# Hovercraft++ TLA+ Specification

Ovidiu-Cristian Marcu[1*]

[1*]Computer Science, University of Luxembourg,  Luxembourg.

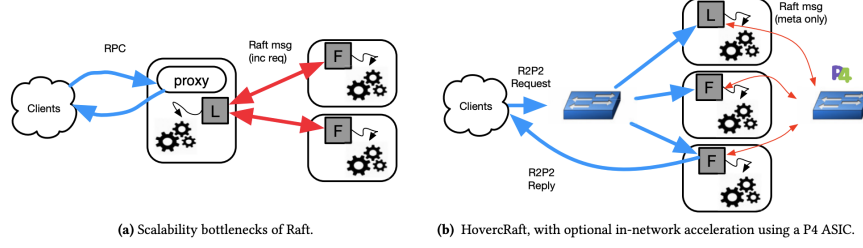Corresponding author(s). E-mail(s): ovidiu21marcu@gmail.com;

## Abstract

We introduce a formal TLA+ specification of the Hovercraft++ consensus protocol, derived from its natural language description in *HovercRaft: Achieving Scalability and Fault-tolerance for microsecond-scale Datacenter Services*. While the original description provides a strong foundation, our TLA+ model facilitates a deeper, rigorous analysis. We focus on the core consensus mechanism's fundamental safety properties, while excluding load balancing components (client replies, read-only operations) and bounded queue optimizations. Model checking the specification allows us to uncover and verify behavior in corner cases not overtly addressed in the original paper. Furthermore, this formalization process itself clarifies the intricacies of Raft-based consensus (upon which our specification, like the original TLA+ Raft specification it extends, is built), offering insights beyond even what standard Raft specifications typically detail. This experience reinforces our recommendation for researchers to invest in formal protocol specification with detailed commentary, as it significantly clarifies contributions, aids potential implementations, and ultimately strengthens protocol design.

**Keywords:** consensus, hovercraft++, TLA+, specification

## 1  Background: the Hovercraft++ Consensus Protocol

Hovercraft++ [1], *"an approach by which adding nodes increases both the resilience and the performance of general-purpose state-machine replication (SMR)"*, is an extension of the Raft [2] consensus protocol, designed to enhance the efficiency and scalability of SMR. SMR and its bottlenecks are described in subsections 2.1.1 and 2.1.2 of [1]. This paper investigates how Hovercraft++ preserves Raft's core safety while liveness guarantees investigation is left for future work. For the purpose of formal specification, its key provisions are:

**(a)** Scalability bottlenecks of Raft.　　　　**(b)** HovercRaft, with optional in-network acceleration using a P4 ASIC.

**Fig. 1** "Eliminating bottlenecks of SMR. Figure 1a shows the leader node bottlenecks for a classic SMR deployment using Raft: (1) the leader acts as the RPC server for all clients; (2) the leader must communicate individually with each follower to replicate messages and ensure their ordering. Figure 1b illustrates Hovercraft++, extending Raft to separate request replication from ordering and using IP multicast (Switch) and in-network accelerators (NetAgg) to convert leader-to-multipoint interactions into point-to-point interactions. Illustration on a 3-node cluster" reproduced from [1].

- **Transparent SMR Integration:** It integrates an extension of Raft directly within the R2P2 [3] transport layer, enabling applications to utilize SMR without modification to their core logic for handling consensus. "Specifically, the SMR layer becomes part of the RPC layer which forwards RPC requests to the application layer only after those requests have been totally ordered and committed by the leader" [1]. Design details provided in Section 3 of [1].
- **Modified Replication and Ordering Mechanics:** Client request (including payload) replication is separated from the leader's ordering task by using IP multicast for request dissemination to all nodes through a Switch. The leader is still responsible on establishing the total order of these received requests. For specific design details we refer the reader to Section 3.2 of [1].
- **In-Network Acceleration of Core SMR Messaging:** Hovercraft++ leverages in-network programmable hardware (e.g., P4 ASICs) to statelessly manage crucial SMR communication patterns. This includes:
  - The fan-out of leader messages (like AppendEntries containing ordering metadata) to followers.
  - The fan-in of follower replies related to these consensus messages.

The in-network aggregator offloading is viewed as a leader extension and it **aims** to improve the efficiency of the consensus protocol's internal communication, particularly as cluster size increases, without altering the fundamental Raft algorithm for achieving agreement. However, the Hovercraft++ protocol was not formally specified nor verified. For specific design details we refer the reader to Sections 4 (Figure 5.b explains communication in Hovercraft++), and 5 of [1]; subsection 6.4 describes important details about the NetAgg aggregator implementation of which we implement the AGG_COMMIT message to ensure Raft servers update their commit index to be able to respond to clients.

　Our TLA+ [1] specification of the Hovercraft++ protocol is based on its design presented in Sections 3, 4, and 5 [1]. We postpone load balancing client replies, load

---

[1] Lamport's TLA+ home page: https://lamport.azurewebsites.net/tla/tla.html

balancing read-only operations and bounded queues (sections 3.3 to 3.6 and partially section 5 of [1]) to future work discussing liveness and client guarantees. Our specification extends the original Raft TLA+ specification [2] to additionally model the Switch and NetAgg components introduced in Hovercraft++, as illustrated in Figure 1.

# 2 Hovercraft TLA+ Specification

---

[2]Raft TLA+ specification https://github.com/ongardie/raft.tla

—————————— MODULE *hovercraftpp* ——————————

In standard Raft the leader is the central hub. A client sends a request only to the leader. The leader must then replicate the entire request *payload* to all followers. The leader gathers acknowledgments from Followers, commits the entry, applies it, and sends a response to the client. Bottleneck: The leader's network bandwidth and processing power for sending the full *payload* to every follower becomes a limiting factor as the cluster size ($N$) or client request size increases. Throughput is limited by *Leader_Bandwidth* / ($(N − 1) ∗ Request\_Payload\_Size$).

We introduce a "Switch" abstraction (representing mechanisms like *IP* Multicast or a dedicated middlebox/programmable switch as used in the *HovercRaft* paper). Clients send requests via *Switch*. The *Switch* is non-crashing. Client requests timeouts due to *Switch* failures will force clients choose another *Switch*. The *Switch*'s responsibility is to deliver the request *payload* to all server nodes (Leader and Followers) "simultaneously". The 'simultaneous' delivery is a model abstraction for efficient broadcast mechanisms like *IP* Multicast, where *payload* dissemination is handled by the network infrastructure. We introduce the "*NetAgg*" network aggregation component of *HovercRaft* ++ . The *NetAgg* is non-crashing and stateless. Client requests will timeout and retry. The leader sends a single message containing the ordering metadata to *NetAgg*. Leader expects a *AggCommit* message in return. If one is not received in due time, the *Leader* declares *NetAgg* failure, another *NetAgg* will be chosen. *NetAgg* is then responsible for disseminating this metadata to followers and collecting their acknowledgments. When a follower receives the ordering metadata message from *NetAgg*, it uses the identifier in the metadata to find the corresponding *payload* in its temporary buffer. Once matched, the follower places the request *payload* into its replicated *log* at the correct index specified by the leader's metadata. Once *NetAgg* receives acknowledgments from a quorum of followers, it sends an *AggCommit* message to all servers, informing them that they can advance their commit index for that entry. *AggCommit* also ensures *NetAgg* failure handling by *Leader* choosing another *NetAgg*. This model assumes that the *Switch* and *NetAgg* components will never fail. We do not model *Switch* and *NetAgg* replicas due to model space constraints. Point-to-point recovery mechanisms between Followers and the *Leader* are partially addressed (assumes no *Follower* crashes).

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

************************* CONSTANTS *****************************************

The set of server *IDs* including the *Switch* and *NetAgg*
CONSTANTS *Server*

The set of client requests that can go into the *log*
Represents the possible values of client requests that are stored in the *log*.
A value with its term represent the request metadata. Same value is used as *payload*. A typical
flow for a request's data/metadata might be :
∗ Client → *Switch* (metadata and *payload*)
∗ *Switch* → *Leader* & Followers (metadata and *payload*)
∗ *Leader* → *NetAgg* (metadata for ordering)
∗ *NetAgg* → Followers → *NetAgg* (metadata for ordering and acknowledgements)
∗ *NetAgg* → All *Servers* (*AggCommit* sent when majority quorum of acks received)
∗ One server of *Servers* will reply back to client (not handled by this model).
CONSTANTS *Value*

Server states for Raft protocol.
CONSTANTS *Follower*, *Candidate*, *Leader*

Fixed states for *Switch* and *NetAgg* indices in *Server*
CONSTANTS *Switch*, *NetAgg*

A reserved value.
CONSTANTS *Nil*

Message types
Standard Raft message types for leader election and *log* replication.
CONSTANTS *RequestVoteRequest*, *RequestVoteResponse*,
            *AppendEntriesRequest*, *AppendEntriesResponse*

*Hovercraft* ++ specific message types for interactions involving the *NetAgg* component.
CONSTANTS *AppendEntriesNetAggRequest*, *AggCommit*

Limits the total number of client requests processed in the model,
used for bounding the state space during model checking.
CONSTANTS *MaxClientRequests*

Limits the number of times any given server can transition to
the *Leader* state, for state space bounding.
CONSTANTS *MaxBecomeLeader*

Defines the maximum value a term number can reach, used for state space bounding.
CONSTANTS *MaxTerm*

*********************** VARIABLES ****************************************

Global variables

A bag of records representing requests and responses sent from one server
to another. This is a function mapping Message to *Nat*.

2

VARIABLE *messages*

An instrumentation variable mapping each server *ID* (from *Server*)
to a natural number, counting how many times it has become a leader.
Used with *MaxBecomeLeader* for state space bounding.

VARIABLE *leaderCount*

An instrumentation variable; a natural number tracking the count of
client requests processed so far. Used with *MaxClientRequests* for
state space bounding.

VARIABLE *maxc*

variable for tracking entry commit message counts
Maps $\langle logIndex, logTerm \rangle$ to a record tracking message counts.
$[\ sentCount \mapsto Nat,\ \backslash * AppendEntriesRequests$ sent for the entry
$\quad ackCount \mapsto Nat,\ \backslash * AppendEntriesResponses$ received for this entry
$\quad committed \mapsto Bool\ ]\ \backslash *$ Flag indicating if the entry is committed

VARIABLE *entryCommitStats*

A tuple grouping all instrumentation-specific variables. Useful for specifying
UNCHANGED *instrumentationVars* in actions that do not modify them.
$instrumentationVars \triangleq \langle leaderCount, maxc, entryCommitStats \rangle$

The unique identifier (*ID* from *Server* set) of the server designated
to act as the *Switch* component.

