

---

MODULE *raft*

---

This is the formal specification for the Raft consensus algorithm.

Modified by *Ovidiu Marcu*. Simplified model and performance invariants added.

Copyright 2014 *Diego Ongaro*.

This work is licensed under the Creative Commons Attribution – 4.0

International License <https://creativecommons.org/licenses/by/4.0/>

EXTENDS *Naturals, FiniteSets, Sequences, TLC*

The set of server *IDs*

CONSTANTS *Server*

The set of requests that can go into the *log*

CONSTANTS *Value*

Server states.

CONSTANTS *Follower, Candidate, Leader*

A reserved value.

CONSTANTS *Nil*

Message types:

CONSTANTS *RequestVoteRequest, RequestVoteResponse,*  
*AppendEntriesRequest, AppendEntriesResponse*

CONSTANTS *MaxClientRequests*

Maximum times a server can become a leader

CONSTANTS *MaxBecomeLeader*

---

Global variables

A bag of records representing requests and responses sent from one server to another. *TLAPS* doesn't support the Bags module, so this is a function mapping Message to *Nat*.

VARIABLE *messages*

maximum client requests so far

VARIABLE *maxc*

A history variable used in the proof. This would not be present in an implementation.

Keeps track of successful elections, including the initial logs of the leader and voters' logs. Set of functions containing various things about successful elections (see *BecomeLeader*).

VARIABLE *elections*

A history variable used in the proof. This would not be present in an implementation.

Keeps track of every *log* ever in the system (set of logs).

VARIABLE *allLogs*

Counter for how many times each server has become leader

VARIABLE *leaderCount*

Maximum term number allowed in the model

CONSTANTS *MaxTerm*

---

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

The server's term number.

VARIABLE *currentTerm*

The server's state (Follower, *Candidate*, or *Leader*).

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. Unfortunately, the Sequence module defines *Head(s)* as the entry with index 1, so be careful not to use that!

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*  $\triangleq \langle nextIndex, matchIndex, elections \rangle$

End of per server variables.

All variables; used for stuttering (asserting state hasn't changed).

*vars*  $\triangleq \langle messages, allLogs, serverVars, candidateVars, leaderVars, logVars, maxc, leaderCount \rangle$

### Helpers

The set of all quorums. This just calculates simple majorities, but the only important property is that every quorum overlaps with every other.

*Quorum*  $\triangleq \{i \in \text{SUBSET}(\text{Server}) : \text{Cardinality}(i) * 2 > \text{Cardinality}(\text{Server})\}$

The term of the last entry in a *log*, or 0 if the *log* is empty.

*LastTerm*(*xlog*)  $\triangleq \text{IF } \text{Len}(\text{xlog}) = 0 \text{ THEN } 0 \text{ ELSE } \text{xlog}[\text{Len}(\text{xlog})].\text{term}$

Helper for *Send* and *Reply*. Given a message *m* and bag of messages, return a new bag of messages with one more *m* in it.

*WithMessage*(*m*, *msgs*)  $\triangleq$   
 IF *m*  $\in$  DOMAIN *msgs* THEN  
   [*msgs* EXCEPT ![*m*] = *msgs*[*m*] + 1]  
 ELSE  
   *msgs* @ @ (*m* :> 1)

Helper for *Discard* and *Reply*. Given a message *m* and bag of messages, return a new bag of messages with one less *m* in it.

*WithoutMessage*(*m*, *msgs*)  $\triangleq$   
 IF *m*  $\in$  DOMAIN *msgs* THEN  
   IF *msgs*[*m*]  $\leq 1$  THEN [*i*  $\in$  DOMAIN *msgs* \ {*m*}  $\mapsto$  *msgs*[*i*]]  
   ELSE [*msgs* EXCEPT ![*m*] = *msgs*[*m*] - 1]  
 ELSE  
   *msgs*

*WithMessage*(*m*, *msgs*)  $\triangleq$

IF *m*  $\in$  DOMAIN *msgs* THEN  
   *msgs*  
   [*msgs* EXCEPT ![*m*] = IF *msgs*[*m*] < 2 THEN *msgs*[*m*] + 1 ELSE 2]  
 ELSE  
   *msgs* @ @ (*m* :> 1)

*WithoutMessage*(*m*, *msgs*)  $\triangleq$

IF  $m \in \text{DOMAIN } msgs$  THEN  
 $[msgs \text{ EXCEPT } ![m] = \text{IF } msgs[m] > 0 \text{ THEN } msgs[m] - 1 \text{ ELSE } 0]$   
 ELSE  
 $msgs$

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.

