

## 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*]

$Replicas \triangleq 0 \dots 2$

$ReplicaOrder \triangleq \langle 0, 1, 2 \rangle$

$Shards \triangleq 0 \dots 2$

$Coords \triangleq 0 \dots 1$

$LatencyBounds \triangleq [c \in Coords \mapsto 1]$

ASSUME  $IsFiniteSet(Replicas)$

ASSUME  $IsFiniteSet(Shards)$

ASSUME  $ReplicaOrder \in Seq(Replicas)$

$Servers \triangleq \{$

$[$

$replicaId \mapsto e[1],$

$shardId \mapsto e[2]$

$] : e \in Replicas \times Shards$

$\}$

---

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

### Constants

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

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

Server Statuses (Within Shard)

$StNormal \triangleq 1$   
 $StViewChange \triangleq 2$   
 $StCrossShardSyncing \triangleq 3$   
 $StRecovering \triangleq 4$   
 $StFailing \triangleq 5$

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  $\langle replicaId, shardId \rangle$ .  $TxnID$  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  $\langle deadline, txnId \rangle$  uniquely identifies one request across replicas

$TxnID = [$   
 $\quad coordId \mapsto i \text{ in } (1 \dots),$   
 $\quad rId \mapsto i \text{ in } (1 \dots)$   
 $]$

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

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

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

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

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

```

InterReplicaSync = [
  mtype      ↦ MInterReplicaSync,
  lView      ↦ 0...y
  sender     ↦ src ∈ Servers,
  dest       ↦ dst ∈ Servers,
  entries    ↦ [LogEntry...]
]

```

*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      ↦ MFastReply,
  sender     ↦ src ∈ Servers,
  dest       ↦ dst ∈ Coords,
  gView      ↦ 0...x
  lView      ↦ 0...x
  txnId     ↦ txnId

```

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

```

hash        ↦ [ entries ↦ log entries so far cv ↦ crashVector ]
deadline    ↦ i ∈ (1 .. MaxTime + MaxBound),
logId       ↦ n ∈ (1 .. )
]

```

```

SlowReply = [
  mtype      ↦ MSlowReply,
  sender     ↦ src ∈ Servers,
  dest       ↦ c ∈ Coords,
  gView      ↦ 0...x
  lView      ↦ 0...x
  txnId     ↦ txnId
  logId      ↦ n ∈ (1 .. )
]

```

```

ViewChangeReq = [
  mtype ↦ MViewChangeReq,
  sender ↦ src ∈ Replicas, (by configManager)
  dest ↦ dst ∈ Servers,
  gView ↦ 0..x
  gVec ↦ the lViews for each shard
]

```

```

ViewChange = [
  mtype      ↦ MViewChange,
  sender     ↦ src ∈ Servers,
  dest       ↦ dst ∈ Servers,
  gView      ↦ 0..x
  gVec       ↦ the lViews for each shard
  lView      ↦ 0...x
  lastNormal ↦ v ∈ ViewIDs,

```

```

lSyncPoint  $\mapsto 0 \dots$ 
entries  $\mapsto l \in vLogs[1 \dots n]$ ,
cv  $\mapsto$  crash vector
]

CrossShardConfirm = [
  mtype  $\mapsto MCrossShardConfirm$ ,
  sender  $\mapsto src \in Servers$ ,
  dest  $\mapsto dst \in Servers$ ,
  lView  $\mapsto 0 \dots x$ 
  gView  $\mapsto 0 \dots$ 
  entries  $\mapsto l \in vLogs[1 \dots n]$ 
]

StartView = [
  mtype  $\mapsto MStartView$ ,
  sender  $\mapsto src \in Servers$ ,
  dest  $\mapsto dst \in Servers$ ,
  lView  $\mapsto 0 \dots x$ 
  gView  $\mapsto 0 \dots x$ 
  gVec  $\mapsto$  the lViews for each shard
  entries  $\mapsto l \in vLogs[1 \dots n]$ ,
  cv  $\mapsto$  crash vector
]

CrashVectorReq = [
  mtype  $\mapsto MCrashVectorReq$ ,
  sender  $\mapsto src \in Servers$ ,
  dest  $\mapsto dst \in Servers$ ,
  nonce  $\mapsto nonce$ 
]

CrashVectorRep = [
  mtype  $\mapsto MCrashVectorRep$ ,
  sender  $\mapsto src \in Servers$ ,
  dest  $\mapsto dst \in Servers$ ,
  nonce  $\mapsto nonce$ ,
  cv  $\mapsto$  vector of counters
]

RecoveryReq = [
  mtype  $\mapsto MRecoveryReq$ ,
  sender  $\mapsto src \in Servers$ ,
  dest  $\mapsto dst \in Servers$ ,
  cv  $\mapsto$  vector of counters
]

RecoveryRep = [
  mtype  $\mapsto MRecoveryRep$ ,
  sender  $\mapsto src \in Servers$ ,
  dest  $\mapsto dst \in Servers$ ,
  gView  $\mapsto 0 \dots x$ 

```

$lView \mapsto 0 \dots x$   
 $cv \mapsto \text{vector of counters}$   
 ]  
 $StartViewReq = [$   
 $mtype \mapsto MStartViewReq,$   
 $sender \mapsto src \in Servers,$   
 $dest \mapsto dst \in Servers,$   
 $lView \mapsto 0 \dots x$   
 $cv \mapsto \text{vector of counters}$   
 $]$

Follower reports to its leader

$LocalSyncStatus = [$   
 $mtype \mapsto MLocalSyncStatus,$   
 $sender \mapsto src \in Servers,$   
 $dest \mapsto dst \in Servers,$   
 $lView \mapsto 0 \dots x$   
 $lSyncPoint \mapsto n \in (1 \dots)$   
 $cv \mapsto \text{vector of counters}$   
 $]$

Leader notifies its followers

$LocalCommit = [$   
 $mtype \mapsto MLocalCommit,$   
 $sender \mapsto src \in Servers,$   
 $dest \mapsto dst \in Servers,$   
 $lView \mapsto 0 \dots x$   
 $entries \mapsto \text{log entries}$   
 $lCommitPoint \mapsto n \in (1 \dots)$   
 $]$

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)

$PeerShardCommitStatus = [$   
 $mtype \mapsto MPeerShardCommitStatus,$   
 $sender \mapsto src \in Servers,$   
 $dest \mapsto dst \in Servers,$   
 $gView \mapsto 0 \dots x$   
 $deadline \mapsto \text{the largest committed deadline}$   
 $]$

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

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

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

```

---

## Network State

The variables record the messages processed by *Replicas*/Clients, so that the *Replicas*/Clients will not process twice

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 *DeadlineQuorum* to collect the deadlines from other servers for agreement

*vDeadlineQuorum*,

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

*vCrossShardConfirmQuorum*,

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

*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

*vCMStatus*,

*vCMView*,

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

*vCMGInfo*,

*vCMPPrepareGInfo*,

Config Manager: quorum of *CMPrepareReplies*

*vCMPPrepareReps*,

*vCMPProcessed*

VARIABLES *ActionName*

*networkVars*  $\triangleq$   $\langle \text{messages} \rangle$

*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*  $\triangleq$   $\langle vCMStatus, vCMView, vCMGInfo,$   
 $vCMPPrepareGInfo, vCMPPrepareReps,$   
 $vCMPProcessed \rangle$

*InitNetworkState*  $\triangleq$   $messages = \{\}$

*InitServerState*  $\triangleq$

$\wedge vServerProcessed = [serverId \in Servers \mapsto \{\}]$   
 $\wedge vLog = [serverId \in Servers \mapsto \langle \rangle]$

$$\begin{aligned}
& \wedge vEarlyBuffer = [serverId \in Servers \mapsto \{\}] \\
& \wedge vLateBuffer = [serverId \in Servers \mapsto \{\}] \\
& \wedge vDeadlineQuorum = [serverId \in Servers \mapsto \{\}] \\
& \wedge vCrossShardConfirmQuorum = [serverId \in Servers \mapsto \{\}] \\
& \wedge vServerStatus = [serverId \in Servers \mapsto StNormal] \\
& \wedge vGView = [serverId \in Servers \mapsto 0] \\
& \wedge vGVec = [ \\
& \quad serverId \in Servers \mapsto [ \\
& \quad \quad shardId \in Shards \mapsto 0 \\
& \quad ] \\
& ] \\
& \wedge vLView = [serverId \in Servers \mapsto 0] \\
& \wedge vServerClock = [serverId \in Servers \mapsto 1] \\
& \wedge vLastNormView = [serverId \in Servers \mapsto 0] \\
& \wedge vViewChange = [serverId \in Servers \mapsto \{\}] \\
& \wedge vLSyncPoint = [serverId \in Servers \mapsto 0] \\
& \wedge vLCommitPoint = [serverId \in Servers \mapsto 0] \\
& \wedge vPeerCommitDeadline = [serverId \in Servers \mapsto \\
& \quad [shardId \in Shards \mapsto 0] \\
& ] \\
& \wedge vLSyncQuorum = [serverId \in Servers \mapsto \{\}] \\
& \wedge vUUIDCounter = [serverId \in Servers \mapsto 0] \\
& \wedge vCrashVector = [ \\
& \quad serverId \in Servers \mapsto [ \\
& \quad \quad rr \in Replicas \mapsto 0 \\
& \quad ] \\
& ] \\
& \wedge vCrashVectorReps = [serverId \in Servers \mapsto \{\}] \\
& \wedge vRecoveryReps = [serverId \in Servers \mapsto \{\}]
\end{aligned}$$

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

$$\begin{aligned}
InitConfigManagerState & \triangleq \\
& \wedge vCMStatus = [ \\
& \quad replicaId \in Replicas \mapsto StNormal \\
& ] \\
& \wedge vCMView = [ \\
& \quad replicaId \in Replicas \mapsto 0 \\
& ] \\
& \wedge vCMGInfo = [ \\
& \quad replicaId \in Replicas \mapsto [
\end{aligned}$$

$$\begin{aligned}
& \quad gView \mapsto 0, \\
& \quad gVec \mapsto [shardId \in Shards \mapsto 0] \\
& \quad ] \\
& \quad ] \\
\wedge \quad vCMPPrepareGInfo = [ \\
& \quad replicaId \in Replicas \mapsto [ \\
& \quad \quad gView \mapsto 0, \\
& \quad \quad gVec \mapsto [shardId \in Shards \mapsto 0] \\
& \quad ] \\
& \quad ] \\
\wedge \quad vCMPPrepareReps = [ \\
& \quad replicaId \in Replicas \mapsto \{\} \\
& \quad ] \\
\wedge \quad vCMPProcessed = [ \\
& \quad replicaId \in Replicas \mapsto \{\} \\
& \quad ]
\end{aligned}$$

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

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

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

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

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

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

$$IsInjective(s) \triangleq$$

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.

$$\text{CHOOSE } f \in [1 \dots Cardinality(S) \rightarrow S] : IsInjective(f)$$

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

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

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

$\text{SetToSortSeq}(S, \text{op}(-, -)) \triangleq$

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  $\text{SetToSeq}$  but with an additional conjunct, because this variant works efficiently without a dedicated *TLC* override.