VARIABLE *switchIndex*

The *Switch*'s internal buffer. Maps a request identifier ($\langle value, term \rangle$)
to the complete request data and *payload*

VARIABLE *switchBuffer*

A per-server variable (maps *Server ID* to a set of request identifiers).
For each server (Leader and Followers), it stores the set of request
identifiers ($\langle value, term \rangle$ tuples) for payloads received from the *Switch*
that are awaiting ordering metadata.

VARIABLE *unorderedRequests*

Records which $\langle value, term \rangle$ pairs the *Switch* has sent to each server.
Maps *Server ID* $\rightarrow$ Set of $\langle Value, Term \rangle$ pairs.

VARIABLE *switchSentRecord*

A tuple grouping variables specific to the *Hovercraft Switch* functionality.
$hovercraftVars \triangleq \langle switchBuffer, unorderedRequests,$
$\qquad\qquad\qquad switchIndex, switchSentRecord \rangle$

*NetAgg* variables

Stores the leader (leader field, type *Server*) and its term (term field, type *Nat*)

3

that *NetAgg* currently recognizes as active.
VARIABLE *netAggCurrentLeaderTerm*

| | |
|---|---|
| VARIABLE *netAggIndex* | Index of the *NetAgg* server |
| VARIABLE *netAggMatchIndex* | *NetAgg*'s view of follower match indices |
| VARIABLE *netAggPendingEntries* | Entries pending aggregation at *NetAgg* |
| VARIABLE *netAggCommitIndex* | *NetAgg*'s view of commit index |

A tuple grouping all variables specific to the *NetAgg* component.
$$netAggVars \triangleq \langle netAggIndex, netAggMatchIndex, netAggPendingEntries,$$
$$netAggCommitIndex, netAggCurrentLeaderTerm \rangle$$

The following variables are all per server (functions with domain *Server*).

Each server's current known term number. (Maps *Server ID* to *Nat*).
VARIABLE *currentTerm*

The server's state (Follower, *Candidate*, *Leader*, *Switch*, or *NetAgg*).
For *Switch* and *NetAgg* entities, this state is fixed.
VARIABLE *state*

The candidate the server voted for in its current term, or
Nil if it hasn't voted for any.
VARIABLE *votedFor*

$$serverVars \triangleq \langle currentTerm, state, votedFor \rangle$$

A Sequence of *log* entries. The index into this sequence is the index of the
*log* entry.
VARIABLE *log*

The index of the latest entry in the *log* the state machine may apply.
VARIABLE *commitIndex*

$$logVars \triangleq \langle log, commitIndex \rangle$$

The following variables are used only on candidates:
The set of servers from which the candidate has received a *RequestVote*
response in its *currentTerm*.
VARIABLE *votesResponded*

The set of servers from which the candidate has received a vote in its
*currentTerm*.
VARIABLE *votesGranted*

A history variable used in the proof. This would not be present in an
implementation.
Function from each server that voted for this candidate in its *currentTerm*
to that voter's *log*.

VARIABLE *voterLog*

$candidateVars \triangleq \langle votesResponded, votesGranted, voterLog \rangle$

The following variables are used only on leaders:
The next entry to send to each follower.
VARIABLE *nextIndex*

The latest entry that each follower has acknowledged is the same as the
leader's. This is used to calculate *commitIndex* on the leader.
VARIABLE *matchIndex*

$leaderVars \triangleq \langle nextIndex, matchIndex \rangle$

The set of server *IDs* participating in the Raft consensus
(*i.e.*, excluding *Switch* and *NetAgg* components).
$Servers \triangleq Server \setminus \{switchIndex, netAggIndex\} \setminus *$ see *Init*
VARIABLE *Servers*

All variables; used for stuttering (asserting state hasn't changed).
*Hovercraft* ++ brings *hovercraftVars* for *Switch* and *netAggVars* for *NetAgg*.
$vars \triangleq \langle messages, serverVars, candidateVars, leaderVars, logVars,$
$\qquad instrumentationVars, hovercraftVars, netAggVars, Servers \rangle$

*********************** HELPERS *********************************************

Defines the set of all possible quorums. A quorum is any subset of *Servers*
(Raft participants) forming a simple majority. The critical property is that
any two quorums must overlap.
$Quorum \triangleq \{i \in \text{SUBSET } (Servers) : Cardinality(i) * 2 > Cardinality(Servers)\}$

The term of the last entry in a *log*, or 0 if the *log* is empty.
$LastTerm(xlog) \triangleq \text{IF } Len(xlog) = 0 \text{ THEN } 0 \text{ ELSE } xlog[Len(xlog)].term$

$WithMessage(m, msgs) \triangleq$
$\quad \text{IF } m \in \text{DOMAIN } msgs \text{ THEN}$
$\qquad msgs$ avoiding duplicates
$\quad \text{ELSE}$
$\qquad msgs @@ (m :> 1)$

to allow duplicates use: $WithMessage(m, msgs) \triangleq$
$[msgs \text{ EXCEPT } ![m] = \text{IF } m \in \text{DOMAIN } msgs \text{ THEN } msgs[m] + 1 \text{ ELSE } 1]$

$WithoutMessage(m, msgs) \triangleq$
$\quad \text{IF } m \in \text{DOMAIN } msgs \text{ THEN}$
$\qquad [msgs \text{ EXCEPT } ![m] = \text{IF } msgs[m] > 0 \text{ THEN } msgs[m] - 1 \text{ ELSE } 0]$
$\quad \text{ELSE}$
$\qquad msgs$

Add a message to the bag of messages.
$Send(m) \triangleq messages' = WithMessage(m, messages)$

Remove a message from the bag of messages. Used when a server is done processing a message.
$Discard(m) \triangleq messages' = WithoutMessage(m, messages)$

Helper for *Send* and *Reply*. Given a message $m$ and bag of messages, return a Combination of *Send* and *Discard*
$Reply(response, request) \triangleq$
$\quad messages' = WithoutMessage(request, WithMessage(response, messages))$

Return the minimum value from a set, or undefined if the set is empty.
$Min(s) \triangleq \text{CHOOSE } x \in s : \forall y \in s : x \leq y$

Return the maximum value from a set, or undefined if the set is empty.
$Max(s) \triangleq \text{CHOOSE } x \in s : \forall y \in s : x \geq y$

Convert a sequence to a set of its elements
$SeqToSet(seq) \triangleq \{seq[i] : i \in \text{DOMAIN } seq\}$

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

$ValidMessage(msgs) \triangleq$
$\quad \{m \in \text{DOMAIN } messages : msgs[m] > 0\}$

The prefix of the *log* of server $i$ that has been committed up to term $x$
$CommittedTermPrefix(i, x) \triangleq$
Only if *log* of $i$ is non-empty, and if there exists an entry up to the term $x$
$\quad \text{IF } Len(log[i]) \neq 0 \land \exists y \in \text{DOMAIN } log[i] : log[i][y].term \leq x$
$\quad \quad \text{THEN}$
then, we use the subsequence up to the maximum committed term of the leader
$\quad \quad \quad \text{LET } maxTermIndex \triangleq$
$\quad \quad \quad \quad \text{CHOOSE } y \in \text{DOMAIN } log[i] :$
$\quad \quad \quad \quad \quad \land log[i][y].term \leq x$
$\quad \quad \quad \quad \quad \land \forall z \in \text{DOMAIN } log[i] : log[i][z].term \leq x \Rightarrow y \geq z$
$\quad \quad \quad \text{IN } \quad SubSeq(log[i], 1, min(maxTermIndex, commitIndex[i]))$
Otherwise the prefix is the empty tuple
$\quad \quad \quad \text{ELSE } \langle\rangle$

$CheckIsPrefix(seq1, seq2) \triangleq$
$\quad \land Len(seq1) \leq Len(seq2)$
$\quad \land \forall i \in 1 .. Len(seq1) : seq1[i] = seq2[i]$

The prefix of the *log* of server $i$ that has been committed
$Committed(i) \triangleq$
$\quad \text{IF } commitIndex[i] = 0$
$\quad \text{THEN } \langle\rangle$

ELSE $SubSeq(log[i], 1, commitIndex[i])$

$MyConstraint \triangleq (\forall\, i \in Servers : currentTerm[i] \leq MaxTerm$
$\qquad\qquad\qquad \wedge\, Len(log[i]) \leq MaxClientRequests)$
$\qquad\qquad\qquad \wedge\, (\forall\, m \in \text{DOMAIN } messages : messages[m] \leq 1)$

$InitHistoryVars \triangleq voterLog\ = [i \in Servers \mapsto [j \in \{\} \mapsto \langle\rangle]]$

$InitServerVars \triangleq\ \wedge currentTerm = [i \in Servers \mapsto 1]$
$\qquad\qquad\qquad\quad \wedge state \qquad\ = [i \in Servers \mapsto Follower]$
$\qquad\qquad\qquad\quad \wedge votedFor \quad\ = [i \in Servers \mapsto Nil]$

$InitCandidateVars \triangleq\ \wedge votesResponded = [i \in Servers \mapsto \{\}]$
$\qquad\qquad\qquad\qquad\ \wedge votesGranted\quad\ = [i \in Servers \mapsto \{\}]$

The values $nextIndex[i][i]$ and $matchIndex[i][i]$ are never read, since the
leader does not send itself messages. It's still easier to include these
in the functions.
$InitLeaderVars \triangleq\ \wedge nextIndex\ \ = [i \in Servers \mapsto [j \in Servers \mapsto 1]]$
$\qquad\qquad\qquad\quad \wedge matchIndex = [i \in Servers \mapsto [j \in Servers \mapsto 0]]$

$InitLogVars \triangleq\ \wedge log \qquad\qquad = [i \in Servers \mapsto \langle\rangle]$
$\qquad\qquad\qquad \wedge commitIndex\ = [i \in Servers \mapsto 0]$

$Init \triangleq$
$\quad \wedge messages = [m \in \{\} \mapsto 0]$
$\quad \wedge switchIndex\ = \text{CHOOSE } s \in Server : \text{TRUE}$   Pick any server as switch
$\quad \wedge netAggIndex = \text{CHOOSE } n \in Server \setminus \{switchIndex\} : \text{TRUE}$   Pick another as $NetAgg$
$\quad \wedge Servers = Server \setminus \{switchIndex,\ netAggIndex\}$   Remaining are Raft servers

$\quad$ Initialize all server state
$\quad \wedge currentTerm = [i \in Server \mapsto 1]$
$\quad \wedge state = [i \in Server \mapsto$
$\qquad\qquad\quad \text{IF } i = switchIndex \text{ THEN } Switch$
$\qquad\qquad\quad\ \text{ELSE } \text{IF } i = netAggIndex \text{ THEN } NetAgg$
$\qquad\qquad\quad\ \text{ELSE } Follower]$
$\quad \wedge votedFor = [i \in Server \mapsto Nil]$

