum*

$\qquad \wedge \ vDeadlineQuorum' = [$

$\qquad\qquad vDeadlineQuorum \text{ EXCEPT } ![mySeverId]$

$\qquad\qquad\quad = \{msg \in vDeadlineQuorum[mySeverId] :$

$\qquad\qquad\qquad msg.txnId \notin syncedTxnIds\}$

$\qquad\qquad ]$

$\qquad \wedge \ vLSyncPoint' = [$

$\qquad\qquad vLSyncPoint \text{ EXCEPT } ![mySeverId] = Len(m.entries)]$

Send slow-replies to coordinators

$\qquad \wedge \ Send(\{[$

$\qquad\qquad mtype \quad \mapsto MSlowReply,$

$\qquad\qquad sender \quad \mapsto mySeverId,$

$\qquad\qquad dest \qquad \mapsto m.entries[i].txnId.coordId,$

$\qquad\qquad gView \quad \mapsto vGView[mySeverId],$

$\qquad\qquad lView \quad\;\; \mapsto vLView[mySeverId],$

$\qquad\qquad txnId \quad\; \mapsto m.entries[i].txnId,$

$\qquad\qquad logId \quad\;\; \mapsto i$

$\qquad\quad ] : i \in (currentSyncPoint + 1) .. Len(m.entries)\})$

$\vee \quad \wedge \ currentSyncPoint \geq Len(m.entries)$

Noting new to sync

$\qquad \wedge \ \text{UNCHANGED } \langle networkVars, vLog, vEarlyBuffer,$

$\qquad\qquad\qquad\qquad vLateBuffer, vDeadlineQuorum, vLSyncPoint \rangle$

$StartLeaderFail(serverId) \triangleq$
  This leader fails
  LET
    $serversInOneShard \triangleq \{$
      $s \in Servers : s.shardId = serverId.shardId$
    $\}$
    $aliveReplicas \triangleq \{$
      $s \in serversInOneShard : \quad \wedge \;\; vServerStatus[s] = StNormal$
      $\qquad\qquad\qquad\qquad\quad \wedge \;\; s \neq serverId$
    $\}$
  IN
    if the current alive replicas are less than $QuorumSize$
    Then no more replicas in this sharding group can fail (by assumption of consensus)
  IF $Cardinality(aliveReplicas) > QuorumSize$ THEN
    $vServerStatus' = [vServerStatus \text{ EXCEPT } ![serverId] = StFailing]$
  ELSE $\quad$ UNCHANGED $\langle vServerStatus \rangle$

$DetectLeaderFail(cmReplicaId) \triangleq$
  $\exists\, shardId \in Shards :$
    LET
      $lView \triangleq vCMGInfo[cmReplicaId].gVec[shardId]$
      $leaderId \triangleq LeaderID(lView)$
      $serverId \triangleq [$
        $replicaId \;\; \mapsto leaderId,$
        $shardId \;\;\;\; \mapsto shardId$
      $]$
    IN
    $vServerStatus[serverId] = StFailing$

$SelectProperLView(currentView,\, shardId) \triangleq$
  LET
    $aliveReplicaId \triangleq$ CHOOSE $replicaId \in Replicas :$
      $\qquad\qquad\qquad\qquad vServerStatus[shardId][replicaId] = StNormal$
  IN
    Ensure 1 the new view is larger than $currentView$
    $* (2)$ its corresponding leader happens to be the selected $aliveReplicaId$
      $(currentView \div Cardinality(Replicas) + 1) * Cardinality(Replicas) + aliveReplicaId$

$PrepareViewChange(cmReplicaId) \triangleq$
  LET
    $newGVec \triangleq [$
      $shardId \in Shards \mapsto$
        $SelectProperLView(vCMGInfo[cmReplicaId].gVec[shardId],\, shardId)$
    $]$
  IN

$\wedge\ vCMPrepareGInfo' = [vCMPrepareGInfo \text{ EXCEPT } ![cmReplicaId] =$
$$[$$
$$gView \quad \mapsto vCMGInfo[cmReplicaId].gView + 1,$$
$$gVec \quad \mapsto\ newGVec$$
$$]$$
$$]$$
$\wedge\ Send(\{[$
$$mtype \quad \mapsto MCMPrepare,$$
$$sender \quad \mapsto cmReplicaId,$$
$$dest \quad \mapsto dstRid,$$
$$cView \quad \mapsto vCMView[cmReplicaId],$$
$$gView \quad \mapsto vCMPrepareGInfo'[cmReplicaId].gView,$$
$$gVec \quad \mapsto\ newGVec$$

$$]: dstRid \in Replicas\})$$

$LaunchViewChange(cmReplicaId)\ \triangleq$
$\quad$ IF $\quad \wedge\ isLeader(cmReplicaId,\ vCMView[cmReplicaId])$
$\quad\quad\quad \wedge\ DetectLeaderFail(cmReplicaId)$
$\quad$ THEN
$\quad\quad PrepareViewChange(cmReplicaId)$
$\quad$ ELSE
$\quad\quad$ UNCHANGED $\quad \langle networkVars \rangle$

$HandleCMPrepare(m)\ \triangleq$
$\quad \wedge\ m.cView = vCMView[m.dest]$
$\quad \wedge\ m.gView > vCMGInfo[m.dest].gView$
$\quad \wedge\ vCMPrepareGInfo' = [vCMPrepareGInfo \text{ EXCEPT } ![m.dest] =$
$$[$$
$$gView \quad \mapsto m.gView,$$
$$gVec \quad \mapsto m.gVec$$
$$]$$
$$]$$
$\quad \wedge\ Send(\{[$
$$mtype \quad \mapsto MCMPrepareReply,$$
$$sender \quad \mapsto m.dest,$$
$$dest \quad \mapsto m.src,$$
$$cView \quad \mapsto m.cView,$$
$$gView \quad \mapsto m.gView$$
$$]\})$$

$HandleCMPrepareReply(m)\ \triangleq$
$\quad \wedge\ m.cView = vCMView[m.dest]$

$\wedge\ isLeader(m.dest,\ vCMView[m.dest])$
$\wedge\ m.gView = vCMPrepareGInfo[m.dest].gView$
$\wedge\ vCMPrepareReps' = [vCMPrepareReps \text{ EXCEPT } ![m.dest] =$
$\qquad vCMPrepareReps[m.dest] \cup \{m\}$
$\quad ]$
$\wedge$ LET
$\qquad quorum\ \triangleq\ \{mm \in vCMPrepareReps[m.dest] : mm.gView = m.gView\}$
$\quad$ IN
$\quad$ IF $Cardinality(quorum) = QuorumSize$ THEN
$\qquad$ Quorum sufficient, the prepared *GInfo* is persisted and can be safely used
$\qquad \wedge\ vCMGInfo' = [vCMGInfo \text{ EXCEPT } ![m.dest] =$
$\qquad\qquad vCMPrepareGInfo[m.dest]$
$\qquad\quad ]$
$\qquad$ notify other follower *CM*, so that they can catch up with the leader
$\qquad \wedge\ Send(\{[$
$\qquad\qquad mtype\quad \mapsto MCMCommit,$
$\qquad\qquad sender\quad \mapsto m.dest,$
$\qquad\qquad dest\qquad \mapsto rid,$
$\qquad\qquad cView\quad \mapsto m.cView,$
$\qquad\qquad gView\quad \mapsto m.gView$
$\qquad\quad ] : rid \in \{r \in Replicas : r \neq m.dest\}\})$
$\qquad$ start view change, broadcast view change request to every server
$\qquad \wedge\ Send(\{[$
$\qquad\qquad mtype\quad \mapsto MViewChangeReq,$
$\qquad\qquad sender\quad \mapsto m.dest,$
$\qquad\qquad dest\qquad \mapsto serverId,$
$\qquad\qquad gView\quad \mapsto vCMGInfo'[m.dest].gView,$
$\qquad\qquad gVec\qquad \mapsto vCMGInfo'[m.dest].gVec$
$\qquad\quad ] : serverId \in Servers\})$
$\quad$ ELSE
$\qquad$ UNCHANGED $\langle networkVars,\ vCMGInfo \rangle$

$HandleCMCommit(m)\ \triangleq$
$\quad \wedge\ m.cView = vCMView[m.dest]$
$\quad \wedge\ \neg isLeader(m.dest,\ vCMView[m.dest])$
$\quad \wedge\ m.gView = vCMPrepareGInfo[m.dest].gView$
$\quad \wedge\ vCMGInfo' = [vCMGInfo \text{ EXCEPT } ![m.dest] =$
$\qquad\qquad\qquad vCMPrepareGInfo[m.dest]$
$\qquad ]$


$HandleViewChangeReq(m)\ \triangleq$
$\quad$ LET
$\qquad myServerId\ \triangleq\ m.dest$
$\qquad myLeader\ \triangleq\ $ CHOOSE $s \in Servers :$

18

$$\wedge\ s.replicaId = LeaderID(m.gVec[myServerId.shardId])$$
$$\wedge\ s.shardId = myServerId.shardId$$

IN

If the *msg*'s view is lower, ignore

$\wedge\ vGView[myServerId] < m.gView$

$\wedge$ IF $vServerStatus[myServerId] = StNormal$ THEN
$\qquad\wedge\ vServerStatus'\quad = [vServerStatus$ EXCEPT $![myServerId]\quad = StViewChange]$
$\qquad\wedge\ vLastNormView' = [vLastNormView$ EXCEPT $![myServerId] = vLView[myServerId]]$
$\quad$ELSE$\qquad$UNCHANGED $\langle vServerStatus, vLastNormView\rangle$