$\text{SortSeq}(\text{SetToSeq}(S), \text{op})$

### View ID Helpers

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

$\text{isLeader}(\text{replicaId}, \text{viewId}) \triangleq (\text{replicaId} = \text{LeaderID}(\text{viewId}))$

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

$\text{ViewGreater}(\text{gv1}, \text{lv1}, \text{gv2}, \text{lv2}) \triangleq$

IF  $\text{gv1} > \text{gv2}$  THEN TRUE

ELSE

IF  $\wedge \text{gv1} = \text{gv2}$

$\wedge \text{lv1} > \text{lv2}$

THEN TRUE

ELSE FALSE

### Client action

Coordinator  $c$  submits a  $\text{txn}$ . We assume Coordinator can only send one  $\text{txn}$  in one tick of time.

If time has reached the bound, this client cannot send request any more

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

$\text{CoordSubmitTxn}(c) \triangleq$

$\wedge \text{vCoordClock}[c] < \text{MaxTime}$

$\wedge \text{Cardinality}(\text{vCoordTxns}[c]) < \text{MaxReqNum}$

$\wedge \text{LET}$

$\text{txnId} \triangleq [$   
 $\text{coordId} \mapsto c,$   
 $\text{rId} \mapsto \text{Cardinality}(\text{vCoordTxns}[c]) + 1$   
 $]$

IN

$\wedge \text{Send}(\{[\text{mtype} \mapsto \text{MTxn},$

$\text{txnId} \mapsto \text{txnId},$

$\text{command} \mapsto \text{""},$

Here we assume involves all *shards*

$$\begin{aligned}
& \text{shards} \mapsto \text{Shards}, \\
& \text{st} \mapsto v\text{CoordClock}[c], \\
& \text{bound} \mapsto \text{LatencyBounds}[c], \\
& \text{sender} \mapsto c, \\
& \text{dest} \mapsto \text{serverId} \\
& ] : \text{serverId} \in \text{Servers}) \\
& \wedge v\text{CoordClock}' = [v\text{CoordClock} \text{ EXCEPT } ![c] = v\text{CoordClock}[c] + 1] \\
& \wedge v\text{CoordTrns}' = [v\text{CoordTrns} \text{ EXCEPT } ![c] = v\text{CoordTrns}[c] \cup \{\text{txnId}\}]
\end{aligned}$$

$$\text{HandleTxn}(m) \triangleq$$

$$\text{LET}$$

$$\text{myServerId} \triangleq m.\text{dest}$$

$$\text{newLog} \triangleq [$$

$$\begin{aligned}
& \text{mtype} \mapsto \text{MLogEntry}, \\
& \text{txnId} \mapsto m.\text{txnId}, \\
& \text{command} \mapsto m.\text{command}, \\
& \text{shards} \mapsto m.\text{shards}, \\
& \text{deadline} \mapsto \text{Max}(\text{LastAppendedDeadline}(v\text{Log}[\text{myServerId}]), m.\text{st} + m.\text{bound})
\end{aligned}$$

$$]$$

$$\text{serversInOneReplica} \triangleq \{s \in \text{Servers} : s.\text{replicaId} = \text{myServerId}.\text{replicaId}\}$$

$$\text{IN}$$

$$\vee \wedge \text{isLeader}(\text{myServerId}.\text{replicaId}, v\text{LView}[\text{myServerId}])$$

$$\begin{aligned}
& \wedge v\text{EarlyBuffer}' = [ \\
& \quad v\text{EarlyBuffer} \text{ EXCEPT } ![myServerId] \\
& \quad = v\text{EarlyBuffer}[\text{myServerId}] \cup \{\text{newLog}\}
\end{aligned}$$

$$\text{Broadcast deadline notifications to other shards}$$

$$\begin{aligned}
& \wedge \text{Send}(\{[ \\
& \quad \text{mtype} \mapsto \text{MDeadlineNotification}, \\
& \quad \text{gView} \mapsto v\text{GView}[\text{myServerId}], \\
& \quad \text{lView} \mapsto v\text{LView}[\text{myServerId}], \\
& \quad \text{sender} \mapsto \text{myServerId}, \\
& \quad \text{dest} \mapsto \text{dstServerId}, \\
& \quad \text{entry} \mapsto \text{newLog} \\
& \quad ] : \text{dstServerId} \in \text{serversInOneReplica}\})
\end{aligned}$$

$$\wedge \text{UNCHANGED } \langle v\text{LateBuffer} \rangle$$

$$\vee \wedge \neg \text{isLeader}(\text{myServerId}.\text{replicaId}, v\text{LView}[\text{myServerId}])$$

$$\begin{aligned}
& \wedge \vee \wedge \text{newLog}.\text{deadline} = (m.\text{st} + m.\text{bound}) \\
& \quad \wedge v\text{EarlyBuffer}' = [ \\
& \quad \quad v\text{EarlyBuffer} \text{ EXCEPT } ![myServerId] \\
& \quad \quad = v\text{EarlyBuffer}[\text{myServerId}] \cup \{\text{newLog}\} \\
& \quad ]
\end{aligned}$$

$$\wedge \text{UNCHANGED } \langle v\text{LateBuffer} \rangle$$

$$\vee \wedge \neg (\text{newLog}.\text{deadline} = (m.\text{st} + m.\text{bound}))$$

$$\wedge v\text{LateBuffer}' = [$$

$$\begin{aligned}
& vLateBuffer \text{ EXCEPT } ![myServerId] \\
& \quad = vLateBuffer[myServerId] \cup \{newLog\} \\
& ] \\
& \wedge \text{ UNCHANGED } \langle vEarlyBuffer \rangle \\
& \wedge \text{ UNCHANGED } \langle networkVars \rangle
\end{aligned}$$

$HandleDeadlineNotification(m) \triangleq$

LET

$myServerId \triangleq m.dest$

$quorum \triangleq \{$

$msg \in vDeadlineQuorum[myServerId]$

$: \wedge msg.entry.txnId = m.entry.txnId$

$\wedge msg.gView = m.gView$

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

$\} \cup \{m\}$

IN

Only leader does deadline agreement

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

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

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

$\wedge vDeadlineQuorum' = [$

$vDeadlineQuorum \text{ EXCEPT } ![myServerId]$

$= vDeadlineQuorum[myServerId] \cup \{m\}$

$]$

$\wedge \text{ IF } Cardinality(quorum) = Cardinality(m.entry.shards)$

THEN

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

LET

$maxDeadlineTxn \triangleq$

CHOOSE  $x \in quorum :$

$\forall y \in quorum :$

$y.entry.deadline \leq x.entry.deadline$

$sequencingTxn \triangleq$

CHOOSE  $x \in vEarlyBuffer[myServerId] :$

$x.txnId = m.entry.txnId$

IN

$\text{ IF } maxDeadlineTxn.entry.deadline > sequencingTxn.deadline$

THEN

$vEarlyBuffer' = [vEarlyBuffer \text{ EXCEPT } ![myServerId]$

$= (vEarlyBuffer[myServerId] \setminus \{sequencingTxn\}) \cup \{maxDeadlineTxn.entry\}]$

ELSE UNCHANGED  $\langle vEarlyBuffer \rangle$

ELSE

Deadline quorum not sufficient so far: do not take further actions

UNCHANGED  $\langle vEarlyBuffer \rangle$

$HandleInterReplicaSync(m) \triangleq$

$\wedge m.lView = vLView[m.dest]$

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

$\wedge m.crashVector[m.sender] = vCrashVector[m.sender]$

$\wedge \neg isLeader(m.dest.replicaId, vLView[m.dest])$

$\wedge$  LET

$myServerId \triangleq m.dest$

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

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

IN

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

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

Kick synced entries out of *earlyBuffer*

$\wedge vEarlyBuffer' = [$   
 $\quad vEarlyBuffer \text{ EXCEPT } ![myServerId]$   
 $\quad = \{msg \in vEarlyBuffer[myServerId] :$   
 $\quad \quad msg.txnId \notin syncedTxnIds\}$   
 $\quad ]$

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 unnecessary transmission in practice.