$\quad$ Initialize empty logs and indices
$\quad \wedge log = [i \in Server \mapsto \langle\rangle]$
$\quad \wedge commitIndex = [i \in Server \mapsto 0]$
$\quad \wedge nextIndex = [i \in Server \mapsto [j \in Server \mapsto 1]]$
$\quad \wedge matchIndex = [i \in Server \mapsto [j \in Server \mapsto 0]]$

$\quad$ Initialize candidate variables
$\quad \wedge votesResponded = [i \in Server \mapsto \{\}]$

$\wedge\ votesGranted = [i \in Server \mapsto \{\}]$
$\wedge\ voterLog = [i \in Server \mapsto [j \in \{\} \mapsto \langle\rangle]]$

$\wedge\ switchBuffer = [vt \in \{\} \mapsto \{\}]$
$\wedge\ unorderedRequests = [s \in Server \mapsto \{\}]$
$\wedge\ switchSentRecord = [s \in Server \mapsto \{\}]$

$\wedge\ netAggCurrentLeaderTerm = Nil$
$\wedge\ netAggMatchIndex = [s \in \{\} \mapsto 0]$
$\wedge\ netAggPendingEntries = \{\}$
$\wedge\ netAggCommitIndex = 0$

$\wedge\ maxc = 0$
$\wedge\ leaderCount = [i\ \ \in Server \mapsto 0]$
$\wedge\ entryCommitStats = [idx\_term \in \{\} \mapsto$
$\qquad\qquad\qquad\qquad [sentCount \mapsto 0,\ ackCount \mapsto 0,\ committed \mapsto \text{FALSE}]]$

$MyInit \triangleq$
$\quad$ LET $ServerSet5 \triangleq$ CHOOSE $S \in$ SUBSET $(Server) : Cardinality(S) = 5$
$\qquad\ TheSwitchId \triangleq$ CHOOSE $s \in ServerSet5 :$ TRUE
$\qquad\ TempSet \triangleq ServerSet5 \setminus \{TheSwitchId\}$
$\qquad\ TheNetAggId \triangleq$ CHOOSE $n \in TempSet :$ TRUE
$\qquad\ TempSet2 \triangleq TempSet \setminus \{TheNetAggId\}$
$\qquad\ TheLeaderId \triangleq$ CHOOSE $l \in TempSet2 :$ TRUE
$\qquad\ FollowerIds \triangleq TempSet2 \setminus \{TheLeaderId\}$

$\qquad\ TheState \triangleq [s \in Server \mapsto$
$\qquad\qquad\qquad$ IF $s = TheSwitchId$ THEN $Switch$
$\qquad\qquad\qquad$ ELSE IF $s = TheNetAggId$ THEN $NetAgg$
$\qquad\qquad\qquad$ ELSE IF $s = TheLeaderId$ THEN $Leader$
$\qquad\qquad\qquad$ ELSE IF $s \in FollowerIds$ THEN $Follower$
$\qquad\qquad\qquad$ ELSE $Follower$
$\qquad\qquad\qquad ]$
$\qquad\ TheSwitchIndex \triangleq TheSwitchId$
$\qquad\ TheNetAggIndex \triangleq TheNetAggId$
$\qquad\ TheServersSet \triangleq Server \setminus \{TheSwitchIndex,\ TheNetAggIndex\}$
$\qquad\ Voters \triangleq TheServersSet \setminus \{TheLeaderId\}$
$\quad$ IN
$\quad \wedge\ Cardinality(Server) \geq 5$
$\quad \wedge\ PrintT(\text{“MyInit: switchIndex=”} \circ ToString(TheSwitchIndex))$

8

$\land$ *PrintT*("MyInit: netAggIndex=" $\circ$ *ToString*(*TheNetAggIndex*))
$\land$ *PrintT*("MyInit: Leader is=" $\circ$ *ToString*(*TheLeaderId*))
$\land$ *PrintT*("MyInit: Servers=" $\circ$ *ToString*(*TheServersSet*))
$\land$ *PrintT*("*MyInit* : *state*[*switchIndex*] = " $\circ$ *ToString*(*TheState*[*TheSwitchIndex*]))
$\land$ *PrintT*("*MyInit* : *state*[*LeaderId*] = " $\circ$ *ToString*(*TheState*[*TheLeaderId*]))
$\land$ *PrintT*("*MyInit* : *switchBuffer Domain* = " $\circ$ *ToString*(DOMAIN [$vt \in \{\} \mapsto \{\}$]))

$\land$ *commitIndex* = [$s \in Server \mapsto 0$]
$\land$ *currentTerm* = [$s \in Server \mapsto 2$]
$\land$ *leaderCount* = [$s \in Server \mapsto$ IF $s = TheLeaderId$ THEN 1 ELSE 0]
$\land$ *log* = [$s \in Server \mapsto \langle\rangle$]
$\land$ *matchIndex* = [$s \in Server \mapsto [t \in Server \mapsto 0]$]
$\land$ *maxc* = 0
$\land$ *messages* = [$m \in \{\} \mapsto 0$]
$\land$ *nextIndex* = [$s \in Server \mapsto [t \in Server \mapsto 1]$]
$\land$ *state* = *TheState*
$\land$ *votedFor* = [$s \in Server \mapsto$
     IF $s = TheLeaderId$ THEN *Nil* ELSE *TheLeaderId*]
$\land$ *voterLog* = [$s \in Server \mapsto$
     IF $s = TheLeaderId$ THEN
     [$v \in Voters \mapsto \langle\rangle$] ELSE [$v \in \{\} \mapsto \langle\rangle$]]]
$\land$ *votesGranted* = [$s \in Server \mapsto$
      IF $s = TheLeaderId$ THEN *Voters* ELSE $\{\}$]
$\land$ *votesResponded* = [$s \in Server \mapsto$
      IF $s = TheLeaderId$ THEN *Voters* ELSE $\{\}$]
$\land$ *entryCommitStats* = [$idx\_term \in \{\} \mapsto$
      [$sentCount \mapsto 0$,
      $ackCount \mapsto 0$,
      $committed \mapsto$ FALSE]]
$\land$ *switchBuffer* = [$vt \in \{\} \mapsto \{\}$]
$\land$ *unorderedRequests* = [$s \in Server \mapsto \{\}$]
$\land$ *switchSentRecord* = [$s \in Server \mapsto \{\}$]
$\land$ *switchIndex* = *TheSwitchIndex*
$\land$ *netAggIndex* = *TheNetAggIndex*
$\land$ *netAggMatchIndex* = [$s \in TheServersSet \mapsto 0$]
$\land$ *netAggPendingEntries* = $\{\}$
$\land$ *netAggCommitIndex* = 0
$\land$ *netAggCurrentLeaderTerm* = [$leader \mapsto TheLeaderId$, $term \mapsto 2$]
$\land$ *Servers* = *TheServersSet*

********************** Actions ******************************

Modified to limit Restarts only for Leaders.
Server $i$ restarts from stable storage.
It loses everything but its *currentTerm*, *votedFor*, and *log*.
Also persists messages, instrumentation and *Switch/NetAgg* variables.

9

$Restart(i) \triangleq$
 $\land\ state[i] = Leader$
 $\land\ (\forall\, srv \in Servers : leaderCount[srv] < MaxBecomeLeader)$
 $\land\ state' \qquad\qquad = [state\ \text{EXCEPT}\ ![i] = Follower]$
 $\land\ votesResponded' = [votesResponded\ \text{EXCEPT}\ ![i] = \{\}]$
 $\land\ votesGranted' \quad = [votesGranted\ \text{EXCEPT}\ ![i] = \{\}]$
 $\land\ voterLog' \qquad\quad = [voterLog\ \text{EXCEPT}\ ![i] = [j \in \{\} \mapsto \langle\rangle]]$
 $\land\ nextIndex' \qquad\ = [nextIndex\ \text{EXCEPT}\ ![i] = [j \in Server \mapsto 1]]$
 $\land\ matchIndex' \qquad = [matchIndex\ \text{EXCEPT}\ ![i] = [j \in Server \mapsto 0]]$
 $\land\ commitIndex' \quad = [commitIndex\ \text{EXCEPT}\ ![i] = 0]$
 $\land\ unorderedRequests' = [unorderedRequests\ \text{EXCEPT}\ ![i] = \{\}]$
 $\land\ switchSentRecord' = [switchSentRecord\ \text{EXCEPT}\ ![i] = \{\}]$
 $\land\ \text{IF}\ netAggCurrentLeaderTerm \neq Nil \land netAggCurrentLeaderTerm.leader = i$
   $\text{THEN}\ \ \land netAggCurrentLeaderTerm' = Nil\ \ $ Deactivate $NetAgg$
      $\land netAggPendingEntries' = \{\}\ \ \ \ $ Flush pending entries
     $netAggCommitIndex$ and $netAggMatchIndex$ could be left or reset;
     they become irrelevant until a new leader activates $NetAgg$.
   $\text{ELSE}\ \ \ \land \text{UNCHANGED}\ netAggCurrentLeaderTerm$
      $\land \text{UNCHANGED}\ netAggPendingEntries$
 $\land\ \text{UNCHANGED}\ \langle messages,\ currentTerm,\ votedFor,\ log,\ instrumentationVars,$
       $switchIndex,\ switchBuffer,\ Servers,$
       $netAggIndex,\ netAggMatchIndex,\ netAggCommitIndex\rangle$

Server $i$ times out and starts a new election. *Follower → Candidate*
$Timeout(i) \triangleq\ \ \land state[i] \in \{Follower,\ Candidate\}$
     $\land (\forall\, srv \in Servers : leaderCount[srv] < MaxBecomeLeader)$
     $\land currentTerm[i] < MaxTerm$
     $\land state' = [state\ \text{EXCEPT}\ ![i] = Candidate]$
     $\land currentTerm' = [currentTerm\ \text{EXCEPT}\ ![i] = currentTerm[i] + 1]$
     Most implementations would probably just set the local vote
     atomically, but messaging localhost for it is weaker.
     $\land votedFor' = [votedFor\ \text{EXCEPT}\ ![i] = Nil]$
     $\land votesResponded' = [votesResponded\ \text{EXCEPT}\ ![i] = \{\}]$
     $\land votesGranted' \quad = [votesGranted\ \text{EXCEPT}\ ![i] = \{\}]$
     $\land voterLog' \qquad\quad = [voterLog\ \text{EXCEPT}\ ![i] = [j \in \{\} \mapsto \langle\rangle]]$
     $\land \text{UNCHANGED}\ \langle messages,\ leaderVars,\ logVars,$
         $instrumentationVars,\ hovercraftVars,$
         $Servers,\ netAggVars\rangle$

