

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

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

$$\begin{aligned}
\text{ceilHalf}F &\triangleq \text{IF } (F \div 2) * 2 = F \text{ THEN } F \div 2 \text{ ELSE } (F + 1) \div 2 \\
\text{floorHalf}F &\triangleq F \div 2 \\
\text{QuorumSize} &\triangleq F + 1 \\
\text{FastQuorumSize} &\triangleq F + \text{ceilHalf}F + 1 \\
\text{RecoveryQuorumSize} &\triangleq \text{ceilHalf}F + 1 \\
\text{FastQuorums} &\triangleq \{R \in \text{SUBSET}(\text{Replicas}) : \\
&\quad \text{Cardinality}(R) \geq \text{FastQuorumSize}\} \\
\text{Quorums} &\triangleq \{R \in \text{SUBSET}(\text{Replicas}) : \\
&\quad \text{Cardinality}(R) * 2 > \text{Cardinality}(\text{Replicas})\}
\end{aligned}$$

Server Status

$$\begin{aligned}
\text{StNormal} &\triangleq 1 \\
\text{StViewChange} &\triangleq 2 \\
\text{StCrossShardSyncing} &\triangleq 3 \\
\text{StRecovering} &\triangleq 4 \\
\text{StFailing} &\triangleq 5
\end{aligned}$$

Message Types

$$\begin{aligned}
\text{MTxn} &\triangleq 1 \\
\text{MLogEntry} &\triangleq 2 \quad \text{Log entry, different from index, it includes command field, which can be large in practice} \\
\text{MDeadlineNotification} &\triangleq 3 \quad \text{Leaders send the message to other leaders for deadline agreement} \\
\text{MInterReplicaSync} &\triangleq 4 \quad \text{Synchronize within shard group (across replicas) to ensure strict serializability} \\
\text{MFastReply} &\triangleq 5 \quad \text{Fast Reply Message} \\
\text{MSlowReply} &\triangleq 6 \quad \text{Slow Reply Message}
\end{aligned}$$

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

$$\begin{aligned}
\text{MViewChangeReq} &\triangleq 7 \quad \text{Sent by config manager when leader/sequencer failure detected} \\
\text{MViewChange} &\triangleq 8 \quad \text{Sent to ACK view change} \\
\text{MStartView} &\triangleq 9 \quad \text{Sent by new leader to start view}
\end{aligned}$$

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

$$\begin{aligned}
\text{MLocalSyncStatus} &\triangleq 10 \quad \text{Sent by the leader to ensure log durability} \\
\text{MLocalCommit} &\triangleq 11 \quad \text{Sent by followers as ACK}
\end{aligned}$$

The following messages are used for periodic sync across sharding groups

This is an optional optimization to enable fast recovery

$$\text{MPeerShardCommitStatus} \triangleq 12$$

The following messages are mainly used for server recovery

$$\begin{aligned}
\text{MCrashVectorReq} &\triangleq 13 \\
\text{MCrashVectorRep} &\triangleq 14 \\
\text{MRecoveryReq} &\triangleq 15 \\
\text{MRecoveryRep} &\triangleq 16 \\
\text{MStartViewReq} &\triangleq 17
\end{aligned}$$

$MCrossShardConfirm \triangleq 19$

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

$MCMPPrepare \triangleq 20$

$MCMPPrepareReply \triangleq 21$

$MCMCommit \triangleq 22$

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 = [
  coordId ↦ i in (1 .. ),
  rId      ↦ i in (1 .. )
]
```

```
Txn = [
  mtype ↦ MTxn
  txnId ↦ TxnID,
  shards ↦ Shards,
  command ↦ command,
  st      ↦ sendTime,
  bound ↦ latencyBound
]
```

```
LogEntry = [
  mtype ↦ MLogEntry
  txnId ↦ TxnID,
  shards ↦ Shards,
  command ↦ command,
  deadline ↦ 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 ↦ MDeadlineNotification,
  gView ↦ 0 ... x
  lView ↦ 0 ... y
  sender ↦ src ∈ Servers,
  dest ↦ dst ∈ Servers,
  entry ↦ 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  $\mapsto$  vector of counters
]

StartViewReq = [
  mtype  $\mapsto$  MStartViewReq,
  sender  $\mapsto src \in Servers$ ,
  dest  $\mapsto dst \in Servers$ ,
  lView  $\mapsto 0 \dots x$ 
  cv  $\mapsto$  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  $\mapsto$  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  $\mapsto$  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  $\mapsto$  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  $\mapsto$   $0 \dots x$ 
    gView  $\mapsto$   $0 \dots x$ 
    gVec  $\mapsto$  [ $shardId \mapsto lView$ ]
  ]

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

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

```

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 *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 we make *CMStatus* always *StNormal*

vCMStatus,
vCMView,

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

vCMGInfo,
vCMPPrepareGInfo,

Config Manager: quorum of *CMPPrepareReplies*

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 $\text{messages} = \{\}$

InitServerState \triangleq

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

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

$$\begin{aligned}
\text{InitCoordState} &\triangleq \\
& \wedge \text{vCoordProcessed} = [c \in \text{Coords} \mapsto \{\}] \\
& \wedge \text{vCoordClock} = [c \in \text{Coords} \mapsto 1] \\
& \wedge \text{vCoordTrns} = [c \in \text{Coords} \mapsto \{\}] \\
\text{InitConfigManagerState} &\triangleq \\
& \wedge \text{vCMStatus} = [\\
& \quad \text{replicaId} \in \text{Replicas} \mapsto \text{StNormal} \\
&] \\
& \wedge \text{vCMView} = [\\
& \quad \text{replicaId} \in \text{Replicas} \mapsto 0 \\
&] \\
& \wedge \text{vCMGInfo} = [\\
& \quad \text{replicaId} \in \text{Replicas} \mapsto [\\
& \quad \quad gView \mapsto 0, \\
& \quad \quad gVec \mapsto [\text{shardId} \in \text{Shards} \mapsto 0] \\
&] \\
&] \\
& \wedge \text{vCMPPrepareGInfo} = [\\
& \quad \text{replicaId} \in \text{Replicas} \mapsto [\\
& \quad \quad gView \mapsto 0, \\
& \quad \quad gVec \mapsto [\text{shardId} \in \text{Shards} \mapsto 0] \\
&]
\end{aligned}$$