$Discard(m) \triangleq messages' = WithoutMessage(m, messages)$

Combination of *Send* and *Discard*

$Reply(response, request) \triangleq$   
 $messages' = WithoutMessage(request, WithMessage(response, messages))$

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.

$Max(s) \triangleq \text{CHOOSE } x \in s : \forall y \in s : x \geq y$

---

Define initial values for all variables

$InitHistoryVars \triangleq$   
 $\wedge elections = \{\}$   
 $\wedge allLogs = \{\}$   
 $\wedge voterLog = [i \in Server \mapsto [j \in \{\} \mapsto \langle \rangle]]$   
 $InitServerVars \triangleq$   
 $\wedge currentTerm = [i \in Server \mapsto 1]$   
 $\wedge state = [i \in Server \mapsto Follower]$   
 $\wedge votedFor = [i \in Server \mapsto Nil]$   
 $InitCandidateVars \triangleq$   
 $\wedge votesResponded = [i \in Server \mapsto \{\}]$   
 $\wedge votesGranted = [i \in Server \mapsto \{\}]$

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.

$InitLeaderVars \triangleq$   
 $\wedge nextIndex = [i \in Server \mapsto [j \in Server \mapsto 1]]$   
 $\wedge matchIndex = [i \in Server \mapsto [j \in Server \mapsto 0]]$   
 $InitLogVars \triangleq$   
 $\wedge log = [i \in Server \mapsto \langle \rangle]$   
 $\wedge commitIndex = [i \in Server \mapsto 0]$   
 $Init \triangleq$   
 $\wedge messages = [m \in \{\} \mapsto 0]$   
 $\wedge InitHistoryVars$   
 $\wedge InitServerVars$   
 $\wedge InitCandidateVars$   
 $\wedge InitLeaderVars$   
 $\wedge InitLogVars$   
 $\wedge maxc = 0$   
 $\wedge leaderCount = [i \in Server \mapsto 0]$

---

Define state transitions

Server  $i$  restarts from stable storage.

It loses everything but its *currentTerm*, *votedFor*, and *log*.

$$\begin{aligned}
\text{Restart}(i) &\triangleq \\
&\wedge \text{state}' = [\text{state} \text{ EXCEPT } ![i] = \text{Follower}] \\
&\wedge \text{votesResponded}' = [\text{votesResponded} \text{ EXCEPT } ![i] = \{\}] \\
&\wedge \text{votesGranted}' = [\text{votesGranted} \text{ EXCEPT } ![i] = \{\}] \\
&\wedge \text{voterLog}' = [\text{voterLog} \text{ EXCEPT } ![i] = [j \in \{\} \mapsto \langle \rangle]] \\
&\wedge \text{nextIndex}' = [\text{nextIndex} \text{ EXCEPT } ![i] = [j \in \text{Server} \mapsto 1]] \\
&\wedge \text{matchIndex}' = [\text{matchIndex} \text{ EXCEPT } ![i] = [j \in \text{Server} \mapsto 0]] \\
&\wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] = 0] \\
&\wedge \text{UNCHANGED } \langle \text{messages}, \text{currentTerm}, \text{votedFor}, \text{log}, \text{elections}, \text{maxc}, \text{leaderCount} \rangle
\end{aligned}$$

Server  $i$  times out and starts a new election.

$$\begin{aligned}
\text{Timeout}(i) &\triangleq \wedge \text{state}[i] \in \{\text{Follower}, \text{Candidate}\} \\
&\wedge \text{currentTerm}[i] < \text{MaxTerm} \\
&\wedge \text{state}' = [\text{state} \text{ EXCEPT } ![i] = \text{Candidate}] \\
&\wedge \text{currentTerm}' = [\text{currentTerm} \text{ EXCEPT } ![i] = \text{currentTerm}[i] + 1] \\
&\text{Most implementations would probably just set the local vote} \\
&\text{atomically, but messaging localhost for it is weaker.} \\
&\wedge \text{votedFor}' = [\text{votedFor} \text{ EXCEPT } ![i] = \text{Nil}] \\
&\wedge \text{votesResponded}' = [\text{votesResponded} \text{ EXCEPT } ![i] = \{\}] \\
&\wedge \text{votesGranted}' = [\text{votesGranted} \text{ EXCEPT } ![i] = \{\}] \\
&\wedge \text{voterLog}' = [\text{voterLog} \text{ EXCEPT } ![i] = [j \in \{\} \mapsto \langle \rangle]] \\
&\wedge \text{UNCHANGED } \langle \text{messages}, \text{leaderVars}, \text{logVars}, \text{maxc}, \text{leaderCount} \rangle
\end{aligned}$$