Modified to restrict *Leader* transitions, bounded by *MaxBecomeLeader*
Candidate $i$ transitions to leader. *Candidate → Leader*
$BecomeLeader(i) \triangleq$
 $\land state[i] = Candidate$
 $\land votesGranted[i] \in Quorum$
 $\land leaderCount[i] < MaxBecomeLeader$

$\land\ state'\qquad = [state\ \text{EXCEPT}\ ![i] = Leader]$

$\land\ nextIndex'\quad = [nextIndex\ \text{EXCEPT}\ ![i] =$
$\qquad\qquad\qquad\qquad [j \in Server \mapsto Len(log[i]) + 1]]$

$\land\ matchIndex' = [matchIndex\ \text{EXCEPT}\ ![i] =$
$\qquad\qquad\qquad\qquad [j \in Server \mapsto 0]]$

$\land\ leaderCount' = [leaderCount\ \text{EXCEPT}\ ![i] = leaderCount[i] + 1]$

$\land\ netAggCurrentLeaderTerm' = [leader \mapsto i,\ term \mapsto currentTerm[i]]$

$\land\ netAggPendingEntries' = \{\}$   Flush pending entries for the new leader

$\land\ netAggCommitIndex' = commitIndex[i]$

$\land\ netAggMatchIndex' = [s \in Servers \mapsto 0]$

$\land\ \text{UNCHANGED}\ \langle messages,\ currentTerm,\ votedFor,\ candidateVars,$
$\qquad\qquad\qquad\ logVars,\ maxc,\ entryCommitStats,\ hovercraftVars,$
$\qquad\qquad\qquad\ Servers,\ netAggIndex\rangle$

Modified up to *MaxTerm*; Back To *Follower*.

Any *RPC* with a newer term causes the recipient to advance its term first.

$UpdateTerm(i,\ j,\ m)\ \triangleq$

$\qquad \land\ state[i] \notin \{Switch,\ NetAgg\} \land state[j] \notin \{Switch,\ NetAgg\}$

$\qquad \land\ m.mterm > currentTerm[i]$

$\qquad \land\ m.mterm < MaxTerm$

$\qquad \land\ \text{LET}\ wasLeader\ \triangleq\ state[i] = Leader$
$\qquad\quad \text{IN}$
$\qquad\quad \land\ currentTerm'\qquad = [currentTerm\ \text{EXCEPT}\ ![i] = m.mterm]$
$\qquad\quad \land\ state'\qquad\qquad\ = [state\qquad\ \text{EXCEPT}\ ![i]\quad = Follower]$
$\qquad\quad \land\ votedFor'\qquad\ = [votedFor\quad\ \text{EXCEPT}\ ![i]\quad = Nil]$
$\qquad\quad \land\ \text{IF}\ wasLeader \land netAggCurrentLeaderTerm \neq Nil \land netAggCurrentLeaderTerm.leader = i$
$\qquad\qquad \text{THEN}\ \land\ netAggCurrentLeaderTerm' = Nil$  Deactivate *NetAgg*
$\qquad\qquad\qquad\quad \land\ netAggPendingEntries' = \{\}$   Flush pending entries
$\qquad\qquad \text{ELSE}\quad \land\ \text{UNCHANGED}\ netAggCurrentLeaderTerm$
$\qquad\qquad\qquad\quad \land\ \text{UNCHANGED}\ netAggPendingEntries$
$\qquad\quad$ messages is unchanged so $m$ can be processed further.
$\qquad \land\ \text{UNCHANGED}\ \langle messages,\ candidateVars,\ leaderVars,\ logVars,$
$\qquad\qquad\qquad\qquad instrumentationVars,\ hovercraftVars,\ Servers,$
$\qquad\qquad\qquad\qquad netAggIndex,\ netAggMatchIndex,\ netAggCommitIndex\rangle$

*************************** REQUEST VOTE *********************************

Message handlers

$i = $ recipient, $j = $ sender, $m = $ message

Candidate $i$ sends $j$ a *RequestVote* request.

$RequestVote(i,\ j)\ \triangleq$

$\qquad \land\ state[i] = Candidate$

$\qquad \land\ state[j] \notin \{Switch,\ NetAgg\}$

$\qquad \land\ j \notin votesResponded[i]$

$\qquad \land\ Send([mtype\qquad\qquad \mapsto RequestVoteRequest,$

$$
\begin{aligned}
&mterm && \mapsto currentTerm[i], \\
&mlastLogTerm && \mapsto LastTerm(log[i]), \\
&mlastLogIndex && \mapsto Len(log[i]), \\
&msource && \mapsto i, \\
&mdest && \mapsto j]) \\
\end{aligned}
$$

$\land$ UNCHANGED $\langle serverVars,\ candidateVars,\ leaderVars,\ logVars,$
$\qquad\qquad\qquad instrumentationVars,\ hovercraftVars,\ Servers,\ netAggVars \rangle$

$HandleRequestVoteRequest(i,\ j,\ m)\ \triangleq$
$\quad$ LET $logOk\ \triangleq\ \lor\ m.mlastLogTerm > LastTerm(log[i])$
$\qquad\qquad\qquad\quad\ \lor\ \land\ m.mlastLogTerm = LastTerm(log[i])$
$\qquad\qquad\qquad\qquad\quad\ \land\ m.mlastLogIndex \geq Len(log[i])$
$\qquad\quad\ grant\ \triangleq\ \land\ m.mterm = currentTerm[i]$
$\qquad\qquad\qquad\qquad\ \land\ logOk$
$\qquad\qquad\qquad\qquad\ \land\ votedFor[i] \in \{Nil,\ j\}$
$\quad$ IN $\quad \land\ m.mterm \leq currentTerm[i]$
$\qquad\qquad \land\ \lor\ grant\quad \land\ votedFor' = [votedFor$ EXCEPT $![i] = j]$
$\qquad\qquad\quad\ \lor\ \neg grant \land$ UNCHANGED $votedFor$
$\qquad\qquad \land\ Reply([mtype \qquad\qquad \mapsto RequestVoteResponse,$
$\qquad\qquad\qquad\quad mterm \qquad\qquad \mapsto currentTerm[i],$
$\qquad\qquad\qquad\quad mvoteGranted \mapsto grant,$

$\qquad\qquad\qquad\quad mlog \qquad\qquad\quad \mapsto log[i],$
$\qquad\qquad\qquad\quad msource \qquad\quad \mapsto i,$
$\qquad\qquad\qquad\quad mdest \qquad\qquad \mapsto j],$
$\qquad\qquad\qquad\quad m)$
$\qquad\qquad \land$ UNCHANGED $\langle state,\ currentTerm,\ candidateVars,\ leaderVars,\ logVars,$
$\qquad\qquad\qquad\qquad\qquad\quad instrumentationVars,\ hovercraftVars,\ Servers,\ netAggVars \rangle$

$HandleRequestVoteResponse(i,\ j,\ m)\ \triangleq$
$\quad \land\ m.mterm = currentTerm[i]$
$\quad \land\ votesResponded' = [votesResponded$ EXCEPT $![i] =$
$\qquad\qquad\qquad\qquad\qquad\qquad votesResponded[i] \cup \{j\}]$
$\quad \land\ \lor\ \land\ m.mvoteGranted$
$\qquad\quad \land\ votesGranted' = [votesGranted$ EXCEPT $![i] =$
$\qquad\qquad\qquad\qquad\qquad\qquad votesGranted[i] \cup \{j\}]$
$\qquad\quad \land\ voterLog' = [voterLog$ EXCEPT $![i] =$
$\qquad\qquad\qquad\qquad\qquad voterLog[i] @@ (j :> m.mlog)]$

$\lor\ \land \neg m.mvoteGranted$
$\qquad \land \text{UNCHANGED}\ \langle votesGranted,\ voterLog \rangle$
$\land\ Discard(m)$
$\land\ \text{UNCHANGED}\ \langle serverVars,\ votedFor,\ leaderVars,\ logVars,$
$\qquad\qquad\qquad instrumentationVars,\ hovercraftVars,\ Servers,\ netAggVars \rangle$

Responses with stale terms are ignored.
$DropStaleResponse(i,\ j,\ m)\ \triangleq$
$\quad \land\ m.mterm < currentTerm[i]$
$\quad \land\ Discard(m)$
$\quad \land\ \text{UNCHANGED}\ \langle serverVars,\ candidateVars,\ leaderVars,\ logVars,$
$\qquad\qquad\qquad instrumentationVars,\ hovercraftVars,\ Servers,\ netAggVars \rangle$

*************************** $AppendEntries$ ********************************

Leader $i$ ingests a request $v$ that has been replicated to its unordered set.
$LeaderIngestHovercRaftRequest(i,\ vt)\ \triangleq$
$\quad \land\ state[i] = Leader$
$\quad \land\ vt \in unorderedRequests[i]$      Request *ID* is pending for the leader
$\quad \land\ vt \in \text{DOMAIN}\ switchBuffer$      use switch buffer to reduce *payload* duplication
$\quad \land\ maxc < MaxClientRequests$
$\quad \land\ \text{LET}\ entryFromBuffer\ \triangleq\ switchBuffer[vt]$
$\qquad\qquad v\ \triangleq\ vt[1]$    Extract value from $\langle value,\ term \rangle$ pair
$\qquad\qquad$ Use leader's current term, keep value and *payload* from buffer
$\qquad\qquad newEntry\ \triangleq\ [term \mapsto currentTerm[i],$
$\qquad\qquad\qquad\qquad\qquad value \mapsto v,$
$\qquad\qquad\qquad\qquad\qquad payload \mapsto entryFromBuffer.payload]$
$\qquad\qquad entryExists\ \triangleq\ \exists\, k \in \text{DOMAIN}\ log[i]:$
$\qquad\qquad\qquad\qquad log[i][k].value = v \land log[i][k].term = newEntry.term$
$\qquad\qquad newLog\ \triangleq\ \text{IF}\ entryExists\ \text{THEN}\ log[i]\ \text{ELSE}\ Append(log[i],\ newEntry)$
$\qquad\qquad newEntryIndex\ \triangleq\ Len(log[i]) + 1$
$\qquad\qquad newEntryKey\ \triangleq\ \langle newEntryIndex,\ newEntry.term \rangle$
$\quad\ \text{IN}$
$\qquad \land\ log' = [log\ \text{EXCEPT}\ ![i] = newLog]$
$\qquad \land\ maxc' = \text{IF}\ entryExists\ \text{THEN}\ maxc\ \text{ELSE}\ maxc + 1$
$\qquad \land\ entryCommitStats' =$
$\qquad\qquad \text{IF}\ \neg entryExists \land newEntryIndex > 0$
$\qquad\qquad\quad \text{THEN}\ entryCommitStats\ @@\ (newEntryKey :> [sentCount \mapsto 0,$
$\qquad\qquad\qquad\qquad\qquad\qquad ackCount \mapsto 0,\ committed \mapsto \text{FALSE}])$
$\qquad\qquad\quad \text{ELSE}\ entryCommitStats$
$\qquad \land\ unorderedRequests' = [unorderedRequests\ \text{EXCEPT}\ ![i] = @ \setminus \{vt\}]$
$\quad \land\ \text{UNCHANGED}\ \langle messages,\ serverVars,\ candidateVars,\ leaderVars,$
$\qquad\qquad\qquad commitIndex,\ leaderCount,\ switchIndex,\ switchBuffer,$
$\qquad\qquad\qquad Servers,\ switchSentRecord,\ netAggVars \rangle$