$\wedge vLateBuffer' = [$   
 $\quad vLateBuffer \text{ EXCEPT } ![myServerId]$   
 $\quad = \{msg \in vLateBuffer[myServerId] :$   
 $\quad \quad msg.txnId \notin syncedTxnIds\}$   
 $\quad ]$

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

$\wedge vDeadlineQuorum' = [$   
 $\quad vDeadlineQuorum \text{ EXCEPT } ![myServerId]$   
 $\quad = \{msg \in vDeadlineQuorum[myServerId] :$   
 $\quad \quad msg.txnId \notin syncedTxnIds\}$   
 $\quad ]$

$\wedge vLSyncPoint' = [$   
 $\quad vLSyncPoint \text{ EXCEPT } ![myServerId] = Len(m.entries)]$

Send slow-replies to coordinators

$\wedge Send(\{[$   
 $\quad mtype \mapsto MSlowReply,$   
 $\quad sender \mapsto myServerId,$   
 $\quad ]$

$$\begin{aligned}
& dest \mapsto m.entries[i].txnId.coordId, \\
& gView \mapsto vGView[myServerId], \\
& lView \mapsto vLView[myServerId], \\
& txnId \mapsto m.entries[i].txnId, \\
& logId \mapsto i \\
& ] : i \in (currentSyncPoint + 1) \dots Len(m.entries)\} \\
\vee & \wedge currentSyncPoint \geq Len(m.entries) \\
& \text{Noting new to sync} \\
& \wedge \text{UNCHANGED } \langle networkVars, vLog, vEarlyBuffer, \\
& \quad vLateBuffer, vDeadlineQuorum, vLSyncPoint \rangle
\end{aligned}$$

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

$$\begin{aligned}
& DetectLeaderFail(cmReplicaId) \triangleq \\
& \quad \exists shardId \in Shards : \\
& \quad \text{LET} \\
& \quad \quad lView \triangleq vCMGInfo[cmReplicaId].gVec[shardId] \\
& \quad \quad leaderId \triangleq LeaderID(lView) \\
& \quad \quad serverId \triangleq [ \\
& \quad \quad \quad replicaId \mapsto leaderId, \\
& \quad \quad \quad shardId \mapsto shardId \\
& \quad \quad ] \\
& \quad \text{IN} \\
& \quad vServerStatus[serverId] = StFailing
\end{aligned}$$

$$\begin{aligned}
& SelectProperLView(currentView, shardId) \triangleq \\
& \quad \text{LET} \\
& \quad \quad aliveReplicaId \triangleq \text{CHOOSE } replicaId \in Replicas :
\end{aligned}$$



$$\begin{aligned}
& vServerStatus[shardId][replicaId] = StNormal \\
& \text{IN} \\
& \quad \text{Ensure 1 the new view is larger than } currentView \\
& \quad * (2) \text{ its corresponding leader happens to be the selected } aliveReplicaId \\
& \quad (currentView \div Cardinality(Replicas) + 1) * Cardinality(Replicas) + aliveReplicaId \\
& PrepareViewChange(cmReplicaId) \triangleq \\
& \quad \text{LET} \\
& \quad \quad newGVec \triangleq [ \\
& \quad \quad \quad shardId \in Shards \mapsto \\
& \quad \quad \quad \quad SelectProperLView(vCMGInfo[cmReplicaId].gVec[shardId], shardId) \\
& \quad \quad ] \\
& \quad \text{IN} \\
& \quad \wedge vCMPPrepareGInfo' = [vCMPPrepareGInfo \text{ EXCEPT } ![cmReplicaId] = \\
& \quad \quad [ \\
& \quad \quad \quad gView \mapsto vCMGInfo[cmReplicaId].gView + 1, \\
& \quad \quad \quad gVec \mapsto newGVec \\
& \quad \quad ] \\
& \quad ] \\
& \quad \wedge Send(\{[ \\
& \quad \quad \quad mtype \mapsto MCMPrepare, \\
& \quad \quad \quad sender \mapsto cmReplicaId, \\
& \quad \quad \quad dest \mapsto dstRid, \\
& \quad \quad \quad cView \mapsto vCMView[cmReplicaId], \\
& \quad \quad \quad gView \mapsto vCMPPrepareGInfo'[cmReplicaId].gView, \\
& \quad \quad \quad gVec \mapsto newGVec \\
& \quad \quad ] : dstRid \in Replicas\}) \\
& LaunchViewChange(cmReplicaId) \triangleq \\
& \quad \text{IF } \wedge isLeader(cmReplicaId, vCMView[cmReplicaId]) \\
& \quad \quad \wedge DetectLeaderFail(cmReplicaId) \\
& \quad \text{THEN} \\
& \quad \quad PrepareViewChange(cmReplicaId) \\
& \quad \text{ELSE} \\
& \quad \quad \text{UNCHANGED } \langle networkVars \rangle \\
& HandleCMPPrepare(m) \triangleq \\
& \quad \wedge m.cView = vCMView[m.dest] \\
& \quad \wedge m.gView > vCMGInfo[m.dest].gView \\
& \quad \wedge vCMPPrepareGInfo' = [vCMPPrepareGInfo \text{ EXCEPT } ![m.dest] = \\
& \quad \quad [ \\
& \quad \quad \quad gView \mapsto m.gView, \\
& \quad \quad \quad gVec \mapsto m.gVec \\
& \quad \quad ] \\
\end{aligned}$$

```

    ]
  ]
  ∧ Send([
    mtype ↦ MCMPrepareReply,
    sender ↦ m.dest,
    dest ↦ m.src,
    cView ↦ m.cView,
    gView ↦ m.gView
  ])

```

```

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

```

```

HandleCMCommit(m)  $\triangleq$ 

```

$$\begin{aligned}
& \wedge m.cView = vCMView[m.dest] \\
& \wedge \neg isLeader(m.dest, vCMView[m.dest]) \\
& \wedge m.gView = vCMPPrepareGInfo[m.dest].gView \\
& \wedge vCMGInfo' = [vCMGInfo \text{ EXCEPT } ![m.dest] = \\
& \quad \quad \quad vCMPPrepareGInfo[m.dest] \\
& ]
\end{aligned}$$

*HandleViewChangeReq*(*m*)  $\triangleq$

LET

$$\begin{aligned}
myServerId & \triangleq m.dest \\
myLeader & \triangleq \text{CHOOSE } s \in Servers : \\
& \quad \wedge s.replicaId = LeaderID(m.gVec[myServerId.shardId]) \\
& \quad \wedge s.shardId = myServerId.shardId
\end{aligned}$$

IN

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

$$\begin{aligned}
& \wedge vGView[myServerId] < m.gView \\
& \wedge \text{IF } vServerStatus[myServerId] = StNormal \text{ THEN} \\
& \quad \wedge vServerStatus' = [vServerStatus \text{ EXCEPT } ![myServerId] = StViewChange] \\
& \quad \wedge vLastNormView' = [vLastNormView \text{ EXCEPT } ![myServerId] = vLView[myServerId]] \\
& \quad \text{ELSE UNCHANGED } \langle vServerStatus, vLastNormView \rangle \\
& \wedge vGView' = [ \\
& \quad \quad vGView \text{ EXCEPT } ![myServerId] = m.vGView \\
& ] \\
& \wedge vGVec' = [ \\
& \quad \quad vGVec \text{ EXCEPT } ![myServerId] = m.gVec \\
& ] \\
& \wedge vLView' = [ \\
& \quad \quad vLView \text{ EXCEPT } ![myServerId] = m.gVec[myServerId.shardId] \\
& ] \\
& \text{Clear early buffer,} \\
& \wedge vEarlyBuffer' = [ \\
& \quad \quad vEarlyBuffer \text{ EXCEPT } ![myServerId] = \{\} \\
& ] \\
& \text{Clear late buffer} \\
& \wedge vLateBuffer' = [ \\
& \quad \quad vLateBuffer \text{ EXCEPT } ![myServerId] = \{\} \\
& ] \\
& \text{Clear deadline quorum} \\
& \wedge vDeadlineQuorum' = [ \\
& \quad \quad vDeadlineQuorum \text{ EXCEPT } ![myServerId] = \{\} \\
& ] \\
& \text{Clear } vCrossShardConfirmQuorum \\
& \wedge vCrossShardConfirmQuorum' = [ \\
& \quad \quad serverId \in Servers \mapsto \{\} \\
& ]
\end{aligned}$$



$$\begin{aligned}
lSyncPoints &\triangleq \{m.lSyncPoint : m \in refinedQuorum\} \\
largestLSyncPointVC &\triangleq \text{CHOOSE } vc \in refinedQuorum : \\
&\quad \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 \\
&\quad \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 : \\
&\quad CountVotes(entry, logSets) \geq RecoveryQuorumSize\} \\
unSyncedLogSeq &\triangleq SetToSortSeq(unSyncedLogs, Compare) \\
\text{IN} \\
syncedLogSeq \circ unSyncedLogSeq \\
SelectEntriesBeyondCommitPoint(entries, deadline) &\triangleq \\
\text{LET} \\
\quad validLogIndices &\triangleq \{ \\
\quad \quad i \in 1 \dots Len(entries) : entries[i].deadline > deadline \\
\quad \} \\
\quad startIndex &\triangleq PickMin(validLogIndices) \\
\text{IN} \\
\quad SubSeq(entries, startIndex, Len(entries)) \\
HandleViewChange(m) &\triangleq \\
\text{LET} \\
\quad myServerId &\triangleq m.dest \\
\quad serversInOneShard &\triangleq \{s \in Servers : s.shardId = myServerId.shardId\} \\
\quad leadersInAllShard &\triangleq \{ \\
\quad \quad s \in Servers : s.replicaId = isLeader(s.replicaId, m.gVec[s.shardId]) \\
\quad \} \\
\text{IN} \\
\quad \wedge \vee ViewGreater(m.gView, m.lView, vGView[myServerId], vLView[myServerId]) \\
\quad \vee \wedge m.gView = vGView[myServerId] \\
\quad \wedge m.lView = vLView[myServerId] \\
\quad \wedge vServerStatus[myServerId] = StViewChange \\
\quad \wedge isLeader(myServerId.replicaId, m.lView) \\
\quad \wedge vGView' = [vGView \text{ EXCEPT } ![myServerId] = m.gView] \\
\quad \wedge vLView' = [vLView \text{ EXCEPT } ![myServerId] = m.lView] \\
\quad \wedge vGVec' = [vGVec \text{ EXCEPT } ![myServerId] = m.gVec] \\
\quad \wedge vViewChange' = [ \\
\quad \quad vViewChange \text{ EXCEPT } ![myServerId] = \{ \\
\quad \quad \quad vc \in vViewChange[myServerId] : \\
\quad \quad \quad vc.lView = m.lView \\
\quad \quad \} \cup \{m\} \\
\quad \quad ]
\end{aligned}$$

```

]
∧ IF Cardinality( $vViewChange'[myServerId]$ ) =  $QuorumSize$  THEN
  ∧  $vLog' = [vLog \text{ EXCEPT } ![myServerId] = ReBuildLogs(vViewChange'[myServerId])]$ 
  ∧  $vServerStatus' = [vServerStatus \text{ EXCEPT } ![myServerId] = StCrossShardSyncing]$ 
  ∧  $vLastNormView' = [vLastNormView \text{ EXCEPT } ![myServerId] = vLView[myServerId]]$ 
  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
  ∧  $vViewChange' = [vViewChange \text{ EXCEPT } ![myServerId] = \{\}]$ 
  ∧  $Send(\{[$ 
     $mtype \mapsto MCrossShardConfirm,$ 
     $sender \mapsto myServerId,$ 
     $dest \mapsto dst,$ 
     $lView \mapsto vLView'[myServerId],$ 
     $gView \mapsto vGView'[myServerId],$ 
     $entries \mapsto SelectEntriesBeyondCommitPoint($ 
       $vLog'[myServerId], vPeerCommitDeadline[dst.shardId])$ 
     $]: dst \in leadersInAllShard\})$ 
  ELSE
    ∧  $vServerStatus' = [vServerStatus \text{ EXCEPT } ![myServerId] = StViewChange]$ 
    ∧ UNCHANGED  $\langle networkVars, vLog, vServerStatus, vViewChange \rangle$ 