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

$$\begin{aligned}
\text{RequestVote}(i, j) &\triangleq \\
&\wedge \text{state}[i] = \text{Candidate} \\
&\wedge j \notin \text{votesResponded}[i] \\
&\wedge \text{Send}([mtype \mapsto \text{RequestVoteRequest}, \\
&\quad mterm \mapsto \text{currentTerm}[i], \\
&\quad mlastLogTerm \mapsto \text{LastTerm}(\text{log}[i]), \\
&\quad mlastLogIndex \mapsto \text{Len}(\text{log}[i]), \\
&\quad msource \mapsto i, \\
&\quad mdest \mapsto j]) \\
&\wedge \text{UNCHANGED } \langle \text{serverVars}, \text{candidateVars}, \text{leaderVars}, \text{logVars}, \text{maxc}, \text{leaderCount} \rangle
\end{aligned}$$

Leader  $i$  sends  $j$  an *AppendEntries* request containing up to 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.

$$\begin{aligned}
\text{AppendEntries}(i, j) &\triangleq \\
&\wedge i \neq j \\
&\wedge \text{state}[i] = \text{Leader}
\end{aligned}$$

$\wedge \text{LET } prevLogIndex \triangleq nextIndex[i][j] - 1$   
 $prevLogTerm \triangleq \text{IF } prevLogIndex > 0 \text{ THEN}$   
 $\quad log[i][prevLogIndex].term$   
 $\quad \text{ELSE}$   
 $\quad 0$   
 $\quad \text{Send up to 1 entry, constrained by the end of the } log.$   
 $lastEntry \triangleq Min(\{Len(log[i]), nextIndex[i][j]\})$   
 $entries \triangleq SubSeq(log[i], nextIndex[i][j], lastEntry)$   
 $\text{IN } Send([mtype \mapsto AppendEntriesRequest,$   
 $\quad mterm \mapsto currentTerm[i],$   
 $\quad mprevLogIndex \mapsto prevLogIndex,$   
 $\quad mprevLogTerm \mapsto prevLogTerm,$   
 $\quad mentries \mapsto entries,$   
 $\quad mlog \text{ is used as a history variable for the proof.}$   
 $\quad \text{It would not exist in a real implementation.}$   
 $\quad mlog \mapsto log[i],$   
 $\quad mcommitIndex \mapsto Min(\{commitIndex[i], lastEntry\}),$   
 $\quad msource \mapsto i,$   
 $\quad mdest \mapsto j])$   
 $\wedge \text{UNCHANGED } \langle serverVars, candidateVars, leaderVars, logVars, maxc, leaderCount \rangle$

Candidate  $i$  transitions to leader.

$BecomeLeader(i) \triangleq$   
 $\wedge state[i] = Candidate$   
 $\wedge votesGranted[i] \in Quorum$   
 $\wedge leaderCount[i] < MaxBecomeLeader$   
 $\wedge state' = [state \text{ EXCEPT } ![i] = Leader]$   
 $\wedge nextIndex' = [nextIndex \text{ EXCEPT } ![i] =$   
 $\quad [j \in Server \mapsto Len(log[i]) + 1]]$   
 $\wedge matchIndex' = [matchIndex \text{ EXCEPT } ![i] =$   
 $\quad [j \in Server \mapsto 0]]$   
 $\wedge elections' = elections \cup$   
 $\quad \{[eterm \mapsto currentTerm[i],$   
 $\quad \quad eleader \mapsto i,$   
 $\quad \quad elog \mapsto log[i],$   
 $\quad \quad evotes \mapsto votesGranted[i],$   
 $\quad \quad evoterLog \mapsto voterLog[i]]\}$   
 $\wedge leaderCount' = [leaderCount \text{ EXCEPT } ![i] = leaderCount[i] + 1]$   
 $\wedge \text{UNCHANGED } \langle messages, currentTerm, votedFor, candidateVars, logVars, maxc \rangle$

Leader  $i$  receives a client request to add  $v$  to the  $log$ .