Client sends a request to the *Switch*, which buffers it,
not yet replicated to servers
$s$ is the *Switch*, $i$ is the leader, $v$ is the request value
(along with term will represent a request *ID* in this model)
$SwitchClientRequest(s, i, v) \triangleq$
 $\wedge\ state[s] = Switch$ Only the switch server can process client requests
 $\wedge\ state[i] = Leader$
 $\wedge$ LET $vt \triangleq \langle v, currentTerm[i]\rangle$ Create $\langle value, term\rangle$ pair
  IN
   $\wedge\ vt \notin$ DOMAIN $switchBuffer$ Only process new requests
   $\wedge$ LET $entryWithPayload \triangleq [term \mapsto currentTerm[i],$
                $value \mapsto v, payload \mapsto v]$
    IN
     $\wedge\ switchBuffer' = switchBuffer @@ (vt :> entryWithPayload)$
     $\wedge\ unorderedRequests' =$
      $[unorderedRequests$ EXCEPT $![s] = unorderedRequests[s] \cup \{vt\}]$
 $\wedge$ UNCHANGED $\langle messages, serverVars, candidateVars, leaderVars, logVars,$
        $leaderCount, entryCommitStats, switchIndex, maxc,$
        $Servers, switchSentRecord, netAggVars\rangle$

check that *Server* $i$ is not a *Switch* or *NetAgg*
$RaftState(i) \triangleq state[i] \notin \{Switch, NetAgg\}$

The *Switch* replicates $vt$ to ALL servers at once (except those that already have it).
This reduces state space by avoiding intermediate states
where only some servers have received the request.
$SwitchClientRequestReplicateAll(s, vt) \triangleq$
 $\wedge\ state[s] = Switch$ Only the switch server can replicate requests
 $\wedge\ vt \in unorderedRequests[s]$ Request must be pending at the switch
 $\wedge$ LET Find all servers that haven't received this $v$/term pair yet

   $targetServers \triangleq \{i \in Server : RaftState(i) \wedge vt \notin switchSentRecord[i]\}$
  IN
  $\wedge\ targetServers \neq \{\}$ At least one server needs the request
  $\wedge\ unorderedRequests' = [i \in Server \mapsto$
   IF $i \in targetServers$
    THEN $unorderedRequests[i] \cup \{vt\}$
    ELSE $unorderedRequests[i]]$
  $\wedge\ switchSentRecord' = [i \in Server \mapsto$
   IF $i \in targetServers$
    THEN $switchSentRecord[i] \cup \{vt\}$
    ELSE $switchSentRecord[i]]$
 $\wedge$ UNCHANGED $\langle messages, serverVars, candidateVars, leaderVars, logVars,$
        $leaderCount, entryCommitStats, switchIndex, switchBuffer,$
        $maxc, Servers, netAggVars\rangle$

Follower $i$ drops/loses one request $vt$ from its unordered requests

This simulates network loss or follower crash scenarios

$FollowerDropRequest(i, vt) \triangleq$

 $\wedge\ state[i] = Follower$ Only followers can drop requests

 $\wedge\ vt \in unorderedRequests[i]$ Request must exist in follower's buffer

 $\wedge\ vt \in \text{DOMAIN } switchBuffer$ Request must still exist in switch buffer

 $\wedge\ unorderedRequests' = [unorderedRequests \text{ EXCEPT } ![i] = @ \setminus \{vt\}]$

 $\wedge\ \text{UNCHANGED } \langle messages,\ serverVars,\ candidateVars,\ leaderVars,\ logVars,$

        $instrumentationVars,\ switchIndex,\ switchBuffer,$

        $switchSentRecord,\ Servers,\ netAggVars\rangle$

Leader $i$ sends $AppendEntries$ to $NetAgg$ instead of directly to followers

$AppendEntriesToNetAgg(i) \triangleq$

 $\wedge\ state[i] = Leader$

 $\wedge\ state[netAggIndex] = NetAgg$

 $\wedge\ Len(log[i]) > 0$

 $\wedge\ \text{LET } nextIndexMin \triangleq Min(\{nextIndex[i][j] : j \in Servers \setminus \{i\}\})$

  $\text{IN}\quad nextIndexMin \leq Len(log[i])$

 $\wedge\ \text{LET } entryIndex \triangleq Min(\{nextIndex[i][j] : j \in Servers \setminus \{i\}\})$

    $entry \triangleq log[i][entryIndex]$

    $entryMetadata \triangleq [term \mapsto entry.term,\ value \mapsto entry.value]$

    $entries \triangleq \langle entryMetadata\rangle$

    $prevLogIndex \triangleq entryIndex - 1$

    $prevLogTerm \triangleq \text{IF } prevLogIndex > 0 \text{ THEN}$

             $log[i][prevLogIndex].term$

           $\text{ELSE } 0$

  $\text{IN}\quad Send([mtype \qquad\qquad\ \ \mapsto AppendEntriesNetAggRequest,$

      $mterm \qquad\qquad\ \ \ \mapsto currentTerm[i],$

      $mprevLogIndex \ \ \mapsto prevLogIndex,$

      $mprevLogTerm \ \ \mapsto prevLogTerm,$

      $mentries \qquad\qquad \mapsto entries,$

      $mentryIndex \qquad \mapsto entryIndex,$

      $mlog \qquad\qquad\quad\ \mapsto log[i],$

      $mcommitIndex \ \ \ \mapsto Min(\{commitIndex[i],\ entryIndex\}),$

      $msource \qquad\qquad \mapsto i,$

      $mdest \qquad\qquad\ \ \mapsto netAggIndex])$

 $\wedge\ \text{UNCHANGED } \langle serverVars,\ candidateVars,\ leaderVars,\ logVars,$

        $instrumentationVars,\ hovercraftVars,\ netAggVars,\ Servers,\ netAggVars\rangle$

$NetAgg$ receives $AppendEntries$ from leader and forwards to ALL followers atomically

$NetAggForwardAppendEntriesAll(m) \triangleq$

 $\wedge\ m.mdest = netAggIndex$

 $\wedge\ m.mtype = AppendEntriesNetAggRequest$

 $\wedge\ netAggCurrentLeaderTerm \neq Nil$    $NetAgg$ must be active

 $\wedge\ netAggCurrentLeaderTerm.leader = m.msource$  Request from current assigned leader

15

$\wedge\ netAggCurrentLeaderTerm.term = m.mterm$

$\wedge\ \text{LET } leaderId\ \triangleq\ m.msource$

$\qquad\quad followers\ \triangleq\ Servers \setminus \{leaderId\}$   All servers except the leader

$\qquad\quad$ Create the set of messages to send to all followers

$\qquad\quad followerMessages\ \triangleq\ \{[mtype \qquad\qquad \mapsto AppendEntriesRequest,$

$\qquad\qquad\qquad\qquad\qquad\qquad\ mterm \qquad\qquad \mapsto m.mterm,$

$\qquad\qquad\qquad\qquad\qquad\qquad\ mprevLogIndex\ \mapsto m.mprevLogIndex,$

$\qquad\qquad\qquad\qquad\qquad\qquad\ mprevLogTerm\ \ \mapsto m.mprevLogTerm,$

$\qquad\qquad\qquad\qquad\qquad\qquad\ mentries \qquad\qquad \mapsto m.mentries,$

$\qquad\qquad\qquad\qquad\qquad\qquad\ mlog \qquad\qquad\qquad \mapsto m.mlog,$

$\qquad\qquad\qquad\qquad\qquad\qquad\ mcommitIndex\ \ \ \mapsto m.mcommitIndex,$

$\qquad\qquad\qquad\qquad\qquad\qquad\ msource \qquad\qquad\ \mapsto netAggIndex,$

$\qquad\qquad\qquad\qquad\qquad\qquad\ mdest \qquad\qquad\quad \mapsto f,$

$\qquad\qquad\qquad\qquad\qquad\qquad\ moriginalLeader \mapsto leaderId]$

$\qquad\qquad\qquad\qquad\qquad\qquad\ : f \in followers\}$

$\qquad\quad$ Remove the processed message and add all new messages

$\qquad\quad RemainingActiveMessages\ \triangleq\ ValidMessage(WithoutMessage(m,\ messages))$

$\qquad\quad$ Update $sentCount$ for this entry

$\qquad\quad entryIndex\ \triangleq\ m.mentryIndex$

$\qquad\quad entryTerm\ \triangleq\ m.mentries[1].term$

$\qquad\quad entryKey\ \triangleq\ \langle entryIndex,\ entryTerm \rangle$

$\quad\text{IN}$

$\quad \wedge\ messages' = [msgRec \in RemainingActiveMessages \cup followerMessages \mapsto 1]$

$\quad \wedge\ netAggPendingEntries' = netAggPendingEntries\ \cup$

$\qquad\quad \{[entryIndex \mapsto m.mentryIndex,$

$\qquad\qquad entryTerm\ \ \mapsto m.mentries[1].term,$

$\qquad\qquad leaderId\ \quad \mapsto leaderId,$

$\qquad\qquad ackCount\ \ \mapsto 0,$

$\qquad\qquad acksFrom\ \ \mapsto \{\}]\}$

$\quad \wedge\ entryCommitStats' =$

$\qquad\quad \text{IF } entryKey \in \text{DOMAIN } entryCommitStats \wedge \neg entryCommitStats[entryKey].committed$

$\qquad\quad \text{THEN } [entryCommitStats \text{ EXCEPT } ![entryKey].sentCount = @ + Cardinality(followers)]$