```

```

BuildGlobalConsistentLog( $serverId, entries$ )  $\triangleq$ 
  LET
     $myEntries \triangleq \{$ 
       $entry \in entries : \wedge serverId \in entry.shards$ 
       $\wedge \forall e \in entries :$ 
        IF  $e.txnId = entry.txnId$  THEN
           $e.deadline \leq entry.deadline$ 
        ELSE TRUE
     $\}$ 
  IN
    SetToSortSeq( $myEntries, Compare$ )

```

```

HandleCrossShardConfirm( $m$ )  $\triangleq$ 
  LET
     $myServerId \triangleq m.dest$ 
  IN
    ∧  $vServerStatus[myServerId] = StCrossShardSyncing$ 
    ∧  $m.gView = vGView[myServerId]$ 
    ∧  $m.lView = vGVec[myServerId][m.sender.shardId]$ 
    ∧  $vCrossShardConfirmQuorum' = [$ 
       $vCrossShardConfirmQuorum \text{ EXCEPT } ![myServerId] = \{$ 
         $mm \in vCrossShardConfirmQuorum[myServerId] :$ 
        ∧  $mm.gView = vGView[myServerId]$ 
       $\}$ 
     $]$ 

```

$$\begin{aligned}
& \wedge mm.lView = vGVec[myServerId][mm.sender.shardId] \\
& \quad \} \cup \{m\} \\
& ] \\
\wedge \text{ IF } & Cardinality(vCrossShardConfirmQuorum'[myServerId]) = Cardinality(Shards) \\
& \text{ THEN} \\
& \quad \text{Check } Txns' \text{ Deadlines to ensure strict serializability is not violated} \\
& \quad \text{In implementation, we should not pass all txns, instead, we should only pass dealines and } txn \text{ indices} \\
& \quad \text{As an optimization, we should also use checkpoint in implementation} \\
& \quad \text{Here for conciseness, we pass all } log \text{ entries} \\
& \text{ LET} \\
& \quad allLogs \triangleq \text{ UNION } \{SeqToSet(mm.entries) : \\
& \quad \quad mm \in vCrossShardConfirmQuorum'[myServerId]\} \\
& \quad serversInOneShard \triangleq \{s \in Servers : s.shardId = myServerId.shardId\} \\
& \text{ IN} \\
& \quad \wedge vLog' = [ \\
& \quad \quad vLog \text{ EXCEPT } ![myServerId] = \\
& \quad \quad \quad BuildGlobalConsistentLog(m.sender, allLogs) \\
& \quad ] \\
& \quad \wedge Send(\{[ \\
& \quad \quad mtype \quad \mapsto MStartView, \\
& \quad \quad sender \quad \mapsto myServerId, \\
& \quad \quad dest \quad \mapsto dst, \\
& \quad \quad lView \quad \mapsto vLView[myServerId], \\
& \quad \quad gView \quad \mapsto vGView[myServerId], \\
& \quad \quad gVec \quad \mapsto vGVec[myServerId], \\
& \quad \quad entries \mapsto vLog'[myServerId], \\
& \quad \quad cv \quad \mapsto vCrashVector[myServerId] \\
& \quad ] : dst \in serversInOneShard\}) \\
& \text{ ELSE} \\
& \quad \text{UNCHANGED } \langle vLog, networkVars \rangle \\
HandleStartView(m) & \triangleq \\
& \text{ LET} \\
& \quad myServerId \triangleq m.dest \\
& \text{ IN} \\
& \quad \wedge \vee ViewGreater(m.gView, m.lView, vGView[myServerId], vLView[myServerId]) \\
& \quad \vee \wedge m.gView = vGView[myServerId] \\
& \quad \wedge m.lView = vLView[myServerId] \\
& \quad \wedge \vee vServerStatus[myServerId] = StViewChange \\
& \quad \quad \vee vServerStatus[myServerId] = StRecovering \\
& \quad \wedge vGView' = [vGView \text{ EXCEPT } ![myServerId] = m.gView] \\
& \quad \wedge vLView' = [vLView \text{ EXCEPT } ![myServerId] = m.gLView] \\
& \quad \wedge vGVec' = [vGVec \text{ EXCEPT } ![myServerId] = m.vGVec] \\
& \quad \wedge vServerStatus' = [vServerStatus \text{ EXCEPT } ![myServerId] = StNormal] \\
& \quad \wedge vLog' = [vLog \text{ EXCEPT } ![myServerId] = m.entries]
\end{aligned}$$

$$\begin{aligned}
& \wedge vEarlyBuffer' = [vEarlyBuffer \text{ EXCEPT } ![myServerId] = \{\}] \\
& \wedge vLateBuffer' = [vLateBuffer \text{ EXCEPT } ![myServerId] = \{\}] \\
& \wedge vDeadlineQuorum' = [vDeadlineQuorum \text{ EXCEPT } ![myServerId] = \{\}] \\
& \wedge vCrossShardConfirmQuorum' = [ \\
& \quad vCrossShardConfirmQuorum \text{ EXCEPT } ![myServerId] = \{\} \\
& \quad ] \\
& \wedge vLSyncPoint' = [vLSyncPoint \text{ EXCEPT } ![myServerId] = Len(vLog'[myServerId])] \\
& \wedge vLastNormView' = [vLastNormView \text{ EXCEPT } ![myServerId] = m.lView] \\
& \wedge vViewChange' = [vViewChange \text{ EXCEPT } ![myServerId] = \{\}] \\
& \wedge vLSyncQuorum' = [vLSyncQuorum \text{ EXCEPT } ![myServerId] = \{\}] \\
& \wedge vCrashVectorReps' = [vCrashVectorReps \text{ EXCEPT } ![myServerId] = \{\}] \\
& \wedge vRecoveryReps' = [vRecoveryReps \text{ EXCEPT } ![myServerId] = \{\}]
\end{aligned}$$

$ResetServerState(serverId) \triangleq$

$$\begin{aligned}
& \wedge vLog' = [vLog \text{ EXCEPT } ![serverId] = \langle \rangle] \\
& \wedge vEarlyBuffer' = [vEarlyBuffer \text{ EXCEPT } ![serverId] = \{\}] \\
& \wedge vLateBuffer' = [vLateBuffer \text{ EXCEPT } ![serverId] = \{\}] \\
& \wedge vDeadlineQuorum' = [vDeadlineQuorum \text{ EXCEPT } ![serverId] = \{\}] \\
& \wedge vCrossShardConfirmQuorum' = [ \\
& \quad vCrossShardConfirmQuorum \text{ EXCEPT } ![serverId] = \{\} \\
& \quad ] \\
& \wedge vGView' = [vGView \text{ EXCEPT } ![serverId] = 0] \\
& \wedge vGVec' = [vGVec \text{ EXCEPT } ![serverId] = [s \in Shards \mapsto 0]] \\
& \wedge vLView' = [vLView \text{ EXCEPT } ![serverId] = 0] \\
& \wedge vLastNormView' = [vLastNormView \text{ EXCEPT } ![serverId] = 0] \\
& \wedge vViewChange' = [vViewChange \text{ EXCEPT } ![serverId] = \{\}] \\
& \wedge vLSyncPoint' = [vLSyncPoint \text{ EXCEPT } ![serverId] = 0] \\
& \wedge vLCommitPoint' = [vLCommitPoint \text{ EXCEPT } ![serverId] = 0] \\
& \wedge vPeerCommitDeadline' = [vPeerCommitDeadline \text{ EXCEPT } ![serverId] = 0] \\
& \wedge vLSyncQuorum' = [vLSyncQuorum \text{ EXCEPT } ![serverId] = \{\}] \\
& \wedge vCrashVector' = [vCrashVector \text{ EXCEPT } ![serverId] = [ \\
& \quad rr \in Replicas \mapsto 0 \\
& \quad ] \\
& \quad ] \\
& \wedge vCrashVectorReps' = [vCrashVectorReps \text{ EXCEPT } ![serverId] = \{\}] \\
& \wedge vRecoveryReps' = [vRecoveryReps \text{ EXCEPT } ![serverId] = \{\}] \\
& \wedge vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![serverId] = \{\}]
\end{aligned}$$