$\wedge\ vGView' = [$
$\qquad vGView$ EXCEPT $![myServerId] = m.vGView$
$\quad]$

$\wedge\ vGVec' = [$
$\qquad vGVec$ EXCEPT $![myServerId] = m.gVec$
$\quad]$

$\wedge\ vLView' = [$
$\qquad vLView$ EXCEPT $![myServerId] = m.gVec[myServerId.shardId]$
$\quad]$

Clear ealry buffer,

$\wedge\ vEarlyBuffer' = [$
$\qquad vEarlyBuffer$ EXCEPT $![myServerId] = \{\}$
$\quad]$

Clear late buffer

$\wedge\ vLateBuffer' = [$
$\qquad vLateBuffer$ EXCEPT $![myServerId] = \{\}$
$\quad]$

Clear deadline quorum

$\wedge\ vDeadlineQuorum' = [$
$\qquad vDeadlineQuorum$ EXCEPT $![myServerId] = \{\}$
$\quad]$

Clear *vCrossShardConfirmQuorum*

$\wedge\ vCrossShardConfirmQuorum' = [$
$\qquad serverId \in Servers \mapsto \{\}$
$\quad]$

Send *ViewChange* to the *myLeader*

$\wedge\ Send(\{[$
$\qquad\qquad mtype\qquad\quad \mapsto MViewChange,$
$\qquad\qquad sender\qquad\quad \mapsto myServerId,$
$\qquad\qquad dest\qquad\quad\ \ \mapsto myLeader,$
$\qquad\qquad gView\qquad\quad \mapsto m.vGView,$
$\qquad\qquad gVec\qquad\quad\ \ \mapsto m.gVec,$
$\qquad\qquad lView\qquad\quad \mapsto vLView'[myServerId],$
$\qquad\qquad lastNormal\ \mapsto vLastNormView'[myServerId],$
$\qquad\qquad lSyncPoint\ \ \mapsto vLSyncPoint[myServerId],$
$\qquad\qquad entries\qquad\ \ \mapsto vLog[myServerId],$

$$cv \qquad\qquad \mapsto vCrashVector[myServerId]$$
$$]\})$$

Define a comparison function based on the key

$Compare(a,\ b) \triangleq$
 $\lor\quad a.deadline < b.deadline$
 $\lor\quad \land\ a.deadline = b.deadline$
   $\land\ a.txnId.coordId < b.txnId.coordId$
 $\lor\quad \land\ a.deadline = b.deadline$
   $\land\ a.txnId.coordId = b.txnId.coordId$
   $\land\ a.txnId.rId < b.txnId.rId$

$isCrashVectorValid(m) \triangleq$
 $\land\ \forall\, rr \in Replicas : vCrashVector[m.dest][rr] \leq m.cv[rr]$
 $\land\ vCrashVector' = [$
   $vCrashVector \text{ EXCEPT } ![m.dest] = [$
    $rr \in Replicas \mapsto Max(m.cv[rr], vCrashVector[m.dest][rr])$
   $]$
  $]$

$CountVotes(entry,\ logSets) \triangleq$
 LET
  $validCandidates \triangleq \{s \in logSets : \exists\, e \in s :$
        $\land\ e.deadline = entry.deadline$
        $\land\ e.txnId = entry.txnId\}$
 IN
  $Cardinality(validCandidates)$

$ReBuildLogs(vcQuorum) \triangleq$
 LET
  $refinedQuorum \triangleq \{m \in vcQuorum :$
         $\forall\, msg \in vcQuorum : msg.lastNormal \leq m.lastNormal\}$
  $lSyncPoints \triangleq \{m.lSyncPoint : m \quad \in refinedQuorum\}$
  $largestLSyncPointVC \triangleq \text{CHOOSE } vc \in refinedQuorum :$
        $\forall\, sp \in lSyncPoints : sp \leq vc.lSyncPoint$
  $syncedLogSeq \triangleq SubSeq(largestLSyncPointVC.entries, 1, largestLSyncPointVC.lSyncPoint)$
  $deadlineBoundary \triangleq \text{IF } largestLSyncPointVC.lSyncPoint = 0 \text{ THEN } 0$
        $\text{ELSE } syncedLogSeq[largestLSyncPointVC.lSyncPoint].deadline$
  $logSets \triangleq \{SeqToSet(m.entries) : m \in refinedQuorum\}$
  $allLogs \triangleq \text{UNION } logSets$
  $allUnSyncedLogs \triangleq \{entry \in allLogs : entry.deadline > deadlineBoundary\}$
  $unSyncedLogs \triangleq \{entry \in allUnSyncedLogs :$
   $CountVotes(entry, logSets) \geq RecoveryQuorumSize\}$
  $unSyncedLogSeq \triangleq SetToSortSeq(unSyncedLogs, Compare)$

IN
$syncedLogSeq \circ unSyncedLogSeq$

$SelectEntriesBeyondCommitPoint(entries, deadline) \triangleq$
    LET
        $validLogIndices \triangleq \{$
            $i \in 1 .. Len(entries) : entries[i].deadline > deadline$
        $\}$
        $startIndex \triangleq PickMin(validLogIndices)$
    IN
    $SubSeq(entries, startIndex, Len(entries))$

$HandleViewChange(m) \triangleq$
    LET
        $myServerId \triangleq m.dest$
        $serversInOneShard \triangleq \{s \in Servers : s.shardId = myServerId.shardId\}$
        $leadersInAllShard \triangleq \{$
            $s \in Servers : s.replicaId = isLeader(s.replicaId, m.gVec[s.shardId])$
        $\}$
    IN
    $\wedge \quad \vee \quad ViewGreater(m.gView, m.lView, vGView[myServerId], vLView[myServerId])$
          $\vee \quad \wedge \quad m.gView = vGView[myServerId]$
              $\wedge \quad m.lView = vLView[myServerId]$
              $\wedge \quad vServerStatus[myServerId] = StViewChange$
    $\wedge \quad isLeader(myServerId.replicaId, m.lView)$
    $\wedge \quad vGView' = [vGView \text{ EXCEPT } ![myServerId] = m.gView]$
    $\wedge \quad vLView' = [vLView \text{ EXCEPT } ![myServerId] = m.lView]$
    $\wedge \quad vGVec' = [vGVec \text{ EXCEPT } ![myServerId] = m.gVec]$
    $\wedge \quad vViewChange' = [$
        $vViewChange \text{ EXCEPT } ![myServerId] = \{$
            $vc \in vViewChange[myServerId] :$
               $vc.lView = m.lView$
        $\} \cup \{m\}$
    $]$
    $\wedge \quad$ IF $Cardinality(vViewChange'[myServerId]) = QuorumSize$ THEN
        $\wedge \quad vLog' = [vLog \text{ EXCEPT } ![myServerId] = ReBuildLogs(vViewChange'[myServerId])]$
        $\wedge \quad vServerStatus' = [vServerStatus \text{ EXCEPT } ![myServerId] = StCrossShardSyncing]$
        $\wedge \quad vLastNormView' = [vLastNormView \text{ EXCEPT } ![myServerId] = vLView[myServerId]]$
      Even after the $log$ is recovered within one shard,
        ∗ The newly elected leader cannot $StartView$
        ∗ It needs to sync with other $shards'$ leaders to ensure strict serializability
        $\wedge \quad vViewChange' = [vViewChange \text{ EXCEPT } ![myServerId] = \{\}]$
        $\wedge \quad Send(\{[$
            $mtype \quad \mapsto MCrossShardConfirm,$
            $sender \quad \mapsto myServerId,$

$$
\begin{aligned}
&dest &&\mapsto dst, \\
&lView &&\mapsto vLView'[myServerId], \\
&gView &&\mapsto vGView'[myServerId], \\
&entries &&\mapsto SelectEntriesBeyondCommitPoint( \\
& && \quad vLog'[myServerId], vPeerCommitDeadline[dst.shardId])
\end{aligned}
$$

$\qquad\qquad\quad ] : dst \in leadersInAllShard\})$

$\qquad\quad$ ELSE

$\qquad\qquad \wedge\ vServerStatus' = [vServerStatus \text{ EXCEPT } ![myServerId] = StViewChange]$

$\qquad\qquad \wedge\ $ UNCHANGED $\langle networkVars, vLog, vServerStatus, vViewChange \rangle$

$BuildGlobalConsistentLog(serverId, entries) \triangleq$

$\quad$ LET

$\qquad myEntries \triangleq \{$

$\qquad\quad entry \in entries : \wedge serverId \in entry.shards$

$\qquad\qquad\qquad\qquad\quad \wedge \forall\, e \in entries :$

$\qquad\qquad\qquad\qquad\qquad\quad$ IF $e.txnId = entry.txnId$ THEN

$\qquad\qquad\qquad\qquad\qquad\qquad\quad e.deadline \leq entry.deadline$

$\qquad\qquad\qquad\qquad\qquad\quad$ ELSE TRUE

$\qquad \}$

$\quad$ IN

$\quad SetToSortSeq(myEntries, Compare)$

$HandleCrossShardConfirm(m) \triangleq$

$\quad$ LET

$\qquad myServerId \triangleq m.dest$

$\quad$ IN

$\quad \wedge\ vServerStatus[myServerId] = StCrossShardSyncing$

$\quad \wedge\ m.gView = vGView[myServerId]$

$\quad \wedge\ m.lView = vGVec[myServerId][m.sender.shardId]$

$\quad \wedge\ vCrossShardConfirmQuorum' = [$

$\qquad vCrossShardConfirmQuorum \text{ EXCEPT } ![myServerId] = \{$

$\qquad\quad mm \in vCrossShardConfirmQuorum[myServerId] :$

$\qquad\qquad \wedge\ mm.gView = vGView[myServerId]$

$\qquad\qquad \wedge\ mm.lView = vGVec[myServerId][mm.sender.shardId]$

$\qquad \} \cup \{m\}$

$\quad ]$

$\quad \wedge\ $ IF $Cardinality(vCrossShardConfirmQuorum'[myServerId]) = Cardinality(Shards)$