$\qquad\quad \text{ELSE } entryCommitStats$

$\quad \wedge\ \text{UNCHANGED } \langle serverVars,\ candidateVars,\ leaderVars,\ logVars,$

$\qquad\qquad\qquad\qquad\quad leaderCount,\ maxc,\ hovercraftVars,$

$\qquad\qquad\qquad\qquad\quad netAggIndex,\ netAggMatchIndex,\ netAggCommitIndex,$

$\qquad\qquad\qquad\qquad\quad netAggCurrentLeaderTerm,\ Servers \rangle$

$NetAgg$ receives $AppendEntries$ response from follower

$NetAggHandleAppendEntriesResponse(m)\ \triangleq$

$\quad \wedge\ m.mtype = AppendEntriesResponse$

$\quad \wedge\ m.mdest = netAggIndex$   Message is for $NetAgg$

$\quad \wedge\ netAggCurrentLeaderTerm \neq Nil$   $NetAgg$ must be active

$\quad \wedge\ m.msuccess$   Process successful $ACKs$, $NACKs$ go to $Leader$ for point to point recovery

$\quad \wedge\ \exists\, pending \in netAggPendingEntries :$

$\wedge\ m.msource \in (Servers \setminus \{pending.leaderId\})$

$\wedge\ m.mmatchIndex \geq pending.entryIndex$

$\wedge\ m.msource \notin pending.acksFrom$

$\wedge\ (\text{IF } m.msource \notin \text{DOMAIN } netAggMatchIndex$

$\quad \text{THEN } PrintT(\text{“DEBUG: m.msource NOT IN DOMAIN netAggMatchIndex”}$

$\qquad\qquad\qquad \circ \text{“\textbackslash n m.msource = ”} \circ ToString(m.msource)$

$\qquad\qquad\qquad \circ \text{“\textbackslash n m = ”} \circ ToString(m)$

$\qquad\qquad\qquad \circ \text{“\textbackslash n pending = ”} \circ ToString(pending)$

$\qquad\qquad\qquad \circ \text{“\textbackslash n netAggMatchIndex = ”} \circ ToString(netAggMatchIndex)$

$\qquad\qquad\qquad \circ \text{“\textbackslash n DOMAIN netAggMatchIndex = ”} \circ ToString(\text{DOMAIN } netAggMatchIndex)$

$\qquad\qquad\qquad \circ \text{“\textbackslash n Servers = ”} \circ ToString(Servers)$

$\qquad\qquad\qquad \circ \text{“\textbackslash n netAggCurrentLeaderTerm = ”} \circ ToString(netAggCurrentLeaderTerm)$

$\qquad\qquad\qquad \circ \text{“\textbackslash n netAggPendingEntries = ”} \circ ToString(netAggPendingEntries)$

$\qquad\qquad\qquad \circ \text{“\textbackslash n currentTerm = ”} \circ ToString(currentTerm)$

$\qquad\qquad\qquad \circ \text{“\textbackslash n state = ”} \circ ToString(state)$

$\qquad\qquad\qquad \circ \text{“\textbackslash n log length for m.msource = ”} \circ ToString(Len(log[m.msource]))$

$\qquad\qquad\qquad \circ \text{“\textbackslash n commitIndex for m.msource = ”} \circ ToString(commitIndex[m.msource])$

$\qquad\qquad\qquad \circ \text{“\textbackslash n All messages = ”} \circ ToString(messages)$

$\qquad\qquad )$

$\quad \text{ELSE TRUE})$

$\wedge\ \text{LET } updatedPending \triangleq [pending \text{ EXCEPT } !.acksFrom = @ \cup \{m.msource\}]$

$\qquad\quad RequiredFollowerAcks \triangleq Cardinality(Servers) \div 2$

$\quad \text{IN}$

$\quad \wedge\ m.msource \in \text{DOMAIN } netAggMatchIndex$

$\quad \wedge\ netAggMatchIndex' = [netAggMatchIndex \text{ EXCEPT } ![m.msource] = m.mmatchIndex]$

$\quad \wedge\ \text{IF } Cardinality(updatedPending.acksFrom) \geq RequiredFollowerAcks$

$\qquad \text{THEN}$

$\qquad \text{LET } AggCommitMsgsSet \triangleq \{[mtype \qquad\qquad \mapsto AggCommit,$

$\qquad\qquad\qquad\qquad\qquad\qquad\quad mcommitIndex \mapsto pending.entryIndex,$

$\qquad\qquad\qquad\qquad\qquad\qquad\quad msource \qquad\qquad \mapsto netAggIndex,$

$\qquad\qquad\qquad\qquad\qquad\qquad\quad mdest \qquad\qquad\ \mapsto srv,$

$\qquad\qquad\qquad\qquad\qquad\qquad\quad mterm \qquad\qquad\ \mapsto pending.entryTerm]$

$\qquad\qquad\qquad\qquad\qquad\qquad : srv \in Servers\}$

$\qquad\qquad RemainingActiveMessages \triangleq ValidMessage(WithoutMessage(m, messages))$

$\qquad \text{IN}$

$\qquad\quad \wedge\ messages' = [msgRec \in RemainingActiveMessages \cup AggCommitMsgsSet \mapsto 1]$

- All messages in the resulting bag have count 1, respecting *MyConstraint*.

$\quad\quad\quad \land \mathit{netAggPendingEntries}' = \mathit{netAggPendingEntries} \setminus \{\mathit{pending}\}$

$\quad\quad\quad$ Remove committed entry from pending

$\quad\quad\quad \land \mathit{netAggCommitIndex}' = \mathit{Max}(\{\mathit{netAggCommitIndex}, \mathit{pending.entryIndex}\})$

$\quad\quad\quad \land$ UNCHANGED $\langle \mathit{serverVars}, \mathit{candidateVars}, \mathit{leaderVars}, \mathit{logVars},$
$\quad\quad\quad\quad\quad\quad \mathit{instrumentationVars}, \mathit{hovercraftVars}, \mathit{netAggIndex},$
$\quad\quad\quad\quad\quad\quad \mathit{netAggCurrentLeaderTerm}, \mathit{Servers}\rangle$

$\quad\quad$ ELSE $\quad$ Majority not yet reached

$\quad\quad\quad \land \mathit{Discard}(m)$ This defines $\mathit{messages}' = \mathit{WithoutMessage}(m, \mathit{messages})$

$\quad\quad\quad \land \mathit{netAggPendingEntries}' = (\mathit{netAggPendingEntries} \setminus \{\mathit{pending}\}) \cup \{\mathit{updatedPending}\}$

$\quad\quad\quad \land$ UNCHANGED $\mathit{netAggCommitIndex}$

$\quad\quad\quad \land$ UNCHANGED $\langle \mathit{serverVars}, \mathit{candidateVars}, \mathit{leaderVars}, \mathit{logVars},$
$\quad\quad\quad\quad\quad\quad \mathit{instrumentationVars}, \mathit{hovercraftVars}, \mathit{netAggIndex}, \mathit{Servers},$
$\quad\quad\quad\quad\quad\quad \mathit{netAggCurrentLeaderTerm}\rangle$

Server receives $\mathit{AGG\_COMMIT}$ from $\mathit{NetAgg}$

$\mathit{HandleAggCommit}(i, m) \triangleq$

$\quad \land m.mtype = \mathit{AggCommit}$

$\quad \land m.mdest = i$

$\quad \land \mathit{state}[i] \in \{\mathit{Leader}, \mathit{Follower}\}$

$\quad \land ((\mathit{state}[i] = \mathit{Leader} \quad \land m.mterm = \mathit{currentTerm}[i]) \lor$

$\quad\quad$ Leader: entry's term must match leader's current term

$\quad\quad\quad (\mathit{state}[i] = \mathit{Follower} \quad \land m.mterm \leq \mathit{currentTerm}[i]))$

$\quad\quad\quad\quad$ Follower: entry's term can be current or older (but not newer)

$\quad \land$ LET $\mathit{receivedCommitIndex} \triangleq m.mcommitIndex \quad$ Added for clarity

$\quad\quad\quad \mathit{currentLogLen} \triangleq \mathit{Len}(\mathit{log}[i]) \quad\quad$ Added: get current *log* length

$\quad\quad\quad \mathit{newAdvancedCommitIndex} \triangleq \mathit{Max}(\{\mathit{commitIndex}[i], \mathit{receivedCommitIndex}\})$

$\quad\quad\quad\quad$ Renamed & Logic: advance if *m.mcommitIndex* is higher

$\quad\quad\quad \mathit{newCommitIndex} \triangleq \mathit{Min}(\{\mathit{newAdvancedCommitIndex}, \mathit{currentLogLen}\})$

$\quad\quad\quad\quad$ Modified: cap at current *log* length

$\quad\quad\quad \mathit{committedIndexes} \triangleq \{k \in \mathit{Nat} : \land k > \mathit{commitIndex}[i]$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \land k \leq \mathit{newCommitIndex}\}$

$\quad\quad\quad \mathit{keysToUpdate} \triangleq$ IF $\mathit{state}[i] = \mathit{Leader}$
$\quad\quad\quad\quad\quad\quad$ THEN $\{\mathit{key} \in$ DOMAIN $\mathit{entryCommitStats} :$
$\quad\quad\quad\quad\quad\quad\quad\quad \mathit{key}[1] \in \mathit{committedIndexes}\}$
$\quad\quad\quad\quad\quad\quad$ ELSE $\{\}$

$\quad\quad$ IN

$\quad\quad \land \mathit{commitIndex}' = [\mathit{commitIndex}$ EXCEPT $![i] = \mathit{newCommitIndex}]$

$\quad\quad \land \mathit{entryCommitStats}' =$ IF $\mathit{state}[i] = \mathit{Leader}$
$\quad\quad\quad\quad\quad\quad$ THEN $[\mathit{key} \in$ DOMAIN $\mathit{entryCommitStats} \mapsto$
$\quad\quad\quad\quad\quad\quad\quad$ IF $\mathit{key} \in \mathit{keysToUpdate}$
$\quad\quad\quad\quad\quad\quad\quad$ THEN $[\mathit{entryCommitStats}[\mathit{key}]$ EXCEPT $!.\mathit{committed} =$ TRUE$]$
$\quad\quad\quad\quad\quad\quad\quad$ ELSE $\mathit{entryCommitStats}[\mathit{key}]]$
$\quad\quad\quad\quad\quad\quad$ ELSE $\mathit{entryCommitStats}$