$StartServerRecovery(serverId) \triangleq$

$$\begin{aligned}
& \text{LET} \\
& \quad serversInOneShard \triangleq \{ \\
& \quad \quad s \in Servers : s.shardId = serverId.shardId \\
& \quad \quad \} \\
& \quad nonce \triangleq vUUIDCounter[serverId] + 1 \\
& \text{IN} \\
& \quad \wedge vServerStatus' = [vServerStatus \text{ EXCEPT } ![serverId] = StRecovering]
\end{aligned}$$



$$\begin{aligned}
& \wedge vUUIDCounter' = [vUUIDCounter \text{ EXCEPT } ![serverId] = vUUIDCounter[serverId] + 1] \\
& \wedge ResetServerState(serverId) \\
& \wedge Send(\{[ \\
& \quad mtype \quad \mapsto MCrashVectorReq, \\
& \quad sender \quad \mapsto serverId, \\
& \quad dest \quad \mapsto dst, \\
& \quad nonce \quad \mapsto nonce \\
& \quad ] : dst \in serversInOneShard\}) \\
\\
HandleCrashVectorReq(m) & \triangleq \\
\text{LET} & \\
\quad myServerId & \triangleq m.dest \\
\text{IN} & \\
& \wedge vServerStatus[myServerId] = StNormal \\
& \wedge Send(\{[ \\
& \quad mtype \quad \mapsto MCrashVectorRep, \\
& \quad sender \quad \mapsto myServerId, \\
& \quad dest \quad \mapsto m.sender, \\
& \quad nonce \quad \mapsto m.nonce, \\
& \quad cv \quad \mapsto vCrashVector[myServerId] \\
& \quad ]\}) \\
\\
AggregateCV(serverId) & \triangleq \\
\text{LET} & \\
\quad cvQuorum & \triangleq \{m.cv : m \in vCrashVectorReps[serverId]\} \\
\quad cvValQuorum & \triangleq [rr \in Replicas \mapsto \{cv[rr] : cv \in cvQuorum\}] \\
\text{IN} & \\
& [rr \in Replicas \mapsto PickMax(cvValQuorum[rr])] \\
\\
HandleCrashVectorRep(m) & \triangleq \\
\text{LET} & \\
\quad myServerId & \triangleq m.dest \\
\quad serversInOneShard & \triangleq \{s \in Servers : s.shardId = myServerId.shardId\} \\
\text{IN} & \\
& \wedge vServerStatus[myServerId] = StRecovering \\
& \wedge vUUIDCounter[myServerId] = m.nonce \\
& \wedge vCrashVectorReps' = [ \\
& \quad vCrashVectorReps \text{ EXCEPT } ![myServerId] = vCrashVectorReps \cup \{m\} \\
& \quad ] \\
& \wedge \text{IF } Cardinality(vCrashVectorReps'[myServerId]) = QuorumSize \text{ THEN} \\
& \quad \text{LET} \\
& \quad \quad acv \triangleq AggregateCV(myServerId) \\
& \quad \quad myCV \triangleq [acv \text{ EXCEPT } ![myServerId] = acv[myServerId] + 1] \\
& \quad \text{IN} \\
& \quad \wedge vCrashVector' = [
\end{aligned}$$

$$\begin{aligned}
& vCrashVector \text{ EXCEPT } ![myServerId] = myCV \\
& ] \\
\wedge & Send(\{[ \\
& \quad mtype \quad \mapsto MRecoveryReq, \\
& \quad sender \quad \mapsto myServerId, \\
& \quad dest \quad \mapsto dst, \\
& \quad cv \quad \mapsto myCV \\
& ] : dst \in serversInOneShard\}) \\
\text{ELSE} & \text{ UNCHANGED } \langle networkVars, vCrashVector \rangle
\end{aligned}$$

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

$$\begin{aligned}
HandleRecoveryRep(m) & \triangleq \\
\text{LET} & \\
& myServerId \triangleq m.dest \\
\text{IN} & \\
\wedge & vServerStatus[myServerId] = StRecovering \\
\wedge & vRecoveryReps' = [ \\
& \quad vRecoveryReps \text{ EXCEPT } ![myServerId] \\
& \quad = vRecoveryReps[myServerId] \cup \{m\} \\
& ] \\
\wedge & \text{ IF } Cardinality(vRecoveryReps[myServerId]) = QuorumSize \text{ THEN} \\
& \text{ LET} \\
& \quad lViewQuorum \triangleq \{mm.lView : mm \in vRecoveryReps[myServerId]\} \\
& \quad gViewQuorum \triangleq \{mm.gView : mm \in vRecoveryReps[myServerId]\} \\
& \text{ IN} \\
& \quad vLView' = [vLView \text{ EXCEPT } ![myServerId] = PickMax(lViewQuorum)] \\
& \quad vGView' = [vGView \text{ EXCEPT } ![myServerId] = PickMax(gViewQuorum)] \\
& \wedge Send(\{[ \\
& \quad mtype \quad \mapsto MStartViewReq, \\
& \quad sender \quad \mapsto myServerId,
\end{aligned}$$

$$\begin{array}{ll} dest & \mapsto [ \\ & replicaId \mapsto LeaderID(vLView[myServerId]), \\ & shardId \mapsto myServerId.shardId \\ & ], \\ lView & \mapsto vLView'[myServerId], \\ cv & \mapsto vCrashVector'[myServerId] \\ \}} \end{array}$$

ELSE UNCHANGED  $\langle networkVars, vLView, vGView \rangle$

$$\begin{array}{l}
\text{HandleStartViewReq}(m) \triangleq \\
\text{LET} \\
\quad myServerId \triangleq m.dest \\
\text{IN} \\
\wedge \ vServerStatus[myServerId] = StNormal \\
\wedge \ vLView[myServerId] = m.lView \\
\wedge \ isLeader(myServerId.replicaId, vLView[myServerId]) \\
\wedge \ Send(\{[ \\
\quad \quad mtype \quad \mapsto MStartView, \\
\quad \quad sender \quad \mapsto myServerId, \\
\quad \quad dest \quad \mapsto m.sender, \\
\quad \quad lView \quad \mapsto vLView[myServerId], \\
\quad \quad gView \quad \mapsto vGView[myServerId], \\
\quad \quad gVec \quad \mapsto vGVec[myServerId], \\
\quad \quad entries \quad \mapsto vLog[myServerId], \\
\quad \quad cv \quad \mapsto vCrashVector[myServerId] \\
\quad \quad \}])
\end{array}$$
$$\begin{array}{l}
StartLocalSync(serverId) \triangleq \\
\quad LET \\
\quad \quad leaderServerId \triangleq [ \\
\quad \quad \quad replicaId \mapsto LeaderID(vLView[serverId]), \\
\quad \quad \quad shardId \mapsto serverId.shardId \\
\quad \quad ] \\
\quad IN \\
\quad \wedge vServerStatus[serverId] = StNormal \\
\quad \wedge Send(\{[ \\
\quad \quad \quad mtype \mapsto MLocalSyncStatus, \\
\quad \quad \quad sender \mapsto serverId, \\
\quad \quad \quad dest \mapsto leaderServerId, \\
\quad \quad \quad lView \mapsto vLView[serverId], \\
\quad \quad \quad lSyncPoint \mapsto vLSyncPoint[serverId], \\
\quad \quad \quad cv \mapsto vCrashVector[serverId] \\
\quad \quad ])
\end{array}$$

$$\begin{aligned}
& \text{HandleLocalSyncStatus}(m) \triangleq \\
& \quad \text{LET} \\
& \quad \quad \text{myServerId} \triangleq m.\text{dest} \\
& \quad \quad \text{lSyncQuorum} \triangleq \text{vLSyncQuorum}[\text{myServerId}] \\
& \quad \text{IN} \\
& \quad \wedge \text{vServerStatus}[\text{myServerId}] = \text{StNormal} \\
& \quad \wedge \text{vLView}[\text{myServerId}] = m.\text{lView} \\
& \quad \wedge \text{isLeader}(\text{myServerId.replicaId}, \text{vLView}[\text{myServerId}]) \\
& \quad \wedge \forall mm \in \text{lSyncQuorum} : \\
& \quad \quad \vee mm.\text{sender} \neq m.\text{sender} \\
& \quad \quad \vee mm.\text{lSyncPoint} < m.\text{lSyncPoint} \\
& \quad \wedge \text{vLSyncQuorum}' = [ \\
& \quad \quad \text{vLSyncQuorum EXCEPT } ![\text{myServerId}] = \\
& \quad \quad \quad \{mm \in \text{lSyncQuorum} : mm.\text{sender} \neq m.\text{sender}\} \cup \{m\} \\
& \quad \quad ] \\
& \quad \wedge \text{IF } \text{Cardinality}(\text{vLSyncQuorum}'[\text{myServerId}]) \geq \text{QuorumSize} \text{ THEN} \\
& \quad \quad \text{LET} \\
& \quad \quad \quad \text{candidateQuorum} \triangleq \{ \\
& \quad \quad \quad \quad R \in \text{SUBSET } (\text{vLSyncQuorum}'[\text{myServerId}]) : \\
& \quad \quad \quad \quad \text{Cardinality}(R) = \text{QuorumSize} \\
& \quad \quad \quad \} \\
& \quad \quad \quad \text{quorumSyncPoints} \triangleq \{ \\
& \quad \quad \quad \quad \{x.\text{lSyncPoint} : x \in R\} : R \in \text{candidateQuorum} \\
& \quad \quad \quad \} \\
& \quad \quad \quad \text{validCommitPoints} \triangleq \{\text{PickMax}(Q) : Q \in \text{quorumSyncPoints}\} \\
& \quad \quad \quad \text{maxCommitPoint} \triangleq \text{PickMax}(\text{validCommitPoints}) \\
& \quad \quad \text{IN} \\
& \quad \quad \wedge \text{vLCommitPoint}' = [\text{vLCommitPoint EXCEPT } ![\text{myServerId}] = \text{maxCommitPoint}] \\
& \quad \quad \wedge \text{Send}(\{[ \\
& \quad \quad \quad \text{mtype} \quad \quad \mapsto \text{MLocalCommit}, \\
& \quad \quad \quad \text{sender} \quad \quad \mapsto \text{myServerId}, \\
& \quad \quad \quad \text{dest} \quad \quad \mapsto m.\text{sender}, \\
& \quad \quad \quad \text{lView} \quad \quad \mapsto \text{vLView}[\text{myServerId}], \\
& \quad \quad \quad \text{lCommitPoint} \mapsto \text{vLCommitPoint}'[\text{myServerId}], \\
& \quad \quad \quad \text{cv} \quad \quad \mapsto \text{vCrashVector}'[\text{myServerId}] \\
& \quad \quad \quad ]\}) \\
& \quad \quad \text{ELSE} \quad \text{UNCHANGED } \langle \text{vLCommitPoint}, \text{networkVars} \rangle
\end{aligned}$$

$$\begin{aligned}
& \text{HandleLocalCommit}(m) \triangleq \\
& \quad \text{LET} \\
& \quad \quad \text{myServerId} \triangleq m.\text{dest} \\
& \quad \text{IN} \\
& \quad \wedge \text{vServerStatus}[\text{myServerId}] = \text{StNormal} \\
& \quad \wedge \text{vLView}[\text{myServerId}] = m.\text{lView}
\end{aligned}$$