```

    ]
  ]
  ∧ vCMPPrepareReps = [
    replicaId ∈ Replicas ↦ {}
  ]
  ∧ vCMPProcessed = [
    replicaId ∈ Replicas ↦ {}
  ]
]

```

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

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

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

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

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

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

$IsInjective(s) \triangleq$

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 $e \notin Range(s)$

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

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

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

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

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

$ViewGreater(gv1, lv1, gv2, lv2) \triangleq$

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

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

$txnId \mapsto txnId,$

$command \mapsto "",$

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

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

$$\begin{aligned}
& HandleDeadlineNotification(m) \triangleq \\
& \text{LET} \\
& \quad myServerId \triangleq m.dest \\
& \quad quorum \triangleq \{ \\
& \quad \quad msg \in vDeadlineQuorum[myServerId]
\end{aligned}$$

```

      :  $\wedge \text{msg.entry.txnId} = m.\text{entry.txnId}$ 
       $\wedge \text{msg.gView} = m.\text{gView}$ 
       $\wedge m.\text{gView} = vGView[\text{myServerId}]$ 
    }  $\cup \{m\}$ 
  IN
    Only leader does deadline agreement
     $\wedge vGView[\text{myServerId}] = m.\text{gView}$ 
     $\wedge vGVec[\text{myServerId}][m.\text{sender.shardId}] = m.lView$ 
     $\wedge isLeader(\text{myServerId.replicaId}, vLView[\text{myServerId}])$ 
     $\wedge vDeadlineQuorum' = [$ 
       $vDeadlineQuorum \text{ EXCEPT } ![\text{myServerId}]$ 
       $= vDeadlineQuorum[\text{myServerId}] \cup \{m\}$ 
    ]
     $\wedge \text{IF } Cardinality(quorum) = Cardinality(m.\text{entry.shards})$ 
      THEN
        Deadline quorum established : Update the deadline of the txn in Sequencer
        LET
           $maxDeadlineTxn \triangleq$ 
          CHOOSE  $x \in quorum :$ 
           $\forall y \in quorum :$ 
             $y.\text{entry.deadline} \leq x.\text{entry.deadline}$ 
           $sequencingTxn \triangleq$ 
          CHOOSE  $x \in vEarlyBuffer[\text{myServerId}] :$ 
             $x.\text{txnId} = m.\text{entry.txnId}$ 
        IN
           $\text{IF } maxDeadlineTxn.\text{entry.deadline} > sequencingTxn.\text{deadline}$ 
            THEN
               $vEarlyBuffer' = [vEarlyBuffer \text{ EXCEPT } ![\text{myServerId}]$ 
               $= (vEarlyBuffer[\text{myServerId}] \setminus \{sequencingTxn\}) \cup \{maxDeadlineTxn.\text{entry}\}]$ 
            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

```

   $\wedge m.\text{crashVector}[m.sender] = vCrashVector[m.sender]$ 
   $\wedge \neg isLeader(m.dest.replicaId, vLView[m.dest])$ 
   $\wedge \text{LET}$ 
     $myServerId \triangleq m.dest$ 
     $syncedTxnIds \triangleq \{m.\text{entries}[i].\text{txnId} : i \in 1 \dots Len(m.\text{entries})\}$ 

```

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

$StartLeaderFail(serverId) \triangleq$

This leader fails

LET

$serversInOneShard \triangleq \{$

$s \in Servers : s.shardId = serverId.shardId$

$\}$

$aliveReplicas \triangleq \{$

$s \in serversInOneShard : \quad \wedge \ vServerStatus[s] = StNormal$
 $\quad \quad \quad \wedge \ s \neq serverId$

$\}$

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 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

$aliveReplicaId \triangleq \text{CHOOSE } replicaId \in Replicas :$

$vServerStatus[shardId][replicaId] = StNormal$

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

$$\begin{aligned}
& \wedge \text{vCMPPrepareGInfo}' = [\text{vCMPPrepareGInfo} \text{ EXCEPT } ![cmReplicaId] = \\
& \quad [\\
& \quad \quad gView \mapsto \text{vCMGInfo}[cmReplicaId].gView + 1, \\
& \quad \quad gVec \mapsto \text{newGVec} \\
& \quad] \\
&] \\
& \wedge \text{Send}(\{[\\
& \quad mtype \mapsto \text{MCMPPrepate}, \\
& \quad sender \mapsto cmReplicaId, \\
& \quad dest \mapsto dstRid, \\
& \quad cView \mapsto \text{vCMView}[cmReplicaId], \\
& \quad gView \mapsto \text{vCMPPrepareGInfo}'[cmReplicaId].gView, \\
& \quad gVec \mapsto \text{newGVec} \\
&] : dstRid \in \text{Replicas}\})
\end{aligned}$$

$$\begin{aligned}
& \text{LaunchViewChange}(cmReplicaId) \triangleq \\
& \text{IF } \wedge \text{isLeader}(cmReplicaId, \text{vCMView}[cmReplicaId]) \\
& \quad \wedge \text{DetectLeaderFail}(cmReplicaId) \\
& \text{THEN} \\
& \quad \text{PrepareViewChange}(cmReplicaId) \\
& \text{ELSE} \\
& \quad \text{UNCHANGED } \langle \text{networkVars} \rangle
\end{aligned}$$

$$\begin{aligned}
& \text{HandleCMPPrepare}(m) \triangleq \\
& \wedge m.cView = \text{vCMView}[m.dest] \\
& \wedge m.gView > \text{vCMGInfo}[m.dest].gView \\
& \wedge \text{vCMPPrepareGInfo}' = [\text{vCMPPrepareGInfo} \text{ EXCEPT } ![m.dest] = \\
& \quad [\\
& \quad \quad gView \mapsto m.gView, \\
& \quad \quad gVec \mapsto m.gVec \\
& \quad] \\
&] \\
& \wedge \text{Send}(\{[\\
& \quad mtype \mapsto \text{MCMPPrepateReply}, \\
& \quad sender \mapsto m.dest, \\
& \quad dest \mapsto m.src, \\
& \quad cView \mapsto m.cView, \\
& \quad gView \mapsto m.gView \\
&]\})
\end{aligned}$$

$$\begin{aligned}
& \text{HandleCMPPrepareReply}(m) \triangleq \\
& \wedge m.cView = \text{vCMView}[m.dest]
\end{aligned}$$

```

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

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

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