18

$\wedge$ IF $state[i] = Leader$
$\quad$ THEN $\wedge nextIndex' = [nextIndex$ EXCEPT $![i] =$
$\qquad\qquad\qquad\qquad [j \in Server \mapsto$
$\qquad\qquad\qquad\qquad\quad$ IF $j \in Servers \setminus \{i\}$
$\qquad\qquad\qquad\qquad\qquad$ THEN $Max(\{nextIndex[i][j],\ newCommitIndex + 1\})$
$\qquad\qquad\qquad\qquad\qquad$ ELSE $nextIndex[i][j]]]$
$\qquad\quad \wedge matchIndex' = [matchIndex$ EXCEPT $![i] =$
$\qquad\qquad\qquad\qquad [j \in Server \mapsto$
$\qquad\qquad\qquad\qquad\quad$ IF $j \in Servers \setminus \{i\}$
$\qquad\qquad\qquad\qquad\qquad$ THEN $Max(\{matchIndex[i][j],\ newCommitIndex\})$
$\qquad\qquad\qquad\qquad\qquad$ ELSE $matchIndex[i][j]]]$
$\quad$ ELSE UNCHANGED $\langle nextIndex,\ matchIndex \rangle$
$\wedge Discard(m)$
$\wedge$ UNCHANGED $\langle serverVars,\ candidateVars,\ log,\ maxc,\ leaderCount,$
$\qquad\qquad\qquad\quad hovercraftVars,\ netAggVars,\ Servers,\ netAggVars \rangle$

Server $i$ receives an *AppendEntries* request from server $j$ with
$m.mterm \leq currentTerm[i]$. This just handles $m.entries$ of length 0 or 1, but
implementations could safely accept more by treating them the same as
multiple independent requests of 1 entry.
fails when *Leader* restarts
$HandleAppendEntriesRequest(i,\ j,\ m) \triangleq$
$\quad$ LET $logOk \triangleq \vee m.mprevLogIndex = 0$
$\qquad\qquad\qquad\quad \vee \wedge m.mprevLogIndex > 0$
$\qquad\qquad\qquad\qquad \wedge m.mprevLogIndex \leq Len(log[i])$
$\qquad\qquad\qquad\qquad \wedge m.mprevLogTerm = log[i][m.mprevLogIndex].term$
$\qquad rejectHovercraftMismatchCondition \triangleq$
$\qquad\quad \wedge m.mentries \neq \langle \rangle$
$\qquad\quad \wedge$ LET $entry \triangleq m.mentries[1]$
$\qquad\qquad\qquad v \triangleq entry.value$
$\qquad\qquad\qquad msgTerm \triangleq entry.term$
$\qquad\qquad$ IN $\neg( \wedge \langle v,\ msgTerm \rangle \in unorderedRequests[i]$
$\qquad\qquad\qquad\quad \wedge \langle v,\ msgTerm \rangle \in$ DOMAIN $switchBuffer$
$\qquad\qquad\qquad\quad \wedge switchBuffer[\langle v,\ msgTerm \rangle].term = msgTerm)$

$\qquad$ Condition that triggers the CHOOSE for the leader; corner case *NACK*
$\qquad isReplyToLeaderCase \triangleq rejectHovercraftMismatchCondition \vee \neg logOk$

$\qquad$ Check if a leader exists to be chosen for the reply
$\qquad canChooseLeaderForReply \triangleq \exists l\_exists \in Servers : state[l\_exists] = Leader$

$\quad$ IN $\quad \wedge m.mterm \leq currentTerm[i]$
$\qquad\quad \wedge \vee \wedge$ reject request branch
$\qquad\qquad\qquad \wedge ($ conditions for rejecting the request
$\qquad\qquad\qquad\qquad \vee m.mterm < currentTerm[i]$
$\qquad\qquad\qquad\qquad \vee \wedge m.mterm = currentTerm[i]$

$$\wedge\ state[i] = Follower$$
$$\wedge\ \neg logOk$$
$$\vee\ \wedge\ m.mterm = currentTerm[i]$$
$$\wedge\ state[i] = Follower$$
$$\wedge\ rejectHovercraftMismatchCondition$$
$$)$$
$$\wedge\ \text{LET}\ respondTo\ \triangleq\ \text{IF}\ isReplyToLeaderCase$$
$$\text{THEN CHOOSE}\ l \in Servers : state[l] = Leader$$
$$\text{ELSE}\quad m.msource$$
$$\text{IN}\quad Reply([mtype\qquad\qquad \mapsto AppendEntriesResponse,$$
$$mterm\qquad\qquad \mapsto currentTerm[i],$$
$$msuccess\qquad\quad\ \mapsto \text{FALSE},$$
$$mmatchIndex\quad\ \mapsto 0,$$
$$msource\qquad\qquad \mapsto i,$$
$$mdest\qquad\qquad\ \mapsto respondTo],$$
$$m)$$
$$\wedge\ \text{UNCHANGED}\ \langle serverVars,\ logVars,\ unorderedRequests\rangle$$

$\vee$ ⬚ return to follower state
$$\wedge\ m.mterm = currentTerm[i]$$
$$\wedge\ state[i] = Candidate$$
$$\wedge\ state' = [state\ \text{EXCEPT}\ ![i] = Follower]$$
$$\wedge\ \text{UNCHANGED}\ \langle currentTerm,\ votedFor,\ logVars,\ messages,$$
$$unorderedRequests\rangle$$

$\vee$ ⬚ accept request
$$\wedge\ m.mterm = currentTerm[i]$$
$$\wedge\ state[i] = Follower$$
$$\wedge\ logOk$$
$$\wedge\ \text{LET}\ index\ \triangleq\ m.mprevLogIndex + 1$$
$$respondToIfAccepted\ \triangleq\ m.msource\ \ \text{respondTo will be }m.msource\text{ here}$$
$$\text{IN}\quad \vee\ \text{already done with request or empty entries}$$
$$\wedge\ \vee\ m.mentries = \langle\rangle$$
$$\vee\ \wedge\ m.mentries \neq \langle\rangle$$
$$\wedge\ Len(log[i]) \geq index$$
$$\wedge\ log[i][index].term = m.mentries[1].term$$
$$\wedge\ commitIndex' = [commitIndex\ \text{EXCEPT}\ ![i] = m.mcommitIndex]$$
$$\wedge\ Reply([mtype\qquad\qquad \mapsto AppendEntriesResponse,$$
$$mterm\qquad\qquad \mapsto currentTerm[i],$$
$$msuccess\qquad\quad\ \mapsto \text{TRUE},$$
$$mmatchIndex\quad\ \mapsto m.mprevLogIndex + Len(m.mentries),$$
$$msource\qquad\qquad \mapsto i,$$
$$mdest\qquad\qquad\ \mapsto respondToIfAccepted],$$
$$m)$$
$$\wedge\ \text{UNCHANGED}\ \langle serverVars,\ log,\ unorderedRequests\rangle$$

$\lor$ conflict: remove 1 entry
    $\land\ m.mentries \neq \langle\rangle$
    $\land\ Len(log[i]) \geq index$
    $\land\ log[i][index].term \neq m.mentries[1].term$
    $\land$ LET $newLog \triangleq SubSeq(log[i],\ 1,\ index - 1)$
      IN   $log' = [log$ EXCEPT $![i] = newLog]$
    $\land$ UNCHANGED $\langle serverVars,\ commitIndex,\ messages,\ unorderedRequests\rangle$

$\lor$ no conflict: append entry
    $\land\ m.mentries \neq \langle\rangle$
    $\land\ Len(log[i]) = m.mprevLogIndex$
    $\land\ \neg rejectHovercraftMismatchCondition$
    $\land$ LET $entryMetadata \triangleq m.mentries[1]$
              $vt \triangleq \langle entryMetadata.value,\ entryMetadata.term\rangle$
              $fullEntryFromCache \triangleq switchBuffer[vt]$
              $entryForLocalLog \triangleq [term \mapsto entryMetadata.term,$
                              $value \mapsto entryMetadata.value,$
                              $payload \mapsto fullEntryFromCache.payload]$
      IN   $log' = [log$ EXCEPT $![i] = Append(log[i],\ entryForLocalLog)]$
          $\land\ unorderedRequests' = [unorderedRequests$ EXCEPT $![i] = @ \setminus \{vt\}]$
    $\land$ UNCHANGED $\langle serverVars,\ commitIndex,\ messages\rangle$
$\land$ UNCHANGED $\langle candidateVars,\ leaderVars,\ instrumentationVars,$
    $switchBuffer,\ switchIndex,\ switchSentRecord,\ Servers,\ netAggVars\rangle$

Server $i$ receives an *AppendEntries* response from server $j$ with
$m.mterm = currentTerm[i]$.
$HandleAppendEntriesResponse(i,\ j,\ m) \triangleq$
    $\land\ m.mterm = currentTerm[i]$
    $\land\ \lor\ \land\ m.msuccess$ successful
        $\land$ LET
            $newMatchIndex \triangleq m.mmatchIndex$
            $entryKey \triangleq$ IF $newMatchIndex > 0 \land newMatchIndex \leq Len(log[i])$
                      THEN $\langle newMatchIndex,\ log[i][newMatchIndex].term\rangle$
                      ELSE $\langle 0,\ 0\rangle$ Invalid index or empty *log*
        IN    $\land\ nextIndex'\ \ = [nextIndex$ EXCEPT $![i][j] = m.mmatchIndex + 1]$
                $\land\ matchIndex' = [matchIndex$ EXCEPT $![i][j] = m.mmatchIndex]$
                $\land\ entryCommitStats' =$
                  IF $\land\ entryKey \neq \langle 0,\ 0\rangle$
                        $\land\ entryKey \in$ DOMAIN $entryCommitStats$
                        $\land\ \neg entryCommitStats[entryKey].committed$
                  THEN $[entryCommitStats$ EXCEPT $![entryKey].ackCount = @ + 1]$
                  ELSE $entryCommitStats$
    $\lor\ \land\ \neg m.msuccess$ not successful
      $\land\ nextIndex' = [nextIndex$ EXCEPT $![i][j] =$
                      $Max(\{nextIndex[i][j] - 1,\ 1\})]$
      $\land$ UNCHANGED $\langle matchIndex,\ entryCommitStats\rangle$

$\land Discard(m)$

$\land$ UNCHANGED $\langle serverVars,\ candidateVars,\ logVars,\ maxc,\ leaderCount,$
$\qquad\qquad\qquad hovercraftVars,\ Servers,\ netAggVars\rangle$

Network state transitions

The network duplicates a message
$DuplicateMessage(m) \triangleq$

$\quad\land Send(m)$

$\quad\land$ UNCHANGED $\langle serverVars,\ candidateVars,\ leaderVars,\ logVars,$
$\qquad\qquad\qquad instrumentationVars,\ hovercraftVars,\ Servers,\ netAggVars\rangle$