$\qquad$ THEN

$\qquad\qquad$ Check $Txns'$ Deadlines to ensure strict serializability is not violated

$\qquad\qquad$ In implementation, we should not pass all txns, instead, we should only pass dealines and $txn$ indices

$\qquad\qquad$ As an optimization, we should also use checkpoint in implementation

$\qquad\qquad$ Here for conciseness, we pass all $log$ entries

$\qquad\qquad$ LET

$\qquad\qquad\quad allLogs \triangleq$ UNION $\{SeqToSet(mm.entries) :$

$\qquad\qquad\qquad\qquad\qquad\qquad mm \in vCrossShardConfirmQuorum'[myServerId]\}$

$$serversInOneShard \;\triangleq\; \{s \in Servers : s.shardId = myServerId.shardId\}$$

IN
$\land\;\; vLog' = [$
    $vLog \text{ EXCEPT } ![myServerId] =$
        $BuildGlobalConsistentLog(m.sender, allLogs)$
$]$

$\land\;\; Send(\{[$

| | | |
|---|---|---|
| $mtype$ | $\mapsto$ | $MStartView,$ |
| $sender$ | $\mapsto$ | $myServerId,$ |
| $dest$ | $\mapsto$ | $dst,$ |
| $lView$ | $\mapsto$ | $vLView[myServerId],$ |
| $gView$ | $\mapsto$ | $vGView[myServerId],$ |
| $gVec$ | $\mapsto$ | $vGVec[myServerId],$ |
| $entries$ | $\mapsto$ | $vLog'[myServerId],$ |
| $cv$ | $\mapsto$ | $vCrashVector[myServerId]$ |

    $] : dst \in serversInOneShard\})$

ELSE
    UNCHANGED $\langle vLog, networkVars \rangle$

$HandleStartView(m) \;\triangleq\;$
  LET
    $myServerId \;\triangleq\; m.dest$
  IN
  $\land\;\; \lor\;\; ViewGreater(m.gView, m.lView, vGView[myServerId], vLView[myServerId])$
      $\lor\;\; \land\;\; m.gView = vGView[myServerId]$
          $\land\;\; m.lView\; = vLView[myServerId]$
          $\land\;\; \lor\;\; vServerStatus[myServerId] = StViewChange$
              $\lor\;\; vServerStatus[myServerId] = StRecovering$
  $\land\;\; vGView' = [vGView \text{ EXCEPT } ![myServerId] = m.gView]$
  $\land\;\; vLView' = [vLView \text{ EXCEPT } ![myServerId] = m.gLView]$
  $\land\;\; vGVec' = [vGVec \text{ EXCEPT } ![myServerId] = m.vGVec]$
  $\land\;\; vServerStatus' = [vServerStatus \text{ EXCEPT } ![myServerId] = StNormal]$
  $\land\;\; vLog' = [vLog \text{ EXCEPT } ![myServerId] = m.entries]$
  $\land\;\; vEarlyBuffer' = [vEarlyBuffer \text{ EXCEPT } ![myServerId] = \{\}]$
  $\land\;\; vLateBuffer' = [vLateBuffer \text{ EXCEPT } ![myServerId] = \{\}]$
  $\land\;\; vDeadlineQuorum' = [vDeadlineQuorum \text{ EXCEPT } ![myServerId] = \{\}]$
  $\land\;\; vCrossShardConfirmQuorum' = [$
    $vCrossShardConfirmQuorum \text{ EXCEPT } ![myServerId] = \{\}$
  $]$
  $\land\;\; vLSyncPoint' = [vLSyncPoint \text{ EXCEPT } ![myServerId] = Len(vLog'[myServerId])]$
  $\land\;\; vLastNormView' = [vLastNormView \text{ EXCEPT } ![myServerId] = m.lView]$
  $\land\;\; vViewChange' = [vViewChange \text{ EXCEPT } ![myServerId] = \{\}]$
  $\land\;\; vLSyncQuorum' = [vLSyncQuorum \text{ EXCEPT } ![myServerId] = \{\}]$
  $\land\;\; vCrashVectorReps' = [vCrashVectorReps \text{ EXCEPT } ![myServerId] = \{\}]$
  $\land\;\; vRecoveryReps' = [vRecoveryReps \text{ EXCEPT } ![myServerId] = \{\}]$

$ResetServerState(serverId) \triangleq$
    $\land\ vLog' = [vLog \text{ EXCEPT } ![serverId] = \langle\rangle]$
    $\land\ vEarlyBuffer' = [vEarlyBuffer \text{ EXCEPT } ![serverId] = \{\}]$
    $\land\ vLateBuffer' = [vLateBuffer \text{ EXCEPT } ![serverId] = \{\}]$
    $\land\ vDeadlineQuorum' = [vDeadlineQuorum \text{ EXCEPT } ![serverId] = \{\}]$
    $\land\ vCrossShardConfirmQuorum' = [$
        $vCrossShardConfirmQuorum \text{ EXCEPT } ![serverId] = \{\}$
    $]$
    $\land\ vGView' = [vGView \text{ EXCEPT } ![serverId] = 0]$
    $\land\ vGVec' = [vGVec \text{ EXCEPT } ![serverId] = [s \in Shards \mapsto 0]]$
    $\land\ vLView' = [vLView \text{ EXCEPT } ![serverId] = 0]$
    $\land\ vLastNormView' = [vLastNormView \text{ EXCEPT } ![serverId] = 0]$
    $\land\ vViewChange' = [vViewChange \text{ EXCEPT } ![serverId] = \{\}]$
    $\land\ vLSyncPoint' = [vLSyncPoint \text{ EXCEPT } ![serverId] = 0]$
    $\land\ vLCommitPoint' = [vLCommitPoint \text{ EXCEPT } ![serverId] = 0]$
    $\land\ vPeerCommitDeadline' = [vPeerCommitDeadline \text{ EXCEPT } ![serverId] = 0]$
    $\land\ vLSyncQuorum' = [vLSyncQuorum \text{ EXCEPT } ![serverId] = \{\}]$
    $\land\ vCrashVector' = [vCrashVector \text{ EXCEPT } ![serverId] = [$
        $rr \in Replicas \mapsto 0$
    $]]$
    $\land\ vCrashVectorReps' = [vCrashVectorReps \text{ EXCEPT } ![serverId] = \{\}]$
    $\land\ vRecoveryReps' = [vRecoveryReps \text{ EXCEPT } ![serverId] = \{\}]$
    $\land\ vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![serverId] = \{\}]$

$StartServerRecovery(serverId) \triangleq$
    LET
        $serversInOneShard \triangleq \{$
            $s \in Servers : s.shardId = serverId.shardId$
        $\}$
        $nonce \triangleq\ vUUIDCounter[serverId] + 1$
    IN
    $\land\ vServerStatus' = [vServerStatus \text{ EXCEPT } ![serverId] = StRecovering]$
    $\land\ vUUIDCounter' = [vUUIDCounter \text{ EXCEPT } ![serverId] = vUUIDCounter[serverId] + 1]$
    $\land\ ResetServerState(serverId)$
    $\land\ Send(\{[$
        $mtype\quad \mapsto MCrashVectorReq,$
        $sender\quad \mapsto serverId,$
        $dest\quad \mapsto dst,$
        $nonce\quad \mapsto nonce$
    $] : dst \in serversInOneShard\})$

$HandleCrashVectorReq(m) \triangleq$
    LET
        $myServerId \triangleq\ m.dest$
    IN

$\wedge$ $vServerStatus[myServerId] = StNormal$
$\wedge$ $Send(\{[$

| | |
|---|---|
| $mtype$ | $\mapsto MCrashVectorRep,$ |
| $sender$ | $\mapsto myServerId,$ |
| $dest$ | $\mapsto m.sender,$ |
| $nonce$ | $\mapsto m.nonce,$ |
| $cv$ | $\mapsto vCrashVector[myServerId]$ |

$]\})$

$AggregateCV(serverId) \triangleq$
    LET
        $cvQuorum \triangleq \{m.cv : m \in vCrashVectorReps[serverId]\}$
        $cvValQuorum \triangleq [rr \in Replicas \mapsto \{cv[rr] : cv \in cvQuorum\}]$
    IN
        $[rr \in Replicas \mapsto PickMax(cvValQuorum[rr])]$

$HandleCrashVectorRep(m) \triangleq$
    LET
        $myServerId \triangleq m.dest$
        $serversInOneShard \triangleq \{s \in Servers : s.shardId = myServerId.shardId\}$
    IN
    $\wedge$ $vServerStatus[myServerId] = StRecovering$
    $\wedge$ $vUUIDCounter[myServerId] = m.nonce$
    $\wedge$ $vCrashVectorReps' = [$
        $vCrashVectorReps$ EXCEPT $![myServerId] = vCrashVectorReps \cup \{m\}$
    $]$
    $\wedge$ IF $Cardinality(vCrashVectorReps'[myServerId]) = QuorumSize$ THEN
        LET
            $acv \triangleq AggregateCV(myServerId)$
            $myCV \triangleq [acv$ EXCEPT $![myServerId] = acv[myServerId] + 1]$
        IN
        $\wedge$ $vCrashVector' = [$
            $vCrashVector$ EXCEPT $![myServerId] = myCV$
        $]$
        $\wedge$ $Send(\{[$

| | |
|---|---|
| $mtype$ | $\mapsto MRecoveryReq,$ |
| $sender$ | $\mapsto myServerId,$ |
| $dest$ | $\mapsto dst,$ |
| $cv$ | $\mapsto myCV$ |

        $] : dst \in serversInOneShard\})$
    ELSE    UNCHANGED $\langle networkVars, vCrashVector \rangle$