```

```

       $\wedge s.replicaId = LeaderID(m.gVec[myServerId.shardId])$ 
       $\wedge s.shardId = myServerId.shardId$ 
IN
  If the msg's view is lower, ignore
   $\wedge vGView[myServerId] < m.gView$ 
   $\wedge$  IF  $vServerStatus[myServerId] = StNormal$  THEN
     $\wedge vServerStatus' = [vServerStatus \text{ EXCEPT } ![myServerId] = StViewChange]$ 
     $\wedge vLastNormView' = [vLastNormView \text{ EXCEPT } ![myServerId] = vLView[myServerId]]$ 
    ELSE UNCHANGED  $\langle vServerStatus, vLastNormView \rangle$ 
   $\wedge vGView' = [$ 
     $vGView \text{ EXCEPT } ![myServerId] = m.vGView$ 
   $]$ 
   $\wedge vGVec' = [$ 
     $vGVec \text{ EXCEPT } ![myServerId] = m.gVec$ 
   $]$ 
   $\wedge vLView' = [$ 
     $vLView \text{ EXCEPT } ![myServerId] = m.gVec[myServerId.shardId]$ 
   $]$ 
  Clear early buffer,
   $\wedge vEarlyBuffer' = [$ 
     $vEarlyBuffer \text{ EXCEPT } ![myServerId] = \{\}$ 
   $]$ 
  Clear late buffer
   $\wedge vLateBuffer' = [$ 
     $vLateBuffer \text{ EXCEPT } ![myServerId] = \{\}$ 
   $]$ 
  Clear deadline quorum
   $\wedge vDeadlineQuorum' = [$ 
     $vDeadlineQuorum \text{ EXCEPT } ![myServerId] = \{\}$ 
   $]$ 
  Clear  $vCrossShardConfirmQuorum$ 
   $\wedge vCrossShardConfirmQuorum' = [$ 
     $serverId \in Servers \mapsto \{\}$ 
   $]$ 
  Send ViewChange to the myLeader
   $\wedge Send(\{$ 
     $mtype \mapsto MViewChange,$ 
     $sender \mapsto myServerId,$ 
     $dest \mapsto myLeader,$ 
     $gView \mapsto m.vGView,$ 
     $gVec \mapsto m.gVec,$ 
     $lView \mapsto vLView'[myServerId],$ 
     $lastNormal \mapsto vLastNormView'[myServerId],$ 
     $lSyncPoint \mapsto vLSyncPoint[myServerId],$ 
     $entries \mapsto vLog[myServerId],$ 

```

$$cv \mapsto vCrashVector[myServerId]$$

Define a comparison function based on the key

$$Compare(a, b) \triangleq$$

$$\begin{aligned} &\vee a.deadline < b.deadline \\ &\vee \wedge a.deadline = b.deadline \\ &\quad \wedge a.txnId.coordId < b.txnId.coordId \\ &\vee \wedge a.deadline = b.deadline \\ &\quad \wedge a.txnId.coordId = b.txnId.coordId \\ &\quad \wedge a.txnId.rId < b.txnId.rId \end{aligned}$$

$$isCrashVectorValid(m) \triangleq$$

$$\begin{aligned} &\wedge \forall rr \in Replicas : vCrashVector[m.dest][rr] \leq m.cv[rr] \\ &\wedge vCrashVector' = [\\ &\quad vCrashVector \text{ EXCEPT } ![m.dest] = [\\ &\quad \quad rr \in Replicas \mapsto \text{Max}(m.cv[rr], vCrashVector[m.dest][rr]) \\ &\quad] \\ &] \end{aligned}$$

$$CountVotes(entry, logSets) \triangleq$$

$$\begin{aligned} &\text{LET} \\ &\quad validCandidates \triangleq \{s \in logSets : \exists e \in s : \\ &\quad \quad \wedge e.deadline = entry.deadline \\ &\quad \quad \wedge e.txnId = entry.txnId\} \\ &\text{IN} \\ &\quad Cardinality(validCandidates) \end{aligned}$$

$$ReBuildLogs(vcQuorum) \triangleq$$

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

IN
 $syncedLogSeq \circ unsyncedLogSeq$
 $SelectEntriesBeyondCommitPoint(entries, deadline) \triangleq$
 LET
 $validLogIndices \triangleq \{$
 $i \in 1 \dots Len(entries) : entries[i].deadline > deadline$
 $\}$
 $startIndex \triangleq PickMin(validLogIndices)$
 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
 $\wedge \vee ViewGreater(m.gView, m.lView, vGView[myServerId], vLView[myServerId])$
 $\vee \wedge m.gView = vGView[myServerId]$
 $\wedge m.lView = vLView[myServerId]$
 $\wedge vServerStatus[myServerId] = StViewChange$
 $\wedge isLeader(myServerId.replicaId, m.lView)$
 $\wedge vGView' = [vGView \text{ EXCEPT } ![myServerId] = m.gView]$
 $\wedge vLView' = [vLView \text{ EXCEPT } ![myServerId] = m.lView]$
 $\wedge vGVec' = [vGVec \text{ EXCEPT } ![myServerId] = m.gVec]$
 $\wedge vViewChange' = [$
 $vViewChange \text{ EXCEPT } ![myServerId] = \{$
 $vc \in vViewChange[myServerId] :$
 $vc.lView = m.lView$
 $\} \cup \{m\}$
 $\]$
 $\wedge \text{ IF } Cardinality(vViewChange'[myServerId]) = QuorumSize \text{ THEN}$
 $\wedge vLog' = [vLog \text{ EXCEPT } ![myServerId] = ReBuildLogs(vViewChange'[myServerId])]$
 $\wedge vServerStatus' = [vServerStatus \text{ EXCEPT } ![myServerId] = StCrossShardSyncing]$
 $\wedge vLastNormView' = [vLastNormView \text{ EXCEPT } ![myServerId] = vLView[myServerId]]$
 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
 $\wedge vViewChange' = [vViewChange \text{ EXCEPT } ![myServerId] = \{\}]$
 $\wedge Send(\{[$
 $mttype \mapsto MCrossShardConfirm,$
 $sender \mapsto myServerId,$

```

    dest      ↦ dst,
    lView     ↦ vLView'[myServerId],
    gView     ↦ vGView'[myServerId],
    entries   ↦ SelectEntriesBeyondCommitPoint(
                  vLog'[myServerId], vPeerCommitDeadline[dst.shardId])
  ] : dst ∈ leadersInAllShard})
ELSE
  ∧ vServerStatus' = [vServerStatus EXCEPT ![myServerId] = StViewChange]
  ∧ UNCHANGED ⟨networkVars, vLog, vServerStatus, vViewChange⟩

BuildGlobalConsistentLog(serverId, entries) ≜
  LET
    myEntries ≜ {
      entry ∈ entries : ∧ serverId ∈ entry.shards
                      ∧ ∀ e ∈ entries :
                        IF e.txnId = entry.txnId THEN
                          e.deadline ≤ entry.deadline
                        ELSE TRUE
    }
  IN
    SetToSortSeq(myEntries, Compare)

HandleCrossShardConfirm(m) ≜
  LET
    myServerId ≜ m.dest
  IN
    ∧ vServerStatus[myServerId] = StCrossShardSyncing
    ∧ m.gView = vGView[myServerId]
    ∧ m.lView = vGVec[myServerId][m.sender.shardId]
    ∧ vCrossShardConfirmQuorum' = [
      vCrossShardConfirmQuorum EXCEPT ![myServerId] = {
        mm ∈ vCrossShardConfirmQuorum[myServerId] :
          ∧ mm.gView = vGView[myServerId]
          ∧ mm.lView = vGVec[myServerId][mm.sender.shardId]
      } ∪ {m}
    ]
  ∧ IF Cardinality(vCrossShardConfirmQuorum'[myServerId]) = Cardinality(Shards)
  THEN
    Check Txns' Deadlines to ensure strict serializability is not violated
    In implementation, we should not pass all txns, instead, we should only pass dealines and txn indices
    As an optimization, we should also use checkpoint in implementation
    Here for conciseness, we pass all log entries
  LET
    allLogs ≜ UNION {SeqToSet(mm.entries) :
                      mm ∈ vCrossShardConfirmQuorum'[myServerId]}