The network drops a message
$DropMessage(m) \triangleq$

$\quad\land Discard(m)$

$\quad\land$ UNCHANGED $\langle serverVars,\ candidateVars,\ leaderVars,\ logVars,$
$\qquad\qquad\qquad instrumentationVars,\ hovercraftVars,\ Servers,\ netAggVars\rangle$

Receive a message.
$Receive(m) \triangleq$

$\quad$ LET $i \triangleq m.mdest$
$\qquad\quad j \triangleq m.msource$

$\quad$ IN $\quad$ Any $RPC$ with a newer term causes the recipient to advance
$\qquad\qquad$ its term first. Responses with stale terms are ignored.
$\qquad\quad \lor UpdateTerm(i,\ j,\ m)$
$\qquad\quad \lor \land m.mtype = RequestVoteRequest$
$\qquad\qquad\ \land HandleRequestVoteRequest(i,\ j,\ m)$
$\qquad\quad \lor \land m.mtype = RequestVoteResponse$
$\qquad\qquad\ \land \lor DropStaleResponse(i,\ j,\ m)$
$\qquad\qquad\qquad \lor HandleRequestVoteResponse(i,\ j,\ m)$
$\qquad\quad \lor \land m.mtype = AppendEntriesRequest$
$\qquad\qquad\ \land HandleAppendEntriesRequest(i,\ j,\ m)$
$\qquad\quad \lor \land m.mtype = AppendEntriesResponse$
$\qquad\qquad\ \land \lor DropStaleResponse(i,\ j,\ m)$
$\qquad\qquad\qquad \lor HandleAppendEntriesResponse(i,\ j,\ m)$

Modified. Leader $i$ sends $j$ an $AppendEntries$ request containing exactly 1 entry.
While implementations may want to send more than 1 at a time, this spec uses
just 1 because it minimizes atomic regions without loss of generality.
Sending empty entries is done for telling followers $Leader$ is alive.
$AppendEntries(i,\ j) \triangleq$

$\quad\land i \neq j$
$\quad\land state[i] = Leader$
$\quad\land Len(log[i]) > 0$

$\quad$ Only proceed if the leader has entries to send
$\quad\land nextIndex[i][j] \leq Len(log[i])$

$\land\ matchIndex[i][j] < nextIndex[i][j]$
$\land\ $ LET $entryIndex\ \triangleq\ nextIndex[i][j]$
$\qquad entry\ \triangleq\ log[i][entryIndex]$
$\qquad entryMetadata\ \triangleq\ [term \mapsto entry.term,\ value \mapsto entry.value]$
$\qquad entries\ \triangleq\ \langle entryMetadata \rangle$
$\qquad entryKey\ \triangleq\ \langle entryIndex,\ entry.term \rangle$
$\qquad prevLogIndex\ \triangleq\ entryIndex - 1$
$\qquad prevLogTerm\ \triangleq\ $ IF $prevLogIndex > 0$ THEN
$\qquad\qquad\qquad\qquad\qquad\qquad log[i][prevLogIndex].term$
$\qquad\qquad\qquad\qquad\quad $ ELSE
$\qquad\qquad\qquad\qquad\qquad 0$

$\qquad lastEntry\ \triangleq\ Min(\{Len(log[i]),\ nextIndex[i][j]\})$
$\qquad entries\ \triangleq\ SubSeq(log[i],\ nextIndex[i][j],\ lastEntry)$

$\quad $ IN $\quad Send([mtype \qquad\qquad\quad \mapsto AppendEntriesRequest,$
$\qquad\qquad\qquad mterm \qquad\qquad\qquad \mapsto currentTerm[i],$
$\qquad\qquad\qquad mprevLogIndex \ \mapsto prevLogIndex,$
$\qquad\qquad\qquad mprevLogTerm \ \ \mapsto prevLogTerm,$
$\qquad\qquad\qquad mentries \qquad\qquad \mapsto entries,$

$\qquad\qquad\qquad mlog \qquad\qquad\qquad\ \mapsto log[i],$
$\qquad\qquad\qquad mcommitIndex \quad\ \mapsto Min(\{commitIndex[i],\ entryIndex\}),$
$\qquad\qquad\qquad msource \qquad\qquad\ \mapsto i,$
$\qquad\qquad\qquad mdest \qquad\qquad\quad\ \mapsto j])$
$\land\quad entryCommitStats' =$
$\qquad $ IF $entryKey \in$ DOMAIN $entryCommitStats$
$\qquad\qquad \land\, \neg entryCommitStats[entryKey].committed$
$\qquad $ THEN $[entryCommitStats$ EXCEPT $![entryKey].sentCount = @ + 1]$
$\qquad $ ELSE $\ entryCommitStats$
$\land$ UNCHANGED $\langle serverVars,\ candidateVars,\ leaderVars,\ logVars,\ maxc,$
$\qquad\qquad\qquad\quad leaderCount,\ hovercraftVars,\ Servers,\ netAggVars \rangle$

$MySwitchPlusPlusNext\ \triangleq$
$\quad \lor\, \exists\, i \in Servers,\ v \in Value :$
$\qquad state[i] = Leader \land SwitchClientRequest(switchIndex,\ i,\ v)$

$\quad \lor\, \exists\, v \in$ DOMAIN $switchBuffer :$
$\qquad SwitchClientRequestReplicateAll(switchIndex,\ v)$

$\quad \lor\, \exists\, i \in Servers,\ v \in$ DOMAIN $switchBuffer :$
$\qquad state[i] = Leader \land LeaderIngestHovercRaftRequest(i,\ v)$

23

$\lor\ \exists\, i \in Servers :$
   $state[i] = Leader \land AppendEntriesToNetAgg(i)$

$\lor\ \exists\, m \in \{msg \in ValidMessage(messages) :$
   $msg.mtype = AppendEntriesNetAggRequest\} : NetAggForwardAppendEntriesAll(m)$

$\lor\ \exists\, m \in \{msg \in ValidMessage(messages) :$
   $msg.mtype \in \{AppendEntriesRequest\}\} :$
   $Receive(m)$

$\lor\ \exists\, m \in \{msg \in ValidMessage(messages) :$
   $msg.mtype = AppendEntriesResponse\ \land$
   $msg.mdest = netAggIndex\} :$
   $NetAggHandleAppendEntriesResponse(m)$

$\lor\ \exists\, i \in Servers,\ m \in \{msg \in ValidMessage(messages) :$
   $msg.mtype = AggCommit\} : m.mdest = i \land HandleAggCommit(i, m)$

$\lor\ \exists\, m \in \{msg \in ValidMessage(messages) :$
   $msg.mtype = AppendEntriesResponse\ \land$
   $msg.mdest \in Servers \land state[msg.mdest] = Leader\} :$
   $\text{LET } i \ \triangleq\ m.mdest$
      $j \ \triangleq\ m.msource$
   $\text{IN }\quad AppendEntries(i, j)$

$\lor\ \exists\, i \in Servers : Restart(i)$

$\lor\ \exists\, i \in Servers : Timeout(i)$
$\lor\ \exists\, i \in Servers : BecomeLeader(i)$
$\lor\ \exists\, i, j \in Servers : RequestVote(i, j)$
$\lor\ \exists\, m \in \{msg \in ValidMessage(messages) :$
   $msg.mtype \in \{RequestVoteRequest,\ RequestVoteResponse\}\} :$
   $Receive(m)$

$\lor\ \exists\, i \in Servers : \exists\, vt \in unorderedRequests[i] :$
   $state[i] = Follower \land FollowerDropRequest(i, vt)$

$MySwitchPlusPlusSpec \triangleq MyInit \wedge \square[MySwitchPlusPlusNext]_{vars}$
$Spec \triangleq Init \wedge \square[MySwitchPlusPlusNext]_{vars}$

─────────────── Invariants ───────────────

$MoreThanOneLeaderInv \triangleq$
    $\forall i, j \in Server :$
      $(\wedge currentTerm[i] = currentTerm[j]$
        $\wedge state[i] = Leader$
        $\wedge state[j] = Leader)$
      $\Rightarrow i = j$

Every (index, term) pair determines a *log* prefix.
From page 8 of the Raft paper: "If two logs contain an entry with the
same index and term, then the logs are identical in all preceding entries."
$LogMatchingInv \triangleq$
    $\forall i, j \in Server : i \neq j \Rightarrow$
      $\forall n \in 1 .. min(Len(log[i]), Len(log[j])) :$
        $log[i][n].term = log[j][n].term \Rightarrow$
        $SubSeq(log[i], 1, n) = SubSeq(log[j], 1, n)$

The committed entries in every *log* are a prefix of the
leader's *log* up to the leader's term (since a next *Leader* may already be
elected without the old leader stepping down yet)
$LeaderCompletenessInv \triangleq$
    $\forall i \in Server :$
      $state[i] = Leader \Rightarrow$
      $\forall j \in Server : i \neq j \Rightarrow$
        $CheckIsPrefix(CommittedTermPrefix(j, currentTerm[i]), log[i])$

Committed *log* entries should never conflict between servers
$LogInv \triangleq$
    $\forall i, j \in Server :$
      $\vee CheckIsPrefix(Committed(i), Committed(j))$
      $\vee CheckIsPrefix(Committed(j), Committed(i))$

Note that *LogInv* checks for safety violations across space
This is a key safety invariant and should always be checked
THEOREM $MySwitchPlusPlusSpec \Rightarrow (\square LogInv \wedge \square LeaderCompletenessInv$
                        $\wedge \square LogMatchingInv \wedge \square MoreThanOneLeaderInv)$

fake inv to obtain a trace and observe progress for client requests advancing to committed.
$LeaderCommitted \triangleq$
    $\exists i \in Servers : commitIndex[i] \neq 2$

# References

[1] Kogias, M., Bugnion, E.: Hovercraft: achieving scalability and fault-tolerance for microsecond-scale datacenter services. In: Proceedings of the Fifteenth European Conference on Computer Systems. EuroSys '20. Association for Computing Machinery, New York, NY, USA (2020). https://doi.org/10.1145/3342195.3387545 . https://doi.org/10.1145/3342195.3387545

[2] Ongaro, D., Ousterhout, J.: In search of an understandable consensus algorithm. In: Proceedings of the 2014 USENIX Conference on USENIX Annual Technical Conference. USENIX ATC'14, pp. 305–320. USENIX Association, USA (2014)

[3] Kogias, M., Prekas, G., Ghosn, A., Fietz, J., Bugnion, E.: R2p2: making rpcs first-class datacenter citizens. In: Proceedings of the 2019 USENIX Conference on Usenix Annual Technical Conference. USENIX ATC '19, pp. 863–879. USENIX Association, USA (2019)