$\wedge \neg isLeader(myServerId.replicaId, vLView[myServerId])$   
 Make sure the *syncPoint* is large enough before updating *CommitPoint*  
 $\wedge$  IF  $\wedge vLSyncPoint[myServerId] \geq m.lCommitPoint$   
 $\wedge vLCommitPoint[myServerId] < m.lCommitPoint$   
 THEN  
 $vLCommitPoint' = [$   
 $\quad vLCommitPoint \text{ EXCEPT } ![myServerId] = m.lCommitPoint$   
 $]$   
 ELSE UNCHANGED  $\langle vLCommitPoint \rangle$

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

*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  
 $\quad \wedge vPeerCommitDeadline[myServerId]' = [$   
 $\quad \quad vPeerCommitDeadline[myServerId]$   
 $\quad \quad \text{EXCEPT } ![m.sender.shardId] = m.deadline$   
 $\quad ]$   
 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)

ReleaseSequencer(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  $\triangleq$ 
      [ i  $\in$  1 .. Len(sortedTxnList)
         $\mapsto$  isCommitting(sortedTxnList[i], vDeadlineQuorum[serverId])
      ]
    canReleaseTxnIndices  $\triangleq$  {
      i  $\in$  1 .. Len(sortedTxnList) :
         $\forall j \in$  1 .. i : committingStatus[j] = TRUE}
  IN
  IF Cardinality(canReleaseTxnIndices) = 0   Nothing to release
  THEN   UNCHANGED  $\langle$ networkVars,
    vLog, vEarlyBuffer, vLateBuffer, vDeadlineQuorum $\rangle$ 
  ELSE
    LET
      releaseUpTo  $\triangleq$  CHOOSE i  $\in$  canReleaseTxnIndices :
         $\forall j \in$  canReleaseTxnIndices : j  $\leq$  i
      releaseSeq  $\triangleq$  SubSeq(sortedTxnList, 1, releaseUpTo)
      releaseTxns  $\triangleq$  {releaseSeq[i] : i  $\in$  1 .. Len(releaseSeq)}
    IN
     $\wedge$  vEarlyBuffer' = [
      vEarlyBuffer EXCEPT ![serverId]
      = vEarlyBuffer[serverId] \ releaseTxns
    ]
     $\wedge$  vDeadlineQuorum' = [
      vDeadlineQuorum EXCEPT ![serverId]
      = {msg  $\in$  vDeadlineQuorum[serverId] :
         $\forall$  txn  $\in$  releaseTxns : txn.txnId  $\neq$  msg.txnId}
      ]
    Append to log
     $\wedge$  vLog' = [vLog EXCEPT ![serverId] = vLog[serverId]  $\circ$  releaseSeq]
     $\wedge$  IF isLeader(serverId.replicaId, vLView[serverId]) THEN
       $\wedge$  vLSyncPoint' = [vLSyncPoint EXCEPT ![serverId] = Len(vLog'[serverId])]
    ELSE   UNCHANGED  $\langle$ vLSyncPoint $\rangle$ 
    Send fast-replies to coordinators
     $\wedge$  Send({[

```

$mtype \mapsto MFastReply,$   
 $sender \mapsto serverId,$   
 $dest \mapsto sortedTxnList[i].txnId.coordId,$   
 $gView \mapsto vGView[serverId],$   
 $lView \mapsto vLView[serverId],$   
 $txnId \mapsto sortedTxnList[i].txnId,$   
 $hash \mapsto [$   
 $\quad log \mapsto vLog'[serverId],$   
 $\quad cv \mapsto vCrashVector$   
 $\quad ],$   
 $logId \mapsto i$   
 $\quad ] : i \in (1 + Len(vLog[serverId])) \dots Len(vLog'[serverId])\}$   
 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  
 $\wedge Send(\{[$   
 $\quad mtype \mapsto MInterReplicaSync,$   
 $\quad lView \mapsto vLView[serverId],$   
 $\quad sender \mapsto serverId,$   
 $\quad dest \mapsto dstServerId,$   
 $\quad entries \mapsto vLog'[serverId]$   
 $\quad ] : dstServerId \in serversInOneShard\})$

$ServerClockMove(serverId) \triangleq$   
 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  
          $\wedge ReleaseSequeuncer(serverId, vServerClock[serverId] + 1)$   
     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  $\langle vCoordClock \rangle$

$$\begin{aligned}
& \vee \wedge vCoordClock[coordId] < MaxTime \\
& \wedge vCoordClock' = [ \\
& \quad vCoordClock \text{ EXCEPT } ![coordId] = vCoordClock[coordId] + 1]
\end{aligned}$$

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

$$\begin{aligned}
Next & \triangleq \\
& \vee \wedge ActionName' = \langle \text{"Next"} \rangle \\
& \wedge \text{UNCHANGED } \langle networkVars, serverStateVars, \\
& \quad coordStateVars, configManagerStateVars \rangle \\
& \vee \exists c \in Coords : \\
& \quad \wedge Cardinality(vCoordTxns[c]) < MaxReqNum \\
& \quad \wedge CoordSubmitTxn(c) \\
& \quad \wedge \text{UNCHANGED } \langle serverStateVars, configManagerStateVars, \\
& \quad \quad vCoordProcessed \rangle \\
& \quad \wedge ActionName' = \langle \text{"CoordSubmitTxn"} \rangle
\end{aligned}$$

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

$$\begin{aligned}
& \vee \exists m \in messages : \\
& \quad \wedge m.mtype = MDeadlineNotification \\
& \quad \wedge vServerStatus[m.dest] = StNormal \\
& \quad \wedge m \notin vServerProcessed[m.dest] \\
& \quad \wedge vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] = \\
& \quad \quad vServerProcessed[m.dest] \cup \{m\}]
\end{aligned}$$



$$\begin{aligned}
& \wedge \text{HandleDeadlineNotification}(m) \\
& \wedge \text{UNCHANGED } \langle \text{networkVars}, \text{coordStateVars}, \text{configManagerStateVars}, \\
& \quad vLog, vCrossShardConfirmQuorum, vLateBuffer, \\
& \quad vServerStatus, vGView, vGVec, \\
& \quad vLView, vServerClock, vLastNormView, \\
& \quad vViewChange, vLSyncPoint, vLCommitPoint, \\
& \quad vPeerCommitDeadline, vLSyncQuorum, \\
& \quad vUUIDCounter, vCrashVector, vCrashVectorReps, \\
& \quad vRecoveryReps \rangle \\
& \wedge \text{ActionName}' = \langle \text{"HandleDeadlineNotification"} \rangle \\
\\
\vee \exists m \in \text{messages} : \\
& \wedge m.mtype = MInterReplicaSync \\
& \wedge vServerStatus[m.dest] = StNormal \\
& \wedge m \notin vServerProcessed[m.dest] \\
& \wedge vServerProcessed' = [vServerProcessed \text{ EXCEPT } ![m.dest] = \\
& \quad vServerProcessed[m.dest] \cup \{m\}] \\
& \wedge \text{HandleInterReplicaSync}(m) \\
& \wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{configManagerStateVars}, \\
& \quad vLog, vCrossShardConfirmQuorum, vLateBuffer, \\
& \quad vServerStatus, vGView, vGVec, \\
& \quad vLView, vServerClock, vLastNormView, \\
& \quad vViewChange, vLCommitPoint, vPeerCommitDeadline, \\
& \quad vLSyncQuorum, vUUIDCounter, vCrashVector, \\
& \quad vCrashVectorReps, vRecoveryReps \rangle \\
& \wedge \text{ActionName}' = \langle \text{"HandleInterReplicaSync"} \rangle \\
\\
\text{Some Leader(s) fail} \\
\vee \exists \text{serverId} \in \text{Servers} : \\
& \wedge vLView[\text{serverId}] < \text{MaxViews} \\
& \wedge \text{isLeader}(\text{serverId.replicaId}, vLView[\text{serverId}]) \\
& \wedge \text{StartLeaderFail}(\text{serverId}) \\
& \wedge \text{UNCHANGED } \langle \text{networkVars}, \text{coordStateVars}, \text{configManagerStateVars}, \\
& \quad vLog, vEarlyBuffer, vLateBuffer, \\
& \quad vDeadlineQuorum, vCrossShardConfirmQuorum, vGView, vGVec, \\
& \quad vLView, vServerClock, vLastNormView, \\
& \quad vViewChange, vLSyncPoint, vLCommitPoint, \\
& \quad vPeerCommitDeadline, vLSyncQuorum, \\
& \quad vUUIDCounter, vCrashVector, vCrashVectorReps, \\
& \quad vRecoveryReps, vServerProcessed \rangle \\
& \wedge \text{ActionName}' = \langle \text{"StartLeaderFail"} \rangle \\
\\
\text{configManager notices some leader(s) fail and launch view change} \\
\vee \exists \text{cmReplicaId} \in \text{Replicas} : \\
& \wedge \text{LaunchViewChange}(\text{cmReplicaId})
\end{aligned}$$

$$\begin{aligned}
& \wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{serverStateVars}, \text{configManagerStateVars} \rangle \\
& \wedge \text{ActionName}' = \langle \text{"LaunchViewChange"} \rangle \\
\vee \exists m \in \text{messages} : \\
& \wedge m.\text{mtype} = \text{MCMPrepare} \\
& \wedge m \notin v\text{CMProcessed}[m.\text{dest}] \\
& \wedge v\text{CMProcessed}' = [v\text{CMProcessed} \text{ EXCEPT } ![m.\text{dest}] = \\
& \quad v\text{CMProcessed}[m.\text{dest}] \cup \{m\}] \\
& \wedge \text{HandleCMPrepare}(m) \\
& \wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{serverStateVars} \rangle \\
& \wedge \text{ActionName}' = \langle \text{"HandleCMPrepare"} \rangle \\
\\
\vee \exists m \in \text{messages} : \\
& \wedge m.\text{mtype} = \text{MCMPrepareReply} \\
& \wedge m \notin v\text{CMProcessed}[m.\text{dest}] \\
& \wedge v\text{CMProcessed}' = [v\text{CMProcessed} \text{ EXCEPT } ![m.\text{dest}] = \\
& \quad v\text{CMProcessed}[m.\text{dest}] \cup \{m\}] \\
& \wedge \text{HandleCMPrepareReply}(m) \\
& \wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{serverStateVars}, \\
& \quad v\text{CMStatus}, v\text{CMView}, v\text{CMPrepareGInfo} \rangle \\
& \wedge \text{ActionName}' = \langle \text{"HandleCMPrepareReply"} \rangle \\
\\
\vee \exists m \in \text{messages} : \\
& \wedge m.\text{mtype} = \text{MCMCommit} \\
& \wedge m \notin v\text{CMProcessed}[m.\text{dest}] \\
& \wedge v\text{CMProcessed}' = [v\text{CMProcessed} \text{ EXCEPT } ![m.\text{dest}] = \\
& \quad v\text{CMProcessed}[m.\text{dest}] \cup \{m\}] \\
& \wedge \text{HandleCMCommit}(m) \\
& \wedge \text{UNCHANGED } \langle \text{networkVars}, \text{coordStateVars}, \text{serverStateVars}, \\
& \quad v\text{CMStatus}, v\text{CMView}, v\text{CMPrepareGInfo}, v\text{CMPrepareReps} \rangle \\
& \wedge \text{ActionName}' = \langle \text{"HandleCMCommit"} \rangle \\
\\
\vee \exists m \in \text{messages} : \\
& \wedge m.\text{mtype} = \text{MViewChangeReq} \\
& \wedge m \notin v\text{ServerProcessed}[m.\text{dest}] \\
& \wedge v\text{ServerProcessed}' = [v\text{ServerProcessed} \text{ EXCEPT } ![m.\text{dest}] = \\