$HandleRecoveryReq(m) \triangleq$

$$\text{LET}$$
$$myServerId \triangleq m.dest$$
$$\text{IN}$$
$$\wedge \ vServerStatus[myServerId] = StNormal$$
$$\wedge \ Send(\{[$$
$$mtype \quad \mapsto MRecoveryRep,$$
$$sender \quad \mapsto myServerId,$$
$$dest \quad \mapsto m.sender,$$
$$gView \quad \mapsto vGView[myServerId],$$
$$lView \quad \mapsto vLView[myServerId],$$
$$cv \quad \mapsto vCrashVector'[myServerId]$$
$$]\})$$

$HandleRecoveryRep(m) \triangleq$
$$\text{LET}$$
$$myServerId \triangleq m.dest$$
$$\text{IN}$$
$$\wedge \ vServerStatus[myServerId] = StRecovering$$
$$\wedge \ vRecoveryReps' = [$$
$$vRecoveryReps \ \text{EXCEPT} \ ![myServerId]$$
$$= vRecoveryReps[myServerId] \cup \{m\}$$
$$]$$
$$\wedge \ \text{IF} \ Cardinality(vRecoveryReps[myServerId]) = QuorumSize \ \text{THEN}$$
$$\text{LET}$$
$$lViewQuorum \triangleq \{mm.lView : mm \in vRecoveryReps[myServerId]\}$$
$$gViewQuorum \triangleq \{mm.gView : mm \in vRecoveryReps[myServerId]\}$$
$$\text{IN}$$
$$\wedge \ vLView' = [vLView \ \text{EXCEPT} \ ![myServerId] = PickMax(lViewQuorum)]$$
$$\wedge \ vGView' = [vLView \ \text{EXCEPT} \ ![myServerId] = PickMax(gViewQuorum)]$$
$$\wedge \ Send(\{[$$
$$mtype \quad \mapsto MStartViewReq,$$
$$sender \quad \mapsto myServerId,$$
$$dest \quad \mapsto [$$
$$replicaId \mapsto LeaderID(vLView[myServerId]),$$
$$shardId \ \mapsto myServerId.shardId$$
$$],$$
$$lView \quad \mapsto vLView'[myServerId],$$
$$cv \quad \mapsto vCrashVector'[myServerId]$$
$$]\})$$
$$\text{ELSE} \quad \text{UNCHANGED} \quad \langle networkVars, vLView, vGView \rangle$$

$HandleStartViewReq(m) \triangleq$
$$\text{LET}$$

$$myServerId \triangleq m.dest$$

IN

$\land\ vServerStatus[myServerId] = StNormal$

$\land\ vLView[myServerId] = m.lView$

$\land\ isLeader(myServerId.replicaId, vLView[myServerId])$

$\land\ Send(\{[$

$$
\begin{aligned}
mtype\ &\mapsto MStartView, \\
sender\ &\mapsto myServerId, \\
dest\ &\mapsto m.sender, \\
lView\ &\mapsto vLView[myServerId], \\
gView\ &\mapsto vGView[myServerId], \\
gVec\ &\mapsto vGVec[myServerId], \\
entries\ &\mapsto vLog[myServerId], \\
cv\ &\mapsto vCrashVector[myServerId]
\end{aligned}
$$

$]\})$

$StartLocalSync(serverId) \triangleq$

LET

$leaderServerId \triangleq\ [$

$replicaId \mapsto LeaderID(vLView[serverId]),$

$shardId\ \mapsto serverId.shardId$

$]$

IN

$\land\ vServerStatus[serverId] = StNormal$

$\land\ Send(\{[$

$$
\begin{aligned}
mtype\ &\mapsto MLocalSyncStatus, \\
sender\ &\mapsto serverId, \\
dest\ &\mapsto leaderServerId, \\
lView\ &\mapsto vLView[serverId], \\
lSyncPoint\ &\mapsto vLSyncPoint[serverId], \\
cv\ &\mapsto vCrashVector[serverId]
\end{aligned}
$$

$]\})$

$HandleLocalSyncStatus(m) \triangleq$

LET

$myServerId \triangleq m.dest$

$lSyncQuorum \triangleq vLSyncQuorum[myServerId]$

IN

$\land\ vServerStatus[myServerId] = StNormal$

$\land\ vLView[myServerId] = m.lView$

$\land\ isLeader(myServerId.replicaId, vLView[myServerId])$

$\land\ \forall mm \in lSyncQuorum :$

$\lor\ mm.sender \neq m.sender$

27

$$\lor \ mm.lSyncPoint < m.lSyncPoint$$
$$\land \ vLSyncQuorum' = [$$
$$vLSyncQuorum \text{ EXCEPT } ![myServerId] =$$
$$\{mm \in lSyncQuorum : mm.sender \neq m.sender\} \cup \{m\}$$
$$]$$
$\land$ IF $Cardinality(vLSyncQuorum'[myServerId]) \geq QuorumSize$ THEN

 LET
$$candidateQuorum \ \triangleq \ \{$$
$$R \in \text{SUBSET } (vLSyncQuorum'[myServerId]) :$$
$$Cardinality(R) = QuorumSize$$
$$\}$$
$$quorumSyncPoints \ \triangleq \ \{$$
$$\{x.lSyncPoint : x \in R\} : R \in candidateQuorum$$
$$\}$$
$$validCommitPoints \ \triangleq \ \{PickMax(Q) : Q \in quorumSyncPoints\}$$
$$maxCommitPoint \ \triangleq \ PickMax(validCommitPoints)$$

 IN
$\land \ vLCommitPoint' = [vLCommitPoint \text{ EXCEPT } ![myServerId] = maxCommitPoint]$
$\land \ Send(\{[$

| | |
|---|---|
| $mtype$ | $\mapsto MLocalCommit,$ |
| $sender$ | $\mapsto myServerId,$ |
| $dest$ | $\mapsto m.sender,$ |
| $lView$ | $\mapsto vLView[myServerId],$ |
| $lCommitPoint$ | $\mapsto vLCommitPoint'[myServerId],$ |
| $cv$ | $\mapsto vCrashVector'[myServerId]$ |

  $]\})$

 ELSE UNCHANGED $\langle vLCommitPoint, networkVars\rangle$

<br/>

$HandleLocalCommit(m) \ \triangleq$

 LET
$$myServerId \ \triangleq \ m.dest$$

 IN
$\land \ vServerStatus[myServerId] = StNormal$
$\land \ vLView[myServerId] = m.lView$
$\land \ \neg isLeader(myServerId.replicaId, vLView[myServerId])$

 Make sure the $syncPoint$ is large enough before updating $CommitPoint$
$\land$ IF $\ \land \ vLSyncPoint[myServerId] \geq m.lCommitPoint$
   $\land \ vLCommitPoint[myServerId] < m.lCommitPoint$

 THEN
$$vLCommitPoint' = [$$
$$vLCommitPoint \text{ EXCEPT } ![myServerId] = m.lCommitPoint$$
$$]$$

 ELSE UNCHANGED $\langle vLCommitPoint\rangle$

<br/><br/>

28

$BroadcastCommitStatusToPeers(serverId) \triangleq$
    LET
        $serversInOneReplica \triangleq \{s \in Servers : s.replicaId = serverId.replicaId\}$
        $commitPoint \triangleq vLCommitPoint[serverId]$
        $commitDeadline \triangleq$
            IF $commitPoint = 0$ THEN $0$
            ELSE  $vLog[commitPoint].deadline$
    IN
    $\wedge$  $vServerStatus[serverId] = StNormal$
    $\wedge$  $Send(\{[$
        $mtype$        $\mapsto MPeerShardCommitStatus,$
        $sender$      $\mapsto serverId,$
        $dest$        $\mapsto dst,$
        $gView$      $\mapsto vGView[serverId],$
        $lView$       $\mapsto vLView[serverId],$
        $deadline$    $\mapsto commitDeadline$
    $] : dst \in serversInOneReplica\})$

$HandlePeerShardCommitStatus(m) \triangleq$
    LET
        $myServerId \triangleq m.dest$
    IN
    $\wedge$  $vServerStatus[myServerId] = StNormal$
    $\wedge$  $vGView[myServerId] = m.gView$
    $\wedge$  $vGVec[myServerId][m.sender.shardId] = m.lView$
    $\wedge$  IF  $m.deadline > vPeerCommitDeadline[myServerId][m.sender.shardId]$ THEN
        $\wedge$  $vPeerCommitDeadline[myServerId]' = [$
            $vPeerCommitDeadline[myServerId]$
                EXCEPT $![m.sender.shardId] = m.deadline$
            $]$
        ELSE  UNCHANGED  $\langle vPeerCommitDeadline \rangle$

$isCommitting(txn, deadlineQ) \triangleq$
    LET $quorum \triangleq \{msg \in deadlineQ : msg.entry.txnId = txn.txnId\}$
    IN    $Cardinality(quorum) = Cardinality(txn.shards)$

$ReleaseSeqeuncer(serverId, currentTime) \triangleq$
    LET
        $serversInOneShard \triangleq \{s \in Servers : s.shardId = serverId.shardId\}$
        $expireTxns \triangleq$
            $\{msg \in vEarlyBuffer[serverId] :$
                $\wedge msg.deadline \leq currentTime\}$

$$sortedTxnList \;\triangleq\; SetToSortSeq(expireTxns,\ Compare)$$

$committingStatus \;\triangleq\;$
$\qquad [i \in 1 \,..\, Len(sortedTxnList)$
$\qquad\quad \mapsto isCommitting(sortedTxnList[i],\ vDeadlineQuorum[serverId])$
$\qquad ]$

$canReleaseTxnIndices \;\triangleq\; \{$
$\qquad i \in 1 \,..\, Len(sortedTxnList) :$
$\qquad\quad \forall j \in 1 \,..\, i : committingStatus[j] = \text{TRUE}\}$