```

$$\begin{aligned}
& serversInOneShard \triangleq \{s \in Servers : s.shardId = myServerId.shardId\} \\
& \text{IN} \\
& \wedge vLog' = [\\
& \quad vLog \text{ EXCEPT } ![myServerId] = \\
& \quad \quad BuildGlobalConsistentLog(m.sender, allLogs) \\
& \quad] \\
& \wedge Send(\{[\\
& \quad mtype \quad \mapsto MStartView, \\
& \quad sender \quad \mapsto myServerId, \\
& \quad dest \quad \mapsto dst, \\
& \quad lView \quad \mapsto vLView[myServerId], \\
& \quad gView \quad \mapsto vGView[myServerId], \\
& \quad gVec \quad \mapsto vGVec[myServerId], \\
& \quad entries \mapsto vLog'[myServerId], \\
& \quad cv \quad \mapsto vCrashVector[myServerId] \\
& \quad] : dst \in serversInOneShard\}) \\
& \text{ELSE} \\
& \quad \text{UNCHANGED } \langle vLog, networkVars \rangle \\
& HandleStartView(m) \triangleq \\
& \quad \text{LET} \\
& \quad \quad myServerId \triangleq m.dest \\
& \quad \text{IN} \\
& \quad \wedge \vee ViewGreater(m.gView, m.lView, vGView[myServerId], vLView[myServerId]) \\
& \quad \quad \vee \wedge m.gView = vGView[myServerId] \\
& \quad \quad \wedge m.lView = vLView[myServerId] \\
& \quad \quad \wedge \vee vServerStatus[myServerId] = StViewChange \\
& \quad \quad \quad \vee vServerStatus[myServerId] = StRecovering \\
& \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.vGVec] \\
& \quad \wedge vServerStatus' = [vServerStatus \text{ EXCEPT } ![myServerId] = StNormal] \\
& \quad \wedge vLog' = [vLog \text{ EXCEPT } ![myServerId] = m.entries] \\
& \quad \wedge vEarlyBuffer' = [vEarlyBuffer \text{ EXCEPT } ![myServerId] = \{\}] \\
& \quad \wedge vLateBuffer' = [vLateBuffer \text{ EXCEPT } ![myServerId] = \{\}] \\
& \quad \wedge vDeadlineQuorum' = [vDeadlineQuorum \text{ EXCEPT } ![myServerId] = \{\}] \\
& \quad \wedge vCrossShardConfirmQuorum' = [\\
& \quad \quad vCrossShardConfirmQuorum \text{ EXCEPT } ![myServerId] = \{\} \\
& \quad] \\
& \quad \wedge vLSyncPoint' = [vLSyncPoint \text{ EXCEPT } ![myServerId] = Len(vLog'[myServerId])] \\
& \quad \wedge vLastNormView' = [vLastNormView \text{ EXCEPT } ![myServerId] = m.lView] \\
& \quad \wedge vViewChange' = [vViewChange \text{ EXCEPT } ![myServerId] = \{\}] \\
& \quad \wedge vLSyncQuorum' = [vLSyncQuorum \text{ EXCEPT } ![myServerId] = \{\}] \\
& \quad \wedge vCrashVectorReps' = [vCrashVectorReps \text{ EXCEPT } ![myServerId] = \{\}] \\
& \quad \wedge vRecoveryReps' = [vRecoveryReps \text{ EXCEPT } ![myServerId] = \{\}]
\end{aligned}$$

$$\begin{aligned}
& \text{ResetServerState}(\text{serverId}) \triangleq \\
& \quad \wedge \text{vLog}' = [\text{vLog} \text{ EXCEPT } ![\text{serverId}] = \langle \rangle] \\
& \quad \wedge \text{vEarlyBuffer}' = [\text{vEarlyBuffer} \text{ EXCEPT } ![\text{serverId}] = \{\}] \\
& \quad \wedge \text{vLateBuffer}' = [\text{vLateBuffer} \text{ EXCEPT } ![\text{serverId}] = \{\}] \\
& \quad \wedge \text{vDeadlineQuorum}' = [\text{vDeadlineQuorum} \text{ EXCEPT } ![\text{serverId}] = \{\}] \\
& \quad \wedge \text{vCrossShardConfirmQuorum}' = [\\
& \quad \quad \text{vCrossShardConfirmQuorum} \text{ EXCEPT } ![\text{serverId}] = \{\} \\
& \quad] \\
& \quad \wedge \text{vGView}' = [\text{vGView} \text{ EXCEPT } ![\text{serverId}] = 0] \\
& \quad \wedge \text{vGVec}' = [\text{vGVec} \text{ EXCEPT } ![\text{serverId}] = [s \in \text{Shards} \mapsto 0]] \\
& \quad \wedge \text{vLView}' = [\text{vLView} \text{ EXCEPT } ![\text{serverId}] = 0] \\
& \quad \wedge \text{vLastNormView}' = [\text{vLastNormView} \text{ EXCEPT } ![\text{serverId}] = 0] \\
& \quad \wedge \text{vViewChange}' = [\text{vViewChange} \text{ EXCEPT } ![\text{serverId}] = \{\}] \\
& \quad \wedge \text{vLSyncPoint}' = [\text{vLSyncPoint} \text{ EXCEPT } ![\text{serverId}] = 0] \\
& \quad \wedge \text{vLCommitPoint}' = [\text{vLCommitPoint} \text{ EXCEPT } ![\text{serverId}] = 0] \\
& \quad \wedge \text{vPeerCommitDeadline}' = [\text{vPeerCommitDeadline} \text{ EXCEPT } ![\text{serverId}] = 0] \\
& \quad \wedge \text{vLSyncQuorum}' = [\text{vLSyncQuorum} \text{ EXCEPT } ![\text{serverId}] = \{\}] \\
& \quad \wedge \text{vCrashVector}' = [\text{vCrashVector} \text{ EXCEPT } ![\text{serverId}] = [\\
& \quad \quad \text{rr} \in \text{Replicas} \mapsto 0 \\
& \quad] \\
& \quad \wedge \text{vCrashVectorReps}' = [\text{vCrashVectorReps} \text{ EXCEPT } ![\text{serverId}] = \{\}] \\
& \quad \wedge \text{vRecoveryReps}' = [\text{vRecoveryReps} \text{ EXCEPT } ![\text{serverId}] = \{\}] \\
& \quad \wedge \text{vServerProcessed}' = [\text{vServerProcessed} \text{ EXCEPT } ![\text{serverId}] = \{\}] \\
& \text{StartServerRecovery}(\text{serverId}) \triangleq \\
& \quad \text{LET} \\
& \quad \quad \text{serversInOneShard} \triangleq \{ \\
& \quad \quad \quad s \in \text{Servers} : s.\text{shardId} = \text{serverId}.\text{shardId} \\
& \quad \quad \} \\
& \quad \quad \text{nonce} \triangleq \text{vUUIDCounter}[\text{serverId}] + 1 \\
& \quad \text{IN} \\
& \quad \wedge \text{vServerStatus}' = [\text{vServerStatus} \text{ EXCEPT } ![\text{serverId}] = \text{StRecovering}] \\
& \quad \wedge \text{vUUIDCounter}' = [\text{vUUIDCounter} \text{ EXCEPT } ![\text{serverId}] = \text{vUUIDCounter}[\text{serverId}] + 1] \\
& \quad \wedge \text{ResetServerState}(\text{serverId}) \\
& \quad \wedge \text{Send}(\{[\\
& \quad \quad \text{mtype} \quad \mapsto \text{MCrashVectorReq}, \\
& \quad \quad \text{sender} \quad \mapsto \text{serverId}, \\
& \quad \quad \text{dest} \quad \mapsto \text{dst}, \\
& \quad \quad \text{nonce} \quad \mapsto \text{nonce} \\
& \quad] : \text{dst} \in \text{serversInOneShard}\}) \\
& \text{HandleCrashVectorReq}(m) \triangleq \\
& \quad \text{LET} \\
& \quad \quad \text{myServerId} \triangleq m.\text{dest} \\
& \quad \text{IN}
\end{aligned}$$