& \quad v\text{ServerProcessed}[m.\text{dest}] \cup \{m\}] \\
& \wedge v\text{ServerStatus}[m.\text{dest}] \neq \text{StFailing} \\
& \wedge \text{HandleViewChangeReq}(m) \\
& \wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{configManagerStateVars}, \\
& \quad v\text{Log}, v\text{ServerClock}, v\text{ViewChange}, v\text{LSyncPoint}, \\
& \quad v\text{LCommitPoint}, v\text{LSyncQuorum}, v\text{PeerCommitDeadline}, \\
& \quad v\text{UUIDCounter}, v\text{CrashVector}, v\text{CrashVectorReps}, v\text{RecoveryReps} \rangle \\
& \wedge \text{ActionName}' = \langle \text{"HandleViewChangeReq"} \rangle
\end{aligned}$$

$\vee \exists m \in \text{messages} :$   
 $\wedge m.\text{mtype} = \text{MViewChange}$   
 $\wedge \text{isCrashVectorValid}(m)$   
 $\wedge m \notin v\text{ServerProcessed}[m.\text{dest}]$   
 $\wedge v\text{ServerProcessed}' = [v\text{ServerProcessed} \text{ EXCEPT } ![m.\text{dest}] =$   
 $\quad v\text{ServerProcessed}[m.\text{dest}] \cup \{m\}]$   
 $\wedge v\text{ServerStatus}[m.\text{dest}] \neq \text{StFailing}$   
 $\wedge \text{HandleViewChange}(m)$   
 $\wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{configManagerStateVars},$   
 $\quad vGVec, v\text{ServerClock}, v\text{LSyncPoint}, v\text{LastNormView},$   
 $\quad v\text{LCommitPoint}, v\text{PeerCommitDeadline}, v\text{LSyncQuorum},$   
 $\quad v\text{UUIDCounter}, v\text{CrashVector}, v\text{CrashVectorReps},$   
 $\quad v\text{RecoveryReps} \rangle$   
 $\wedge \text{ActionName}' = \langle \text{"HandleViewChange"} \rangle$

$\vee \exists m \in \text{messages} :$   
 $\wedge m.\text{mtype} = \text{MCrossShardConfirm}$   
 $\wedge m \notin v\text{ServerProcessed}[m.\text{dest}]$   
 $\wedge v\text{ServerProcessed}' = [v\text{ServerProcessed} \text{ EXCEPT } ![m.\text{dest}] =$   
 $\quad v\text{ServerProcessed}[m.\text{dest}] \cup \{m\}]$   
 $\wedge v\text{ServerStatus}[m.\text{dest}] = \text{StViewChange}$   
 $\wedge \text{HandleCrossShardConfirm}(m)$   
 $\wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{configManagerStateVars},$   
 $\quad vGVec, v\text{ServerClock}, v\text{LSyncPoint}, v\text{LastNormView},$   
 $\quad v\text{LCommitPoint}, v\text{PeerCommitDeadline}, v\text{LSyncQuorum},$   
 $\quad v\text{UUIDCounter}, v\text{CrashVector}, v\text{CrashVectorReps}, v\text{RecoveryReps} \rangle$   
 $\wedge \text{ActionName}' = \langle \text{"HandleCrossShardConfirm"} \rangle$

$\vee \exists m \in \text{messages} :$   
 $\wedge m.\text{mtype} = \text{MStartView}$   
 $\wedge \text{isCrashVectorValid}(m)$   
 $\wedge m \notin v\text{ServerProcessed}[m.\text{dest}]$   
 $\wedge v\text{ServerProcessed}' = [v\text{ServerProcessed} \text{ EXCEPT } ![m.\text{dest}] =$   
 $\quad v\text{ServerProcessed}[m.\text{dest}] \cup \{m\}]$   
 $\wedge \text{HandleStartView}(m)$   
 $\wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{configManagerStateVars},$   
 $\quad v\text{ServerClock}, v\text{LCommitPoint}, v\text{PeerCommitDeadline},$   
 $\quad v\text{UUIDCounter}, v\text{CrashVector} \rangle$   
 $\wedge \text{ActionName}' = \langle \text{"HandleStartView"} \rangle$

Failed server rejoin  
 $\vee \exists \text{serverId} \in \text{Servers} :$   
 $\wedge v\text{ServerStatus}[\text{serverId}] = \text{StFailing}$   
 $\wedge v\text{ServerStatus}' = [v\text{ServerStatus} \text{ EXCEPT } ![\text{serverId}] = \text{StRecovering}]$   
 $\wedge \text{ResetServerState}(\text{serverId})$   
 $\wedge \text{StartServerRecovery}(\text{serverId})$

$$\begin{aligned}
& \wedge \text{UNCHANGED } \langle \text{networkVars}, \text{coordStateVars}, \text{coordStateVars} \rangle \\
& \wedge \text{ActionName}' = \langle \text{"StartReplicaRecovery"} \rangle \\
\vee \quad & \exists m \in \text{messages} : \\
& \wedge m.\text{mtype} = \text{MCrashVectorReq} \\
& \wedge m \notin v\text{ServerProcessed}[m.\text{dest}] \\
& \wedge v\text{ServerProcessed}' = [v\text{ServerProcessed} \text{ EXCEPT } ![m.\text{dest}] = \\
& \quad v\text{ServerProcessed}[m.\text{dest}] \cup \{m\}] \\
& \wedge \text{HandleCrashVectorReq}(m) \\
& \wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{configManagerStateVars}, \\
& \quad v\text{Log}, v\text{EarlyBuffer}, v\text{LateBuffer}, v\text{DeadlineQuorum}, \\
& \quad v\text{CrossShardConfirmQuorum}, v\text{ServerStatus}, \\
& \quad v\text{GView}, v\text{GVec}, v\text{LView}, v\text{ServerClock}, v\text{LastNormView}, \\
& \quad v\text{ViewChange}, v\text{LSyncPoint}, v\text{LCommitPoint}, \\
& \quad v\text{PeerCommitDeadline}, v\text{LSyncQuorum}, v\text{UUIDCounter}, \\
& \quad v\text{CrashVector}, v\text{CrashVectorReps}, v\text{RecoveryReps} \rangle \\
& \wedge \text{ActionName}' = \langle \text{"HandleCrashVectorReq"} \rangle \\
\vee \quad & \exists m \in \text{messages} : \\
& \wedge m.\text{mtype} = \text{MCrashVectorRep} \\
& \wedge m \notin v\text{ServerProcessed}[m.\text{dest}] \\
& \wedge v\text{ServerProcessed}' = [v\text{ServerProcessed} \text{ EXCEPT } ![m.\text{dest}] = \\
& \quad v\text{ServerProcessed}[m.\text{dest}] \cup \{m\}] \\
& \wedge \text{HandleCrashVectorRep}(m) \\
& \wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{configManagerStateVars}, \\
& \quad v\text{Log}, v\text{EarlyBuffer}, v\text{LateBuffer}, \\
& \quad v\text{DeadlineQuorum}, v\text{CrossShardConfirmQuorum}, v\text{ServerStatus}, \\
& \quad v\text{GView}, v\text{GVec}, v\text{LView}, v\text{ServerClock}, v\text{LastNormView}, \\
& \quad v\text{ViewChange}, v\text{LSyncPoint}, v\text{LCommitPoint}, \\
& \quad v\text{PeerCommitDeadline}, v\text{LSyncQuorum}, \\
& \quad v\text{UUIDCounter}, v\text{CrashVectorReps}, v\text{RecoveryReps} \rangle \\
& \wedge \text{ActionName}' = \langle \text{"HandleCrashVectorRep"} \rangle \\
\vee \quad & \exists m \in \text{messages} : \\
& \wedge m.\text{mtype} = \text{MRecoveryReq} \\
& \wedge m \notin v\text{ServerProcessed}[m.\text{dest}] \\
& \wedge v\text{ServerProcessed}' = [v\text{ServerProcessed} \text{ EXCEPT } ![m.\text{dest}] = \\
& \quad v\text{ServerProcessed}[m.\text{dest}] \cup \{m\}] \\
& \wedge \text{isCrashVectorValid}(m) \\
& \wedge \text{HandleRecoveryReq}(m) \\
& \wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{configManagerStateVars}, \\
& \quad v\text{Log}, v\text{EarlyBuffer}, v\text{LateBuffer}, \\
& \quad v\text{DeadlineQuorum}, v\text{CrossShardConfirmQuorum}, v\text{ServerStatus}, \\
& \quad v\text{GView}, v\text{GVec}, v\text{LView}, v\text{ServerClock}, v\text{LastNormView}, \\
& \quad v\text{ViewChange}, v\text{LSyncPoint}, v\text{LCommitPoint}, \\
& \quad v\text{PeerCommitDeadline}, v\text{LSyncQuorum},
\end{aligned}$$

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

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

$$\begin{aligned}
& vDeadlineQuorum, vCrossShardConfirmQuorum, \\
& vGView, vGVec, vLView, vServerClock, vLastNormView, \\
& vViewChange, vLSyncPoint, vLCommitPoint, \\
& vPeerCommitDeadline, vLSyncQuorum, vUUIDCounter, \\
& vCrashVector, vCrashVectorReps, vRecoveryReps\} \\
& \wedge \text{ActionName}' = \langle \text{"HandlePeerShardCommitStatus"} \rangle
\end{aligned}$$

Clock Move

$$\begin{aligned}
& \vee \exists \text{serverId} \in \text{Servers} : \\
& \quad \wedge \text{ServerClockMove}(\text{serverId}) \\
& \quad \wedge \text{UNCHANGED} \langle \text{coordStateVars}, \text{configManagerStateVars} \rangle \\
& \quad \wedge \text{ActionName}' = \langle \text{"ServerClockMove"} \rangle \\
& \vee \exists \text{coordId} \in \text{Coords} : \\
& \quad \wedge \text{CoordClockMove}(\text{coordId}) \\
& \quad \wedge \text{UNCHANGED} \langle \text{networkVars}, \text{serverStateVars}, \text{configManagerStateVars}, \\
& \quad \quad vCoordTxns, vCoordProcessed \rangle \\
& \quad \wedge \text{ActionName}' = \langle \text{"CoordClockMove"} \rangle
\end{aligned}$$

$$\text{Spec} \triangleq \text{Init} \wedge \Box[\text{Next}] \langle \text{networkVars}, \text{serverStateVars}, \text{coordStateVars}, \text{configManagerStateVars}, \text{ActionName} \rangle$$

$$\text{ShardRecovered}(\text{shardId}, lViewID) \triangleq$$

$$\begin{aligned}
& \text{LET} \\
& \quad \text{serversInOneShard} \triangleq \{s \in \text{Servers} : s.\text{shardId} = \text{shardId}\} \\
& \quad \text{leaderServer} \triangleq [ \\
& \quad \quad \text{replicaId} \mapsto \text{LeaderID}(lViewID), \\
& \quad \quad \text{shardId} \mapsto \text{shardId} \\
& \quad ] \\
& \text{IN} \\
& \quad \wedge \exists RM \in \text{SUBSET}(\text{serversInOneShard}) : \\
& \quad \quad \wedge \text{Cardinality}(RM) \geq \text{QuorumSize} \\
& \quad \quad \wedge \text{leaderServer} \in RM \\
& \quad \quad \wedge \forall r \in RM : vServerStatus[r] = \text{StNormal} \\
& \quad \quad \wedge \forall r \in RM : vLastNormView[r] \geq lViewID
\end{aligned}$$

$$\text{CommittedInView}(v, \text{shardId}, \text{txnId}) \triangleq$$

$$\begin{aligned}
& \text{LET} \\
& \quad \text{serversInOneShard} \triangleq \{s \in \text{Servers} : s.\text{shardId} = \text{shardId}\} \\
& \quad \text{leaderServer} \triangleq [ \\
& \quad \quad \text{replicaId} \mapsto \text{LeaderID}(v), \\
& \quad \quad \text{shardId} \mapsto \text{shardId} \\
& \quad ]
\end{aligned}$$