IN

IF $\;Cardinality(canReleaseTxnIndices) = 0$   Nothing to release

THEN    UNCHANGED $\;\langle networkVars,$
$\qquad\qquad\qquad vLog,\ vEarlyBuffer,\ vLateBuffer,\ vDeadlineQuorum\rangle$

ELSE

LET

$\qquad releaseUpTo \;\triangleq\; \text{CHOOSE}\ i \in canReleaseTxnIndices :$
$\qquad\qquad\qquad\qquad\qquad \forall j \in canReleaseTxnIndices : j \le i$
$\qquad releaseSeq \;\triangleq\; SubSeq(sortedTxnList,\ 1,\ releaseUpTo)$
$\qquad releaseTxns \;\triangleq\; \{releaseSeq[i] : i \in 1 \,..\, Len(releaseSeq)\}$

IN

$\wedge\ vEarlyBuffer' = [$
$\quad vEarlyBuffer\ \text{EXCEPT}\ ![serverId]$
$\qquad = vEarlyBuffer[serverId] \setminus releaseTxns]$
$\wedge\ vDeadlineQuorum' = [$
$\quad vDeadlineQuorum\ \text{EXCEPT}\ ![serverId]$
$\qquad = \{msg \in vDeadlineQuorum[serverId] :$
$\qquad\qquad \forall txn \in releaseTxns : txn.txnId \ne msg.entry.txnId\}$
$\quad ]$

Append to $log$
$\wedge\ vLog' = [vLog\ \text{EXCEPT}\ ![serverId] = vLog[serverId] \circ releaseSeq]$
$\wedge\ \text{IF}\ isLeader(serverId.replicaId,\ vLView[serverId])\ \text{THEN}$
$\qquad\quad \wedge\ vLSyncPoint' = [vLSyncPoint\ \text{EXCEPT}\ ![serverId] = Len(vLog'[serverId])]$
$\qquad \text{ELSE} \qquad \text{UNCHANGED}\ \langle vLSyncPoint\rangle$

Send fast-replies to coordinators
$\wedge\ Send(\{[$
$\quad mtype \quad \mapsto MFastReply,$
$\quad sender \quad \mapsto serverId,$
$\quad dest \qquad \mapsto sortedTxnList[i].txnId.coordId,$
$\quad gView \quad \mapsto vGView[serverId],$
$\quad lView \quad \mapsto vLView[serverId],$
$\quad txnId \quad \mapsto sortedTxnList[i].txnId,$
$\quad hash \qquad \mapsto [$
$\qquad\qquad log \mapsto vLog'[serverId],$
$\qquad\qquad cv \mapsto vCrashVector$
$\qquad\quad ],$
$\quad logId \quad \mapsto i$

$] : i \in (1 + Len(vLog[serverId])) .. Len(vLog'[serverId])\})$

 Send *InterReplicaSync* to the other servers in the same sharding group

 In real implementation, we send the *log* indices incrementally (*i.e.*, consider it as an optimization)

 Here for clarity and simplicity, we always send the whole *log* list

 $\wedge\ Send(\{[$

$$
\begin{aligned}
mtype\ &\mapsto MInterReplicaSync,\\
lView\ &\mapsto vLView[serverId],\\
sender\ &\mapsto serverId,\\
dest\ &\mapsto dstServerId,\\
entries\ &\mapsto vLog'[serverId]
\end{aligned}
$$

$] : dstServerId \in serversInOneShard\})$

$ServerClockMove(serverId) \triangleq$
 IF $vServerClock[serverId] \geq MaxTime$  THEN
  UNCHANGED $\langle networkVars,\ serverStateVars \rangle$
 ELSE
  $\wedge\ \ vServerClock' = [$
   $vServerClock$ EXCEPT $![serverId] = vServerClock[serverId] + 1]$
  $\wedge\ $ IF $vServerStatus[serverId] = StNormal$ THEN
   $\wedge\ ReleaseSeqeuncer(serverId,\ vServerClock[serverId] + 1)$
   ELSE
   UNCHANGED $\langle networkVars,\ vLog,\ vEarlyBuffer,$
    $vLateBuffer,\ vDeadlineQuorum \rangle$
  $\wedge\ $ UNCHANGED $\langle vCrossShardConfirmQuorum,$
   $vServerStatus,\ vGView,\ vGVec,\ vLView,\ vLastNormView,$
   $vViewChange,\ vLSyncPoint,\ vLCommitPoint,$
   $vPeerCommitDeadline,\ vLSyncQuorum,$
   $vUUIDCounter,\ vCrashVector,\ vCrashVectorReps,$
   $vRecoveryReps,\ vServerProcessed \rangle$

$CoordClockMove(coordId) \triangleq$
 $\vee\ \ \wedge vCoordClock[coordId] \geq MaxTime$
  $\wedge$ UNCHANGED $\langle vCoordClock \rangle$
 $\vee\ \ \wedge vCoordClock[coordId] < MaxTime$
  $\wedge vCoordClock' = [$
   $vCoordClock$ EXCEPT $![coordId] = vCoordClock[coordId] + 1]$

$Init \triangleq$
 $\wedge InitNetworkState$
 $\wedge InitServerState$
 $\wedge InitCoordState$
 $\wedge InitConfigManagerState$
 $\wedge ActionName = \langle \text{"Init"} \rangle$

$Next \triangleq$
$\quad \lor \;\; \land ActionName' = \langle \text{``Next''} \rangle$
$\qquad \land \text{UNCHANGED} \;\langle networkVars, \, serverStateVars,$
$\qquad\qquad\qquad\qquad coordStateVars, \, configManagerStateVars \rangle$
$\quad \lor \;\; \exists\, c \in Coords :$
$\qquad \land Cardinality(vCoordTxns[c]) < MaxReqNum$
$\qquad \land CoordSubmitTxn(c)$
$\qquad \land \text{UNCHANGED} \;\langle serverStateVars, \, configManagerStateVars,$
$\qquad\qquad\qquad vCoordProcessed \rangle$
$\qquad \land ActionName' = \langle \text{``CoordSubmitTxn''} \rangle$

$\quad \lor \exists\, m \in messages :$
$\qquad \land m.mtype = MTxn$
$\qquad \land vServerStatus[m.dest] = StNormal$
$\qquad \land m \notin vServerProcessed[m.dest]$
$\qquad \land vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] =$
$\qquad\quad vServerProcessed[m.dest] \cup \{m\}]$
$\qquad \land HandleTxn(m)$
$\qquad \land \text{UNCHANGED} \;\; \langle coordStateVars, \, configManagerStateVars,$
$\qquad\quad vLog, \, vDeadlineQuorum, \, vCrossShardConfirmQuorum,$
$\qquad\quad vServerStatus, \, vGView, \, vGVec,$
$\qquad\quad vLView, \, vServerClock, \, vLastNormView,$
$\qquad\quad vViewChange, \, vLSyncPoint, \, vLCommitPoint,$
$\qquad\quad vPeerCommitDeadline, \, vLSyncQuorum,$
$\qquad\quad vUUIDCounter, \, vCrashVector,$
$\qquad\quad vCrashVectorReps, \, vRecoveryReps \rangle$
$\qquad \land ActionName' = \langle \text{``HandleTxn''} \rangle$

$\quad \lor \exists\, m \in messages :$
$\qquad \land m.mtype = MDeadlineNotification$
$\qquad \land vServerStatus[m.dest] = StNormal$
$\qquad \land m \notin vServerProcessed[m.dest]$
$\qquad \land vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] =$
$\qquad\quad vServerProcessed[m.dest] \cup \{m\}]$
$\qquad \land HandleDeadlineNotification(m)$
$\qquad \land \text{UNCHANGED} \;\; \langle networkVars, \, coordStateVars, \, configManagerStateVars,$
$\qquad\qquad vLog, \, vCrossShardConfirmQuorum, \, vLateBuffer,$
$\qquad\qquad vServerStatus, \, vGView, \, vGVec,$
$\qquad\qquad vLView, \, vServerClock, \, vLastNormView,$
$\qquad\qquad vViewChange, \, vLSyncPoint, \, vLCommitPoint,$
$\qquad\qquad vPeerCommitDeadline, \, vLSyncQuorum,$
$\qquad\qquad vUUIDCounter, \, vCrashVector, \, vCrashVectorReps,$
$\qquad\qquad vRecoveryReps \rangle$
$\qquad \land ActionName' = \langle \text{``HandleDeadlineNotification''} \rangle$

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

$\land m.mtype = MInterReplicaSync$
$\land vServerStatus[m.dest] = StNormal$
$\land m \notin vServerProcessed[m.dest]$
$\land vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] =$
$\quad vServerProcessed[m.dest] \cup \{m\}]$
$\land HandleInterReplicaSync(m)$
$\land \text{UNCHANGED } \langle coordStateVars, configManagerStateVars,$
$\qquad vLog, \quad vCrossShardConfirmQuorum, vLateBuffer,$
$\qquad vServerStatus, vGView, vGVec,$
$\qquad vLView, vServerClock, vLastNormView,$
$\qquad vViewChange, vLCommitPoint, vPeerCommitDeadline,$
$\qquad vLSyncQuorum, vUUIDCounter, vCrashVector,$
$\qquad vCrashVectorReps, vRecoveryReps \rangle$
$\land ActionName' = \langle \text{"HandleInterReplicaSync"} \rangle$