$$\begin{aligned}
& \wedge vServerStatus[myServerId] = StNormal \\
& \wedge Send(\{[\\
& \quad mtype \quad \mapsto MCrashVectorRep, \\
& \quad sender \quad \mapsto myServerId, \\
& \quad dest \quad \mapsto m.sender, \\
& \quad nonce \quad \mapsto m.nonce, \\
& \quad cv \quad \mapsto vCrashVector[myServerId] \\
& \quad]\})
\end{aligned}$$

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

$$HandleRecoveryReq(m) \triangleq$$

```

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

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

HandleStartViewReq(m)  $\triangleq$ 
LET

```

$$\begin{aligned}
& myServerId \triangleq m.dest \\
& \text{IN} \\
& \wedge vServerStatus[myServerId] = StNormal \\
& \wedge vLView[myServerId] = m.lView \\
& \wedge isLeader(myServerId.replicaId, vLView[myServerId]) \\
& \wedge Send(\{[\\
& \quad mtype \quad \mapsto MStartView, \\
& \quad sender \quad \mapsto myServerId, \\
& \quad dest \quad \mapsto m.sender, \\
& \quad lView \quad \mapsto vLView[myServerId], \\
& \quad gView \quad \mapsto vGView[myServerId], \\
& \quad gVec \quad \mapsto vGVec[myServerId], \\
& \quad entries \quad \mapsto vLog[myServerId], \\
& \quad cv \quad \mapsto vCrashVector[myServerId] \\
& \quad]\}) \\
\\
& StartLocalSync(serverId) \triangleq \\
& \quad LET \\
& \quad \quad leaderServerId \triangleq [\\
& \quad \quad \quad replicaId \mapsto LeaderID(vLView[serverId]), \\
& \quad \quad \quad shardId \mapsto serverId.shardId \\
& \quad \quad] \\
& \quad \text{IN} \\
& \quad \wedge vServerStatus[serverId] = StNormal \\
& \quad \wedge Send(\{[\\
& \quad \quad mtype \quad \mapsto MLocalSyncStatus, \\
& \quad \quad sender \quad \mapsto serverId, \\
& \quad \quad dest \quad \mapsto leaderServerId, \\
& \quad \quad lView \quad \mapsto vLView[serverId], \\
& \quad \quad lSyncPoint \mapsto vLSyncPoint[serverId], \\
& \quad \quad cv \quad \mapsto vCrashVector[serverId] \\
& \quad \quad]\}) \\
\\
& HandleLocalSyncStatus(m) \triangleq \\
& \quad LET \\
& \quad \quad myServerId \triangleq m.dest \\
& \quad \quad lSyncQuorum \triangleq vLSyncQuorum[myServerId] \\
& \quad \text{IN} \\
& \quad \wedge vServerStatus[myServerId] = StNormal \\
& \quad \wedge vLView[myServerId] = m.lView \\
& \quad \wedge isLeader(myServerId.replicaId, vLView[myServerId]) \\
& \quad \wedge \forall mm \in lSyncQuorum : \\
& \quad \quad \vee mm.sender \neq m.sender
\end{aligned}$$

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

$$\begin{aligned}
& HandleLocalCommit(m) \triangleq \\
& \quad \text{LET} \\
& \quad \quad myServerId \triangleq m.dest \\
& \quad \text{IN} \\
& \quad \wedge vServerStatus[myServerId] = StNormal \\
& \quad \wedge vLView[myServerId] = m.lView \\
& \quad \wedge \neg isLeader(myServerId.replicaId, vLView[myServerId]) \\
& \quad \text{Make sure the } syncPoint \text{ is large enough before updating } CommitPoint \\
& \quad \wedge \text{IF } \wedge vLSyncPoint[myServerId] \geq m.lCommitPoint \\
& \quad \quad \wedge vLCommitPoint[myServerId] < m.lCommitPoint \\
& \quad \quad \text{THEN} \\
& \quad \quad \quad vLCommitPoint' = [\\
& \quad \quad \quad \quad vLCommitPoint \text{ EXCEPT } ![myServerId] = m.lCommitPoint \\
& \quad \quad \quad] \\
& \quad \quad \text{ELSE} \quad \text{UNCHANGED } \langle vLCommitPoint \rangle
\end{aligned}$$

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

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

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

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

```

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.entry.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  $\mapsto$  [
        log  $\mapsto$  vLog'[serverId],
        cv  $\mapsto$  vCrashVector
      ],
      logId  $\mapsto$  i
    ]}