```

replySet  $\triangleq$  {
  m  $\in$  messages :  $\wedge$   $\vee$  m.mtype = MFastReply
                   $\vee$  m.mtype = MSlowReply
                   $\wedge$  m.txnId = txnId
                   $\wedge$  m.sender  $\in$  serversInOneShard
                   $\wedge$  m.lView = v
}
IN
IF  $\forall$  reply  $\in$  replySet :
   $\vee$  reply.mtype  $\neq$  MFastReply
   $\vee$  reply.sender  $\neq$  leaderServer
THEN No leader's fast reply  $\rightarrow$  This txn is not committed
  FALSE
ELSE
  LET
    leaderReply  $\triangleq$  CHOOSE reply  $\in$  replySet :
       $\wedge$  reply.mtype = MFastReply
       $\wedge$  reply.sender = leaderServer
  IN
    Committed in Fast Path
     $\vee \exists$  fastQuorum  $\in$  SUBSET replySet :
       $\wedge$  leaderReply  $\in$  fastQuorum
       $\wedge$  Cardinality(fastQuorum) = FastQuorumSize
      All replies have the same hash (or it is a slow reply)
       $\wedge \forall$  reply  $\in$  fastQuorum :
         $\vee \wedge$  reply.mtype = MFastReply
         $\wedge$  reply.hash = leaderReply.hash
        Slow Reply can be used as fast reply
         $\vee$  reply.mtype = MSlowReply
    Committed in Slow Path
     $\vee \exists$  slowQuorum  $\in$  SUBSET replySet :
       $\wedge$  leaderReply  $\in$  slowQuorum
       $\wedge$  Cardinality(slowQuorum) = QuorumSize
       $\wedge \forall$  reply  $\in$  slowQuorum  $\setminus$  {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



$$\begin{aligned}
Durability &\triangleq \\
&\forall shardId \in Shards : \\
&\quad \forall v1, v2 \in 0 \dots MaxViews : \\
&\quad \quad \text{If a } txn \text{ is committed in lower view (} v1 \text{),} \\
&\quad \quad \text{it is impossible to make this request uncommitted in higher view} \\
&\quad \neg( \wedge v1 < v2 \\
&\quad \quad \wedge ShardRecovered(shardId, v2) \\
&\quad \quad \wedge \exists c \in Coords : \\
&\quad \quad \quad \exists txnId \in vCoordTxns[c] : \\
&\quad \quad \quad \quad \wedge CommittedInView(v1, shardId, txnId) \\
&\quad \quad \quad \quad \wedge \neg CommittedInView(v2, shardId, txnId) \\
&\quad )
\end{aligned}$$

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

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

```

        v1LeaderReply.hash ≠ 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 ≜
    LET
        allTxns ≜ UNION {vCoordTxns[c] : c ∈ Coords}
    IN
        ∀ shardId ∈ Shards :
            ∀ txnId1, txnId2 ∈ allTxns :
                IF txnId1 = txnId2 THEN TRUE
                ELSE
                    ∀ v1, v2 ∈ 1 .. MaxViews :
                        IF ∧ CommittedInView(v1, shardId, txnId1)
                           ∧ CommittedInView(v1, shardId, txnId2)
                        THEN
                            LET
                                v1LeaderReply ≜ CHOOSE m ∈ messages :
                                    ∧ m.mtype = MFastReply
                                    ∧ m.txnId = txnId1
                                    ∧ m.lView = v1
                                    ∧ m.sender.shardId = shardId
                                    ∧ m.sender.replicaId = LeaderID(v1)
                                v2LeaderReply ≜ CHOOSE m ∈ messages :
                                    ∧ m.mtype = MFastReply
                                    ∧ m.txnId = txnId2
                                    ∧ m.lView = v2
                                    ∧ m.sender.shardId = shardId
                                    ∧ m.sender.replicaId = LeaderID(v2)
                            IN
                                They cannot be committed in the same log position, regardless of the view
                                v1LeaderReply.logId ≠ 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 ≜
    LET
        allTxns ≜ UNION {vCoordTxns[c] : c ∈ Coords}

```

```

IN
 $\forall txnId1, txnId2 \in allTxns :$ 
  IF  $txnId1 = txnId2$  THEN TRUE
  ELSE
     $\forall v \in 1 \dots MaxViews :$ 
       $\forall shardId1, shardId2 \in Shards :$ 
        IF  $shardId1 = shardId2$  THEN TRUE
        ELSE
          IF  $\wedge CommittedInView(v, shardId1, txnId1)$ 
              $\wedge CommittedInView(v, shardId1, txnId2)$ 
              $\wedge CommittedInView(v, shardId2, txnId1)$ 
              $\wedge CommittedInView(v, shardId2, txnId2)$ 
          THEN
            LET
               $txn1\_LeaderReplyOnShard1 \triangleq$  CHOOSE  $m \in messages :$ 
                 $\wedge m.mtype = MFastReply$ 
                 $\wedge m.txnId = txnId1$ 
                 $\wedge m.lView = v$ 
                 $\wedge m.sender.shardId = shardId1$ 
                 $\wedge m.sender.replicaId = LeaderID(v)$ 
               $txn2\_LeaderReplyOnShard1 \triangleq$  CHOOSE  $m \in messages :$ 
                 $\wedge m.mtype = MFastReply$ 
                 $\wedge m.txnId = txnId2$ 
                 $\wedge m.lView = v$ 
                 $\wedge m.sender.shardId = shardId1$ 
                 $\wedge m.sender.replicaId = LeaderID(v)$ 
               $txn1\_LeaderReplyOnShard2 \triangleq$  CHOOSE  $m \in messages :$ 
                 $\wedge m.mtype = MFastReply$ 
                 $\wedge m.txnId = txnId1$ 
                 $\wedge m.lView = v$ 
                 $\wedge m.sender.shardId = shardId2$ 
                 $\wedge m.sender.replicaId = LeaderID(v)$ 
               $txn2\_LeaderReplyOnShard2 \triangleq$  CHOOSE  $m \in messages :$ 
                 $\wedge m.mtype = MFastReply$ 
                 $\wedge m.txnId = txnId2$ 
                 $\wedge m.lView = v$ 
                 $\wedge m.sender.shardId = shardId2$ 
                 $\wedge m.sender.replicaId = LeaderID(v)$ 
            IN
               $\vee \wedge txn1\_LeaderReplyOnShard1.logId > txn2\_LeaderReplyOnShard1.logId$ 
               $\wedge txn1\_LeaderReplyOnShard2.logId > txn2\_LeaderReplyOnShard2.logId$ 
               $\vee \wedge txn1\_LeaderReplyOnShard1.logId < txn2\_LeaderReplyOnShard1.logId$ 
               $\wedge txn1\_LeaderReplyOnShard2.logId < txn2\_LeaderReplyOnShard2.logId$ 
          ELSE TRUE

```

[