Some $Leader(s)$ fail
$\lor \exists serverId \in Servers :$
$\quad \land vLView[serverId] < MaxViews$
$\quad \land isLeader(serverId.replicaId, vLView[serverId])$
$\quad \land StartLeaderFail(serverId)$
$\quad \land \text{UNCHANGED } \langle networkVars, coordStateVars, configManagerStateVars,$
$\qquad vLog, vEarlyBuffer, vLateBuffer,$
$\qquad vDeadlineQuorum, vCrossShardConfirmQuorum, vGView, vGVec,$
$\qquad vLView, vServerClock, vLastNormView,$
$\qquad vViewChange, vLSyncPoint, vLCommitPoint,$
$\qquad vPeerCommitDeadline, vLSyncQuorum,$
$\qquad vUUIDCounter, vCrashVector, vCrashVectorReps,$
$\qquad vRecoveryReps, vServerProcessed \rangle$
$\quad \land ActionName' = \langle \text{"StartLeaderFail"} \rangle$

Config Manager notices some $leader(s)$ fail and launch view change
$\lor \quad \exists cmReplicaId \in Replicas :$
$\quad \land LaunchViewChange(cmReplicaId)$
$\quad \land \text{UNCHANGED } \langle coordStateVars, serverStateVars, configManagerStateVars \rangle$
$\quad \land ActionName' = \langle \text{"LaunchViewChange"} \rangle$

$\lor \exists m \in messages :$
$\quad \land m.mtype = MCMPrepare$
$\quad \land m \notin vCMProcessed[m.dest]$
$\quad \land vCMProcessed' = [vCMProcessed \text{ EXCEPT } ![m.dest] =$
$\qquad vCMProcessed[m.dest] \cup \{m\}]$
$\quad \land HandleCMPrepare(m)$
$\quad \land \text{UNCHANGED } \langle coordStateVars, serverStateVars \rangle$
$\quad \land ActionName' = \langle \text{"HandleCMPrepare"} \rangle$

$\lor\ \exists\,m \in messages:$
$\quad \land\ m.mtype = MCMPrepareReply$
$\quad \land\ m \notin vCMProcessed[m.dest]$
$\quad \land\ vCMProcessed' = [vCMProcessed \text{ EXCEPT } ![m.dest] =$
$\qquad\quad vCMProcessed[m.dest] \cup \{m\}]$
$\quad \land\ HandleCMPrepareReply(m)$
$\quad \land\ \text{UNCHANGED}\ \ \langle coordStateVars,\ serverStateVars,$
$\qquad\qquad\qquad vCMStatus,\ vCMView,\ vCMPrepareGInfo \rangle$
$\quad \land\ ActionName' = \langle \text{``HandleCMPrepareReply''} \rangle$

$\lor\ \exists\,m \in messages:$
$\quad \land\ m.mtype = MCMCommit$
$\quad \land\ m \notin vCMProcessed[m.dest]$
$\quad \land\ vCMProcessed' = [vCMProcessed \text{ EXCEPT } ![m.dest] =$
$\qquad\quad vCMProcessed[m.dest] \cup \{m\}]$
$\quad \land\ HandleCMCommit(m)$
$\quad \land\ \text{UNCHANGED}\ \ \langle networkVars,\ coordStateVars,\ serverStateVars,$
$\qquad\qquad\qquad vCMStatus,\ vCMView,\ vCMPrepareGInfo,\ vCMPrepareReps \rangle$
$\quad \land\ ActionName' = \langle \text{``HandleCMCommit''} \rangle$


$\lor\ \exists\,m \in messages:$
$\quad \land\ m.mtype = MViewChangeReq$
$\quad \land\ m \notin vServerProcessed[m.dest]$
$\quad \land\ vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] =$
$\qquad\quad vServerProcessed[m.dest] \cup \{m\}]$
$\quad \land\ vServerStatus[m.dest] \neq StFailing$
$\quad \land\ HandleViewChangeReq(m)$
$\quad \land\ \text{UNCHANGED}\ \ \langle coordStateVars,\ configManagerStateVars,$
$\qquad\quad vLog,\ vServerClock,\ vViewChange,\ vLSyncPoint,$
$\qquad\quad vLCommitPoint,\ vLSyncQuorum,\ vPeerCommitDeadline,$
$\qquad\quad vUUIDCounter,\ vCrashVector,\ vCrashVectorReps,\ vRecoveryReps \rangle$
$\quad \land\ ActionName' = \langle \text{``HandleViewChangeReq''} \rangle$


$\lor\ \exists\,m \in messages:$
$\quad \land\ m.mtype = MViewChange$
$\quad \land\ isCrashVectorValid(m)$
$\quad \land\ m \notin vServerProcessed[m.dest]$
$\quad \land\ vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] =$
$\qquad\quad vServerProcessed[m.dest] \cup \{m\}]$
$\quad \land\ vServerStatus[m.dest] \neq StFailing$
$\quad \land\ HandleViewChange(m)$
$\quad \land\ \text{UNCHANGED}\ \ \langle coordStateVars,\ configManagerStateVars,$
$\qquad\quad vGVec,\ vServerClock,\ vLSyncPoint,\ vLastNormView,$
$\qquad\quad vLCommitPoint,\ vPeerCommitDeadline,\ vLSyncQuorum,$

$$vUUIDCounter, vCrashVector, vCrashVectorReps,$$
$$vRecoveryReps\rangle$$
$$\wedge\ ActionName' = \langle\text{``HandleViewChange''}\rangle$$

$\vee\ \exists\, m \in messages :$
  $\wedge\ m.mtype = MCrossShardConfirm$
  $\wedge\ m \notin vServerProcessed[m.dest]$
  $\wedge\ vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] =$
       $vServerProcessed[m.dest] \cup \{m\}]$
  $\wedge\ vServerStatus[m.dest] = StViewChange$
  $\wedge\ HandleCrossShardConfirm(m)$
  $\wedge\ \text{UNCHANGED }\ \langle coordStateVars, configManagerStateVars,$
       $vGVec, vServerClock, vLSyncPoint, vLastNormView,$
       $vLCommitPoint, vPeerCommitDeadline, vLSyncQuorum,$
       $vUUIDCounter, vCrashVector, vCrashVectorReps, vRecoveryReps\rangle$
  $\wedge\ ActionName' = \langle\text{``HandleCrossShardConfirm''}\rangle$

$\vee\ \exists\, m \in messages :$
  $\wedge\ m.mtype = MStartView$
  $\wedge\ isCrashVectorValid(m)$
  $\wedge\ m \notin vServerProcessed[m.dest]$
  $\wedge\ vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] =$
       $vServerProcessed[m.dest] \cup \{m\}]$
  $\wedge\ HandleStartView(m)$
  $\wedge\ \text{UNCHANGED }\ \langle coordStateVars, configManagerStateVars,$
       $vServerClock, vLCommitPoint, vPeerCommitDeadline,$
       $vUUIDCounter, vCrashVector\rangle$
  $\wedge\ ActionName' = \langle\text{``HandleStartView''}\rangle$

Failed server rejoin
$\vee\ \exists\, serverId \in Servers :$
  $\wedge\ vServerStatus[serverId] = StFailing$
  $\wedge\ vServerStatus' = [vServerStatus \text{ EXCEPT } ![serverId] = StRecovering]$
  $\wedge\ ResetServerState(serverId)$
  $\wedge\ StartServerRecovery(serverId)$
  $\wedge\ \text{UNCHANGED }\ \langle networkVars, coordStateVars, coordStateVars\rangle$
  $\wedge\ ActionName' = \langle\text{``StartReplicaRecovery''}\rangle$

$\vee\ \exists\, m \in messages :$
  $\wedge\ m.mtype = MCrashVectorReq$
  $\wedge\ m \notin vServerProcessed[m.dest]$
  $\wedge\ vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] =$
       $vServerProcessed[m.dest] \cup \{m\}]$
  $\wedge\ HandleCrashVectorReq(m)$
  $\wedge\ \text{UNCHANGED }\ \langle coordStateVars, configManagerStateVars,$
       $vLog, vEarlyBuffer, vLateBuffer, vDeadlineQuorum,$

$$
\begin{aligned}
&\quad\quad vCrossShardConfirmQuorum,\ vServerStatus, \\
&\quad\quad vGView,\ vGVec,\ vLView,\ vServerClock,\ vLastNormView, \\
&\quad\quad vViewChange,\ vLSyncPoint,\ vLCommitPoint, \\
&\quad\quad vPeerCommitDeadline,\ vLSyncQuorum,\ vUUIDCounter, \\
&\quad\quad vCrashVector,\ vCrashVectorReps,\ vRecoveryReps\rangle \\
&\land ActionName' = \langle\,\text{``HandleCrashVectorReq"}\,\rangle
\end{aligned}
$$

$$
\begin{aligned}
\lor\ &\exists\, m \in messages : \\
&\land\ m.mtype = MCrashVectorRep \\
&\land\ m \notin vServerProcessed[m.dest] \\
&\land\ vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] = \\
&\quad\quad vServerProcessed[m.dest] \cup \{m\}] \\
&\land\ HandleCrashVectorRep(m) \\
&\land\ \text{UNCHANGED } \langle coordStateVars,\ configManagerStateVars, \\
&\quad\quad vLog,\ vEarlyBuffer,\ vLateBuffer, \\
&\quad\quad vDeadlineQuorum,\ vCrossShardConfirmQuorum,\ vServerStatus, \\
&\quad\quad vGView,\ vGVec,\ vLView,\ vServerClock,\ vLastNormView, \\
&\quad\quad vViewChange,\ vLSyncPoint,\ vLCommitPoint, \\
&\quad\quad vPeerCommitDeadline,\ vLSyncQuorum, \\
&\quad\quad vUUIDCounter,\ vCrashVectorReps,\ vRecoveryReps\rangle \\
&\land\ ActionName' = \langle\,\text{``HandleCrashVectorRep"}\,\rangle
\end{aligned}
$$