```

$] : i \in (1 + \text{Len}(v\text{Log}[\text{serverId}])) \dots \text{Len}(v\text{Log}'[\text{serverId}]))\}$
 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
 $\wedge \text{Send}(\{[$
 $\text{mtype} \mapsto M\text{InterReplicaSync},$
 $\text{lView} \mapsto v\text{LView}[\text{serverId}],$
 $\text{sender} \mapsto \text{serverId},$
 $\text{dest} \mapsto \text{dstServerId},$
 $\text{entries} \mapsto v\text{Log}'[\text{serverId}]$
 $] : \text{dstServerId} \in \text{serversInOneShard}\})$

$\text{ServerClockMove}(\text{serverId}) \triangleq$
 IF $v\text{ServerClock}[\text{serverId}] \geq \text{MaxTime}$ THEN
 UNCHANGED $\langle \text{networkVars}, \text{serverStateVars} \rangle$
 ELSE
 $\wedge v\text{ServerClock}' = [$
 $v\text{ServerClock} \text{ EXCEPT } ![\text{serverId}] = v\text{ServerClock}[\text{serverId}] + 1]$
 \wedge IF $v\text{ServerStatus}[\text{serverId}] = \text{StNormal}$ THEN
 $\wedge \text{ReleaseSequeuncer}(\text{serverId}, v\text{ServerClock}[\text{serverId}] + 1)$
 ELSE
 UNCHANGED $\langle \text{networkVars}, v\text{Log}, v\text{EarlyBuffer},$
 $v\text{LateBuffer}, v\text{DeadlineQuorum} \rangle$
 \wedge UNCHANGED $\langle v\text{CrossShardConfirmQuorum},$
 $v\text{ServerStatus}, v\text{GView}, v\text{GVec}, v\text{LView}, v\text{LastNormView},$
 $v\text{ViewChange}, v\text{LSyncPoint}, v\text{LCommitPoint},$
 $v\text{PeerCommitDeadline}, v\text{LSyncQuorum},$
 $v\text{UUIDCounter}, v\text{CrashVector}, v\text{CrashVectorReps},$
 $v\text{RecoveryReps}, v\text{ServerProcessed} \rangle$

$\text{CoordClockMove}(\text{coordId}) \triangleq$
 $\vee \wedge v\text{CoordClock}[\text{coordId}] \geq \text{MaxTime}$
 \wedge UNCHANGED $\langle v\text{CoordClock} \rangle$
 $\vee \wedge v\text{CoordClock}[\text{coordId}] < \text{MaxTime}$
 $\wedge v\text{CoordClock}' = [$
 $v\text{CoordClock} \text{ EXCEPT } ![\text{coordId}] = v\text{CoordClock}[\text{coordId}] + 1]$

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

$$\begin{aligned}
Next &\triangleq \\
&\vee \wedge ActionName' = \langle \text{"Next"} \rangle \\
&\quad \wedge \text{UNCHANGED } \langle networkVars, serverStateVars, \\
&\quad \quad \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 \quad vCoordProcessed \rangle \\
&\quad \wedge ActionName' = \langle \text{"CoordSubmitTxn"} \rangle \\
&\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 \quad vServerProcessed[m.dest] \cup \{m\}] \\
&\quad \wedge HandleTxn(m) \\
&\quad \wedge \text{UNCHANGED } \langle coordStateVars, configManagerStateVars, \\
&\quad \quad \quad vLog, vDeadlineQuorum, vCrossShardConfirmQuorum, \\
&\quad \quad \quad vServerStatus, vGView, vGVec, \\
&\quad \quad \quad vLView, vServerClock, vLastNormView, \\
&\quad \quad \quad vViewChange, vLSyncPoint, vLCommitPoint, \\
&\quad \quad \quad vPeerCommitDeadline, vLSyncQuorum, \\
&\quad \quad \quad vUUIDCounter, vCrashVector, \\
&\quad \quad \quad vCrashVectorReps, vRecoveryReps \rangle \\
&\quad \wedge ActionName' = \langle \text{"HandleTxn"} \rangle \\
&\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 \quad vServerProcessed[m.dest] \cup \{m\}] \\
&\quad \wedge HandleDeadlineNotification(m) \\
&\quad \wedge \text{UNCHANGED } \langle networkVars, coordStateVars, configManagerStateVars, \\
&\quad \quad \quad vLog, vCrossShardConfirmQuorum, vLateBuffer, \\
&\quad \quad \quad vServerStatus, vGView, vGVec, \\
&\quad \quad \quad vLView, vServerClock, vLastNormView, \\
&\quad \quad \quad vViewChange, vLSyncPoint, vLCommitPoint, \\
&\quad \quad \quad vPeerCommitDeadline, vLSyncQuorum, \\
&\quad \quad \quad vUUIDCounter, vCrashVector, vCrashVectorReps, \\
&\quad \quad \quad vRecoveryReps \rangle \\
&\quad \wedge ActionName' = \langle \text{"HandleDeadlineNotification"} \rangle \\
&\vee \exists m \in messages :
\end{aligned}$$

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

Some *Leader(s)* fail

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

Config Manager notices some *leader(s)* fail and launch view change

$$\begin{aligned}
& \vee \exists cmReplicaId \in Replicas : \\
& \quad \wedge LaunchViewChange(cmReplicaId) \\
& \quad \wedge \text{UNCHANGED } \langle coordStateVars, serverStateVars, configManagerStateVars \rangle \\
& \quad \wedge ActionName' = \langle \text{"LaunchViewChange"} \rangle
\end{aligned}$$

$$\begin{aligned}
& \vee \exists m \in messages : \\
& \quad \wedge m.mtype = MCMPPrepare \\
& \quad \wedge m \notin vCMPProcessed[m.dest] \\
& \quad \wedge vCMPProcessed' = [vCMPProcessed \text{ EXCEPT } ![m.dest] = \\
& \quad \quad vCMPProcessed[m.dest] \cup \{m\}] \\
& \quad \wedge HandleCMPPrepare(m) \\
& \quad \wedge \text{UNCHANGED } \langle coordStateVars, serverStateVars \rangle \\
& \quad \wedge ActionName' = \langle \text{"HandleCMPPrepare"} \rangle
\end{aligned}$$

$$\begin{aligned}
& \vee \exists m \in \text{messages} : \\
& \quad \wedge m.\text{mtype} = \text{MCMPPrepareReply} \\
& \quad \wedge m \notin v\text{CMPProcessed}[m.\text{dest}] \\
& \quad \wedge v\text{CMPProcessed}' = [v\text{CMPProcessed} \text{ EXCEPT } ![m.\text{dest}] = \\
& \quad \quad v\text{CMPProcessed}[m.\text{dest}] \cup \{m\}] \\
& \quad \wedge \text{HandleCMPPrepareReply}(m) \\
& \quad \wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{serverStateVars}, \\
& \quad \quad v\text{CMStatus}, v\text{CMView}, v\text{CMPPrepareGInfo} \rangle \\
& \quad \wedge \text{ActionName}' = \langle \text{"HandleCMPPrepareReply"} \rangle \\
\\
& \vee \exists m \in \text{messages} : \\
& \quad \wedge m.\text{mtype} = \text{MCMCommit} \\
& \quad \wedge m \notin v\text{CMPProcessed}[m.\text{dest}] \\
& \quad \wedge v\text{CMPProcessed}' = [v\text{CMPProcessed} \text{ EXCEPT } ![m.\text{dest}] = \\
& \quad \quad v\text{CMPProcessed}[m.\text{dest}] \cup \{m\}] \\
& \quad \wedge \text{HandleCMCommit}(m) \\
& \quad \wedge \text{UNCHANGED } \langle \text{networkVars}, \text{coordStateVars}, \text{serverStateVars}, \\
& \quad \quad v\text{CMStatus}, v\text{CMView}, v\text{CMPPrepareGInfo}, v\text{CMPPrepareReps} \rangle \\
& \quad \wedge \text{ActionName}' = \langle \text{"HandleCMCommit"} \rangle \\
\\
& \vee \exists m \in \text{messages} : \\
& \quad \wedge m.\text{mtype} = \text{MViewChangeReq} \\
& \quad \wedge m \notin v\text{ServerProcessed}[m.\text{dest}] \\
& \quad \wedge v\text{ServerProcessed}' = [v\text{ServerProcessed} \text{ EXCEPT } ![m.\text{dest}] = \\
& \quad \quad v\text{ServerProcessed}[m.\text{dest}] \cup \{m\}] \\
& \quad \wedge v\text{ServerStatus}[m.\text{dest}] \neq \text{StFailing} \\
& \quad \wedge \text{HandleViewChangeReq}(m) \\
& \quad \wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{configManagerStateVars}, \\
& \quad \quad v\text{Log}, v\text{ServerClock}, v\text{ViewChange}, v\text{LSyncPoint}, \\
& \quad \quad v\text{LCommitPoint}, v\text{LSyncQuorum}, v\text{PeerCommitDeadline}, \\
& \quad \quad v\text{UUIDCounter}, v\text{CrashVector}, v\text{CrashVectorReps}, v\text{RecoveryReps} \rangle \\
& \quad \wedge \text{ActionName}' = \langle \text{"HandleViewChangeReq"} \rangle \\
\\
& \vee \exists m \in \text{messages} : \\
& \quad \wedge m.\text{mtype} = \text{MViewChange} \\
& \quad \wedge \text{isCrashVectorValid}(m) \\
& \quad \wedge m \notin v\text{ServerProcessed}[m.\text{dest}] \\
& \quad \wedge v\text{ServerProcessed}' = [v\text{ServerProcessed} \text{ EXCEPT } ![m.\text{dest}] = \\
& \quad \quad v\text{ServerProcessed}[m.\text{dest}] \cup \{m\}] \\
& \quad \wedge v\text{ServerStatus}[m.\text{dest}] \neq \text{StFailing} \\
& \quad \wedge \text{HandleViewChange}(m) \\
& \quad \wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{configManagerStateVars}, \\
& \quad \quad v\text{GVec}, v\text{ServerClock}, v\text{LSyncPoint}, v\text{LastNormView}, \\
& \quad \quad v\text{LCommitPoint}, v\text{PeerCommitDeadline}, v\text{LSyncQuorum},
\end{aligned}$$

$$\begin{aligned}
& vUUIDCounter, vCrashVector, vCrashVectorReps, \\
& vRecoveryReps) \\
& \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] = \\
& \quad vServerProcessed[m.dest] \cup \{m\}] \\
& \wedge vServerStatus[m.dest] = StViewChange \\
& \wedge HandleCrossShardConfirm(m) \\
& \wedge \text{UNCHANGED } \langle coordStateVars, configManagerStateVars, \\
& \quad vGVec, vServerClock, vLSyncPoint, vLastNormView, \\
& \quad vLCommitPoint, vPeerCommitDeadline, vLSyncQuorum, \\
& \quad 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] = \\
& \quad vServerProcessed[m.dest] \cup \{m\}] \\
& \wedge HandleStartView(m) \\
& \wedge \text{UNCHANGED } \langle coordStateVars, configManagerStateVars, \\
& \quad vServerClock, vLCommitPoint, vPeerCommitDeadline, \\
& \quad vUUIDCounter, vCrashVector \rangle \\
& \wedge ActionName' = \langle \text{"HandleStartView"} \rangle \\
\text{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] = \\
& \quad vServerProcessed[m.dest] \cup \{m\}] \\
& \wedge HandleCrashVectorReq(m) \\
& \wedge \text{UNCHANGED } \langle coordStateVars, configManagerStateVars, \\
& \quad vLog, vEarlyBuffer, vLateBuffer, vDeadlineQuorum,
\end{aligned}$$

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

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

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

Periodic Sync

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

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

$\wedge \text{ActionName}' = \langle \text{"HandleLocalSyncStatus"} \rangle$

$\vee \exists m \in \text{messages} :$
 $\wedge m.\text{mtype} = \text{MLocalCommit}$
 $\wedge m \notin v\text{ServerProcessed}[m.\text{dest}]$
 $\wedge v\text{ServerProcessed}' = [v\text{ServerProcessed} \text{ EXCEPT } ![m.\text{dest}] =$
 $\quad v\text{ServerProcessed}[m.\text{dest}] \cup \{m\}]$
 $\wedge \text{isCrashVectorValid}(m)$
 $\wedge \text{HandleLocalCommit}(m)$
 $\wedge \text{UNCHANGED } \langle \text{coordStateVars}, \text{configManagerStateVars},$
 $\quad \text{networkVars}, v\text{Log}, v\text{EarlyBuffer}, v\text{LateBuffer},$
 $\quad v\text{DeadlineQuorum}, v\text{CrossShardConfirmQuorum},$
 $\quad v\text{ServerStatus}, v\text{ServerClock},$
 $\quad v\text{GView}, v\text{GVec}, v\text{LView}, v\text{LastNormView},$
 $\quad v\text{ViewChange}, v\text{LSyncPoint}, v\text{PeerCommitDeadline},$
 $\quad v\text{LSyncQuorum}, v\text{UUIDCounter},$
 $\quad v\text{CrashVectorReps}, v\text{RecoveryReps} \rangle$
 $\wedge \text{ActionName}' = \langle \text{"HandleLocalCommit"} \rangle$