$$
\begin{aligned}
\lor\ &\exists\, m \in messages : \\
&\land\ m.mtype = MRecoveryReq \\
&\land\ m \notin vServerProcessed[m.dest] \\
&\land\ vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] = \\
&\quad\quad vServerProcessed[m.dest] \cup \{m\}] \\
&\land\ isCrashVectorValid(m) \\
&\land\ HandleRecoveryReq(m) \\
&\land\ \text{UNCHANGED } \langle coordStateVars,\ configManagerStateVars, \\
&\quad\quad vLog,\ vEarlyBuffer,\ vLateBuffer, \\
&\quad\quad vDeadlineQuorum,\ vCrossShardConfirmQuorum,\ vServerStatus, \\
&\quad\quad vGView,\ vGVec,\ vLView,\ vServerClock,\ vLastNormView, \\
&\quad\quad vViewChange,\ vLSyncPoint,\ vLCommitPoint, \\
&\quad\quad vPeerCommitDeadline,\ vLSyncQuorum, \\
&\quad\quad vUUIDCounter,\ vCrashVectorReps,\ vRecoveryReps\rangle \\
&\land\ ActionName' = \langle\,\text{``HandleRecoveryReq"}\,\rangle
\end{aligned}
$$

$$
\begin{aligned}
\lor\ &\exists\, m \in messages : \\
&\land\ m.mtype = MRecoveryRep \\
&\land\ m \notin vServerProcessed[m.dest] \\
&\land\ vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] = \\
&\quad\quad vServerProcessed[m.dest] \cup \{m\}] \\
&\land\ isCrashVectorValid(m) \\
&\land\ HandleRecoveryRep(m) \\
&\land\ \text{UNCHANGED } \langle coordStateVars,\ configManagerStateVars,
\end{aligned}
$$

$$
\begin{aligned}
&\quad\quad vLog,\ vEarlyBuffer,\ vLateBuffer, \\
&\quad\quad vDeadlineQuorum,\ vCrossShardConfirmQuorum,\ vServerStatus, \\
&\quad\quad vGVec,\ vServerClock,\ vLastNormView, \\
&\quad\quad vViewChange,\ vLSyncPoint,\ vLCommitPoint, \\
&\quad\quad vPeerCommitDeadline,\ vLSyncQuorum, \\
&\quad\quad vUUIDCounter,\ vCrashVectorReps,\ vRecoveryReps\rangle \\
&\quad \wedge\ ActionName' = \langle\text{``HandleRecoveryRep"}\rangle
\end{aligned}
$$

$$
\begin{aligned}
\vee\ &\exists\, m \in messages: \\
&\wedge\ m.mtype = MStartViewReq \\
&\wedge\ m \notin vServerProcessed[m.dest] \\
&\wedge\ vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] = \\
&\quad\quad vServerProcessed[m.dest] \cup \{m\}] \\
&\wedge\ isCrashVectorValid(m) \\
&\wedge\ HandleStartViewReq(m) \\
&\wedge\ \text{UNCHANGED }\ \langle coordStateVars,\ configManagerStateVars, \\
&\quad\quad vLog,\ vEarlyBuffer,\ vLateBuffer,\ vDeadlineQuorum, \\
&\quad\quad vCrossShardConfirmQuorum,\ vServerStatus, \\
&\quad\quad vGView,\ vGVec,\ vLView,\ vServerClock, \\
&\quad\quad vLastNormView,\ vViewChange,\ vLSyncPoint, \\
&\quad\quad vLCommitPoint,\ vPeerCommitDeadline,\ vLSyncQuorum, \\
&\quad\quad vUUIDCounter,\ vCrashVector, \\
&\quad\quad vCrashVectorReps,\ vRecoveryReps\rangle \\
&\wedge\ ActionName' = \langle\text{``HandleStartViewReq"}\rangle
\end{aligned}
$$

Periodic *Sync*
$$
\begin{aligned}
\vee\ &\exists\, serverId \in Servers: \\
&\wedge\ StartLocalSync(serverId) \\
&\wedge\ \text{UNCHANGED }\ \langle coordStateVars, \\
&\quad\quad serverStateVars,\ configManagerStateVars\rangle \\
&\wedge\ ActionName' = \langle\text{``StartLocalSync"}\rangle
\end{aligned}
$$

$$
\begin{aligned}
\vee\ &\exists\, m \in messages: \\
&\wedge\ m.mtype = MLocalSyncStatus \\
&\wedge\ m \notin vServerProcessed[m.dest] \\
&\wedge\ vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] = \\
&\quad\quad vServerProcessed[m.dest] \cup \{m\}] \\
&\wedge\ isCrashVectorValid(m) \\
&\wedge\ HandleLocalSyncStatus(m) \\
&\wedge\ \text{UNCHANGED }\ \langle coordStateVars,\ configManagerStateVars, \\
&\quad\quad vLog,\ vEarlyBuffer,\ vLateBuffer, \\
&\quad\quad vDeadlineQuorum,\ vCrossShardConfirmQuorum, \\
&\quad\quad vServerClock,\ vViewChange,\ vGVec,\ vGView, \\
&\quad\quad vLSyncPoint,\ vLView,\ vLastNormView, \\
&\quad\quad vServerStatus,\ vPeerCommitDeadline, \\
&\quad\quad vUUIDCounter,\ vCrashVectorReps,\ vRecoveryReps\rangle
\end{aligned}
$$

$\wedge\ ActionName' = \langle$ "HandleLocalSyncStatus" $\rangle$

$\vee\ \exists\, m \in messages :$
  $\wedge\ m.mtype = MLocalCommit$
  $\wedge\ m \notin vServerProcessed[m.dest]$
  $\wedge\ vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] =$
      $vServerProcessed[m.dest] \cup \{m\}]$
  $\wedge\ isCrashVectorValid(m)$
  $\wedge\ HandleLocalCommit(m)$
  $\wedge\ \text{UNCHANGED}\ \langle coordStateVars,\ configManagerStateVars,$
      $networkVars,\ vLog,\ vEarlyBuffer,\ vLateBuffer,$
      $vDeadlineQuorum,\ vCrossShardConfirmQuorum,$
      $vServerStatus,\ vServerClock,$
      $vGView,\ vGVec,\ vLView,\ vLastNormView,$
      $vViewChange,\ vLSyncPoint,\ vPeerCommitDeadline,$
      $vLSyncQuorum,\ vUUIDCounter,$
      $vCrashVectorReps,\ vRecoveryReps\rangle$
  $\wedge\ ActionName' = \langle$ "HandleLocalCommit" $\rangle$

$\vee\ \exists\, serverId \in Servers :$
  $\wedge\ BroadcastCommitStatusToPeers(serverId)$
  $\wedge\ \text{UNCHANGED}\ \langle coordStateVars,\ serverStateVars,$
      $configManagerStateVars\rangle$
  $\wedge\ ActionName' = \langle$ "BroadcastCommitStatusToPeers" $\rangle$

$\vee\ \exists\, m \in messages :$
  $\wedge\ m.mtype = MPeerShardCommitStatus$
  $\wedge\ m \notin vServerProcessed[m.dest]$
  $\wedge\ vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] =$
      $vServerProcessed[m.dest] \cup \{m\}]$
  $\wedge\ HandlePeerShardCommitStatus(m)$
  $\wedge\ \text{UNCHANGED}\ \langle networkVars,\ coordStateVars,\ configManagerStateVars,$
      $vLog,\ vEarlyBuffer,\ vLateBuffer,\ vServerStatus,$
      $vDeadlineQuorum,\ vCrossShardConfirmQuorum,$
      $vGView,\ vGVec,\ vLView,\ vServerClock,\ vLastNormView,$
      $vViewChange,\ vLSyncPoint,\ vLCommitPoint,$
      $vPeerCommitDeadline,\ vLSyncQuorum,\ vUUIDCounter,$
      $vCrashVector,\ vCrashVectorReps,\ vRecoveryReps\rangle$
  $\wedge\ ActionName' = \langle$ "HandlePeerShardCommitStatus" $\rangle$

Clock Move
$\vee\ \exists\, serverId \in Servers :$
  $\wedge\ ServerClockMove(serverId)$
  $\wedge\ \text{UNCHANGED}\ \langle coordStateVars,\ configManagerStateVars\rangle$
  $\wedge\ ActionName' = \langle$ "ServerClockMove" $\rangle$

$\lor \exists\, coordId \in Coords :$
$\quad \land\, CoordClockMove(coordId)$
$\quad \land\, \textsc{unchanged}\ \langle networkVars,\ serverStateVars,\ configManagerStateVars,$
$\qquad\quad vCoordTxns,\ vCoordProcessed\rangle$
$\quad \land\, ActionName' = \langle\, \text{``CoordClockMove''}\,\rangle$

$Spec\ \triangleq\ Init \land \square[Next]\langle networkVars,$
$\qquad\qquad\qquad\qquad\quad serverStateVars,\ coordStateVars,\ configManagerStateVars,$
$\qquad\qquad\qquad\qquad\quad ActionName\rangle$

$ShardRecovered(shardId,\ lViewID)\ \triangleq$
$\quad \textsc{let}$
$\qquad serversInOneShard\ \triangleq\ \{s \in Servers : s.shardId = shardId\}$
$\qquad leaderServer\ \triangleq\ [$
$\qquad\quad replicaId \mapsto LeaderID(lViewID),$
$\qquad\quad shardId \mapsto shardId$
$\qquad ]$
$\quad \textsc{in}$
$\quad \land\, \exists\, RM \in \textsc{subset}\ (serversInOneShard) :$
$\qquad \land\ Cardinality(RM) \geq QuorumSize$
$\qquad \land\ leaderServer \in RM$
$\qquad \land\ \forall\ r \in RM : vServerStatus[r]\quad = StNormal$
$\qquad \land\ \forall\ r \in RM : vLastNormView[r] \geq lViewID$