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

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

Clock Move

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

$$\begin{aligned}
& \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 \text{vCoordTxns}, \text{vCoordProcessed} \rangle \\
& \quad \wedge \text{ActionName}' = \langle \text{"CoordClockMove"} \rangle \\
\text{Spec} & \triangleq \text{Init} \wedge \Box[\text{Next}] \langle \text{networkVars}, \\
& \quad \text{serverStateVars}, \text{coordStateVars}, \text{configManagerStateVars}, \\
& \quad \text{ActionName} \rangle \\
\\
\text{ShardRecovered}(\text{shardId}, \text{lViewID}) & \triangleq \\
\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}(\text{lViewID}), \\
& \quad \quad \text{shardId} \mapsto \text{shardId} \\
& \quad] \\
\text{IN} & \\
& \quad \wedge \exists \text{RM} \in \text{SUBSET}(\text{serversInOneShard}) : \\
& \quad \quad \wedge \text{Cardinality}(\text{RM}) \geq \text{QuorumSize} \\
& \quad \quad \wedge \text{leaderServer} \in \text{RM} \\
& \quad \quad \wedge \forall r \in \text{RM} : \text{vServerStatus}[r] = \text{StNormal} \\
& \quad \quad \wedge \forall r \in \text{RM} : \text{vLastNormView}[r] \geq \text{lViewID} \\
\\
\text{CommittedInView}(v, \text{shardId}, \text{txnId}) & \triangleq \\
\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] \\
& \quad \text{replySet} \triangleq \{ \\
& \quad \quad m \in \text{messages} : \wedge \vee m.\text{mtype} = \text{MFastReply} \\
& \quad \quad \quad \vee m.\text{mtype} = \text{MSlowReply} \\
& \quad \quad \quad \wedge m.\text{txnId} = \text{txnId} \\
& \quad \quad \quad \wedge m.\text{sender} \in \text{serversInOneShard} \\
& \quad \quad \quad \wedge m.\text{lView} = v \\
& \quad \quad \} \\
\text{IN} & \\
& \quad \text{IF } \forall \text{reply} \in \text{replySet} : \\
& \quad \quad \vee \text{reply}.\text{mtype} \neq \text{MFastReply} \\
& \quad \quad \vee \text{reply}.\text{sender} \neq \text{leaderServer} \\
& \quad \text{THEN } \text{No leader's fast reply} \rightarrow \text{This txn is not committed} \\
& \quad \text{FALSE}
\end{aligned}$$

```

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

```

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

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

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

$$\begin{aligned}
&Linearizability \triangleq \\
&\quad \text{LET} \\
&\quad \quad allTxns \triangleq \text{UNION } \{vCoordTxns[c] : c \in Coords\} \\
&\quad \text{IN} \\
&\quad \quad \forall shardId \in Shards : \\
&\quad \quad \quad \forall txnId1, txnId2 \in allTxns :
\end{aligned}$$

```

IF  $txnId1 = txnId2$  THEN TRUE
ELSE
   $\forall v1, v2 \in 1 \dots MaxViews :$ 
    IF  $\wedge CommittedInView(v1, shardId, txnId1)$ 
        $\wedge CommittedInView(v1, shardId, txnId2)$ 
    THEN
      LET
         $v1LeaderReply \triangleq$  CHOOSE  $m \in messages :$ 
           $\wedge m.mtype = MFastReply$ 
           $\wedge m.txnId = txnId1$ 
           $\wedge m.lView = v1$ 
           $\wedge m.sender.shardId = shardId$ 
           $\wedge m.sender.replicaId = LeaderID(v1)$ 
         $v2LeaderReply \triangleq$  CHOOSE  $m \in messages :$ 
           $\wedge m.mtype = MFastReply$ 
           $\wedge m.txnId = txnId2$ 
           $\wedge m.lView = v2$ 
           $\wedge m.sender.shardId = shardId$ 
           $\wedge m.sender.replicaId = LeaderID(v2)$ 
      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
       $\forall v \in 1 \dots MaxViews :$ 
         $\forall shardId1, shardId2 \in Shards :$ 
          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

```

$$\begin{aligned}
& \text{txn1_LeaderReplyOnShard1} \triangleq \text{CHOOSE } m \in \text{messages} : \\
& \quad \wedge m.\text{mtype} = \text{MFastReply} \\
& \quad \wedge m.\text{txnId} = \text{txnId1} \\
& \quad \wedge m.\text{lView} = v \\
& \quad \wedge m.\text{sender.shardId} = \text{shardId1} \\
& \quad \wedge m.\text{sender.replicaId} = \text{LeaderID}(v) \\
& \text{txn2_LeaderReplyOnShard1} \triangleq \text{CHOOSE } m \in \text{messages} : \\
& \quad \wedge m.\text{mtype} = \text{MFastReply} \\
& \quad \wedge m.\text{txnId} = \text{txnId2} \\
& \quad \wedge m.\text{lView} = v \\
& \quad \wedge m.\text{sender.shardId} = \text{shardId1} \\
& \quad \wedge m.\text{sender.replicaId} = \text{LeaderID}(v) \\
& \text{txn1_LeaderReplyOnShard2} \triangleq \text{CHOOSE } m \in \text{messages} : \\
& \quad \wedge m.\text{mtype} = \text{MFastReply} \\
& \quad \wedge m.\text{txnId} = \text{txnId1} \\
& \quad \wedge m.\text{lView} = v \\
& \quad \wedge m.\text{sender.shardId} = \text{shardId2} \\
& \quad \wedge m.\text{sender.replicaId} = \text{LeaderID}(v) \\
& \text{txn2_LeaderReplyOnShard2} \triangleq \text{CHOOSE } m \in \text{messages} : \\
& \quad \wedge m.\text{mtype} = \text{MFastReply} \\
& \quad \wedge m.\text{txnId} = \text{txnId2} \\
& \quad \wedge m.\text{lView} = v \\
& \quad \wedge m.\text{sender.shardId} = \text{shardId2} \\
& \quad \wedge m.\text{sender.replicaId} = \text{LeaderID}(v) \\
& \text{IN} \\
& \vee \wedge \text{txn1_LeaderReplyOnShard1.logId} > \text{txn2_LeaderReplyOnShard1.logId} \\
& \quad \wedge \text{txn1_LeaderReplyOnShard2.logId} > \text{txn2_LeaderReplyOnShard2.logId} \\
& \vee \wedge \text{txn1_LeaderReplyOnShard1.logId} < \text{txn2_LeaderReplyOnShard1.logId} \\
& \quad \wedge \text{txn1_LeaderReplyOnShard2.logId} < \text{txn2_LeaderReplyOnShard2.logId} \\
& \text{ELSE TRUE}
\end{aligned}$$