$CommittedInView(v,\ shardId,\ txnId)\ \triangleq$
$\quad \textsc{let}$
$\qquad serversInOneShard\ \triangleq\ \{s \in Servers : s.shardId = shardId\}$
$\qquad leaderServer\ \triangleq\ [$
$\qquad\quad replicaId \mapsto LeaderID(v),$
$\qquad\quad shardId \mapsto shardId$
$\qquad ]$
$\qquad replySet\ \triangleq\ \{$
$\qquad\quad m \in messages : \land\ \ \lor\ m.mtype = MFastReply$
$\qquad\qquad\qquad\qquad\qquad\qquad\ \lor\ m.mtype = MSlowReply$
$\qquad\qquad\qquad\qquad\quad \land\ m.txnId = txnId$
$\qquad\qquad\qquad\qquad\quad \land\ m.sender \in serversInOneShard$
$\qquad\qquad\qquad\qquad\quad \land\ m.lView = v$
$\qquad \}$
$\quad \textsc{in}$
$\quad \textsc{if}\ \forall\, reply \in replySet :$
$\qquad \lor\ reply.mtype \neq MFastReply$
$\qquad \lor\ reply.sender \neq leaderServer$
$\quad\ \textsc{then}\quad$ No leader's fast $reply \rightarrow$ This $txn$ is not committed
$\qquad\quad \textsc{false}$

39

```
ELSE
    LET
        leaderReply  ≜  CHOOSE reply ∈ replySet :
                                ∧  reply.mtype = MFastReply
                                ∧  reply.sender = leaderServer
    IN
        Committed in Fast Path
    ∨  ∃ fastQuorum ∈ SUBSET replySet :
            ∧  leaderReply ∈ fastQuorum
            ∧  Cardinality(fastQuorum) = FastQuorumSize
            All replies have the same hash (or it is a slow reply)
            ∧  ∀ reply ∈ fastQuorum :
                    ∨  ∧  reply.mtype = MFastReply
                       ∧  reply.hash = leaderReply.hash
                    Slow Reply can be used as fast reply
                    ∨  reply.mtype = MSlowReply
        Committed in Slow Path
    ∨  ∃ slowQuorum ∈ SUBSET  replySet :
            ∧  leaderReply ∈ slowQuorum
            ∧  Cardinality(slowQuorum) = QuorumSize
            ∧  ∀ reply ∈ slowQuorum \ {leaderReply} :
                    reply.mtype = MSlowReply
```

## Invariants

Durability [In-Shard-Property]: Committed txns always survive failure *i.e.* If a *txn* is committed (to be more precise, locally 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  ≜
    ∀ shardId ∈ Shards :
        ∀ v1, v2 ∈ 0 .. MaxViews :
            If a txn is committed in lower view (v1,),
            it is impossible to make this request uncommited in higher vie
            ¬(  ∧  v1 < v2
                ∧  ShardRecovered(shardId, v2)
                ∧  ∃ c ∈ Coords :
                        ∃ txnId ∈ vCoordTxns[c] :
                            ∧  CommittedInView(v1, shardId, txnId)
                            ∧  ¬CommittedInView(v2, shardId, txnId)
            )
```

Consistency [In-Shard-Property]: Committed txns 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), then it remains committed in higher view ($v2$)

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

$Consistency \triangleq$
    $\forall\, shardId \in Shards :$
       $\forall\, v1,\, v2 \in 1\, .. \, MaxViews :$
          $\neg(\ \ \wedge\, v1 < v2$
               To check *Consistency* of txns in higher views,
               the shard should have entered the higher views
             $\wedge\, ShardRecovered(shardId,\, v2)$
             $\wedge\, \exists\, c \in Coords :$
               $\exists\, txnId \in vCoordTxns[c] :$
                   Durability has been checked in another invariant
                 IF   $\wedge\, CommittedInView(v1,\, shardId,\, txnId)$
                     $\wedge\, CommittedInView(v2,\, shardId,\, txnId)$
                THEN
                  LET
                    $v1LeaderReply \triangleq$ CHOOSE $m \in messages :$
                            $\wedge$     $m.mtype = MFastReply$
                            $\wedge$     $m.txnId\ = txnId$
                            $\wedge$     $m.lView = v1$
                            $\wedge$     $m.sender.shardId = shardId$
                            $\wedge$     $m.sender.replicaId = LeaderID(v1)$
                    $v2LeaderReply \triangleq$ CHOOSE $m \in messages :$
                            $\wedge$     $m.mtype = MFastReply$
                            $\wedge$     $m.txnId\ = txnId$
                            $\wedge$     $m.lView = v2$
                            $\wedge$     $m.sender.shardId = shardId$
                            $\wedge$     $m.sender.replicaId = LeaderID(v2)$
                 IN
                    $v1LeaderReply.hash \neq v2LeaderReply.hash$
                ELSE   FALSE
       $)$

*Linearizability* [In-Shard-Property]: Only one *txn* can be committed for a given position, *i.e.* If one *txn* has committed at position $i$, then no contrary observation can be made

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

$Linearizability \triangleq$
    LET
       $allTxns \triangleq$ UNION $\{vCoordTxns[c] : c \in Coords\}$
    IN
    $\forall\, shardId \in Shards :$
       $\forall\, txnId1,\, txnId2 \in allTxns :$

IF $txnId1 = txnId2$ THEN TRUE
ELSE
 $\forall\, v1,\, v2 \in 1\,..\, MaxViews :$
  IF $\land\;\; CommittedInView(v1, shardId, txnId1)$
    $\land\;\; CommittedInView(v1, shardId, txnId2)$
  THEN
   LET
    $v1LeaderReply \;\triangleq\;$ CHOOSE $m \in messages :$
          $\land\;\;\;\; m.mtype = MFastReply$
          $\land\;\;\;\; m.txnId\; = txnId1$
          $\land\;\;\;\; m.lView = v1$
          $\land\;\;\;\; m.sender.shardId = shardId$
          $\land\;\;\;\; m.sender.replicaId = LeaderID(v1)$
    $v2LeaderReply \;\triangleq\;$ CHOOSE $m \in messages :$
          $\land\;\;\;\; m.mtype = MFastReply$
          $\land\;\;\;\; m.txnId\; = txnId2$
          $\land\;\;\;\; m.lView = v2$
          $\land\;\;\;\; m.sender.shardId = shardId$
          $\land\;\;\;\; m.sender.replicaId = LeaderID(v2)$
   IN
    They cannot be committed in the same *log* position, regardless of the view
    $v1LeaderReply.logId \neq v2LeaderReply.logId$
  ELSE Not both are committed, so no need to check
   TRUE

Serializability [Cross-Shard-Property]: Given two txns and two *shards*: If they are both committed in both *shards*, then they should be committed in the same order, *i.e.*, if $txn-1$ committed before $txn-2$ on Shard$-1$, then $txn-1$ is also committed before $txn-2$ on Shard$-2$

$Serializability \;\triangleq\;$
 LET
  $allTxns \;\triangleq\;$ UNION $\{vCoordTxns[c] : c \in Coords\}$
 IN
 $\forall\, txnId1,\, txnId2 \in allTxns :$
  IF $txnId1 = txnId2$ THEN TRUE
  ELSE
   $\forall\, v \in 1\,..\, MaxViews :$
    $\forall\, shardId1,\, shardId2 \in Shards :$
     IF $shardId1 = shardId2$ THEN TRUE
     ELSE
      IF $\land\;\; CommittedInView(v, shardId1, txnId1)$
        $\land\;\; CommittedInView(v, shardId1, txnId2)$
        $\land\;\; CommittedInView(v, shardId2, txnId1)$
        $\land\;\; CommittedInView(v, shardId2, txnId2)$
      THEN
       LET

$txn1\_LeaderReplyOnShard1 \triangleq$ CHOOSE $m \in messages :$
$\qquad \wedge\ m.mtype = MFastReply$
$\qquad \wedge\ m.txnId\ = txnId1$
$\qquad \wedge\ m.lView = v$
$\qquad \wedge\ m.sender.shardId = shardId1$
$\qquad \wedge\ m.sender.replicaId = LeaderID(v)$
$txn2\_LeaderReplyOnShard1 \triangleq$ CHOOSE $m \in messages :$
$\qquad \wedge\ m.mtype = MFastReply$
$\qquad \wedge\ m.txnId\ = txnId2$
$\qquad \wedge\ m.lView = v$
$\qquad \wedge\ m.sender.shardId = shardId1$
$\qquad \wedge\ m.sender.replicaId = LeaderID(v)$
$txn1\_LeaderReplyOnShard2 \triangleq$ CHOOSE $m \in messages :$
$\qquad \wedge\ m.mtype = MFastReply$
$\qquad \wedge\ m.txnId\ = txnId1$
$\qquad \wedge\ m.lView = v$
$\qquad \wedge\ m.sender.shardId = shardId2$
$\qquad \wedge\ m.sender.replicaId = LeaderID(v)$
$txn2\_LeaderReplyOnShard2 \triangleq$ CHOOSE $m \in messages :$
$\qquad \wedge\ m.mtype = MFastReply$
$\qquad \wedge\ m.txnId\ = txnId2$
$\qquad \wedge\ m.lView = v$
$\qquad \wedge\ m.sender.shardId = shardId2$
$\qquad \wedge\ m.sender.replicaId = LeaderID(v)$
IN
$\vee\ \wedge\ txn1\_LeaderReplyOnShard1.logId > txn2\_LeaderReplyOnShard1.logId$
$\quad\ \wedge\ txn1\_LeaderReplyOnShard2.logId > txn2\_LeaderReplyOnShard2.logId$
$\vee\ \wedge\ txn1\_LeaderReplyOnShard1.logId < txn2\_LeaderReplyOnShard1.logId$
$\quad\ \wedge\ txn1\_LeaderReplyOnShard2.logId < txn2\_LeaderReplyOnShard2.logId$

ELSE  TRUE