$ClientRequest(i, v) \triangleq$   
 $\wedge state[i] = Leader$   
 $\wedge \text{LET } entry \triangleq [term \mapsto currentTerm[i],$   
 $\quad value \mapsto v]$

$$\begin{aligned}
& newLog \triangleq Append(log[i], entry) \\
& \text{IN } log' = [log \text{ EXCEPT } ![i] = newLog] \\
& \wedge \text{UNCHANGED } \langle messages, serverVars, candidateVars, \\
& \quad leaderVars, commitIndex \rangle
\end{aligned}$$

$$\begin{aligned}
ValidMessage(msgs) & \triangleq \\
& \{m \in \text{DOMAIN } messages : msgs[m] > 0\}
\end{aligned}$$

$$\begin{aligned}
ClientRequest(i, v) & \triangleq \\
& \wedge state[i] = Leader \\
& \wedge maxc < MaxClientRequests \\
& \wedge \text{LET } entry \triangleq [term \mapsto currentTerm[i], \\
& \quad value \mapsto v] \\
& \quad entryExists \triangleq \\
& \quad \exists s \in Server : \\
& \quad \quad \exists j \in \text{DOMAIN } log[s] : \\
& \quad \quad \quad log[s][j].value = v \\
& \quad \quad \exists j \in \text{DOMAIN } log[i] : log[i][j].value = v \wedge log[i][j].term = currentTerm[i] \\
& newLog \triangleq \text{IF } entryExists \\
& \quad \text{THEN } log[i] \quad \text{Keep } log \text{ unchanged if entry exists} \\
& \quad \text{ELSE } Append(log[i], entry) \quad \text{Otherwise append the new entry}
\end{aligned}$$

$$\begin{aligned}
& \text{IN} \\
& \wedge log' = [log \text{ EXCEPT } ![i] = newLog] \\
& \wedge maxc' = \text{IF } entryExists \text{ THEN } maxc \text{ ELSE } maxc + 1 \\
& \wedge \text{UNCHANGED } \langle messages, serverVars, candidateVars, \\
& \quad leaderVars, commitIndex, leaderCount \rangle
\end{aligned}$$

Leader  $i$  advances its  $commitIndex$ .

This is done as a separate step from handling  $AppendEntries$  responses, in part to minimize atomic regions, and in part so that leaders of single-server clusters are able to mark entries committed.

$$\begin{aligned}
AdvanceCommitIndex(i) & \triangleq \\
& \wedge state[i] = Leader \\
& \wedge \text{LET } \text{The set of servers that agree up through index.} \\
& \quad Agree(index) \triangleq \{i\} \cup \{k \in Server : \\
& \quad \quad \quad matchIndex[i][k] \geq index\} \\
& \quad \text{The maximum indexes for which a quorum agrees} \\
& \quad agreeIndexes \triangleq \{index \in 1 \dots Len(log[i]) : \\
& \quad \quad \quad Agree(index) \in Quorum\} \\
& \quad \text{New value for } commitIndex'[i] \\
& newCommitIndex \triangleq \\
& \quad \text{IF } \wedge agreeIndexes \neq \{\} \\
& \quad \quad \wedge log[i][Max(agreeIndexes)].term = currentTerm[i] \\
& \quad \text{THEN} \\
& \quad \quad Max(agreeIndexes) \\
& \quad \text{ELSE}
\end{aligned}$$

$commitIndex[i]$   
 IN  $commitIndex' = [commitIndex \text{ EXCEPT } ![i] = newCommitIndex]$   
 $\wedge \text{UNCHANGED } \langle messages, serverVars, candidateVars, leaderVars, log, maxc, leaderCount \rangle$

---

Message handlers

$i = \text{recipient}, j = \text{sender}, m = \text{message}$

Server  $i$  receives a *RequestVote* request from server  $j$  with

$m.mterm \leq currentTerm[i]$ .

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

LET  $logOk \triangleq \vee m.mlastLogTerm > LastTerm(log[i])$   
 $\vee \wedge m.mlastLogTerm = LastTerm(log[i])$   
 $\wedge m.mlastLogIndex \geq Len(log[i])$

$grant \triangleq \wedge m.mterm = currentTerm[i]$   
 $\wedge logOk$

$\wedge votedFor[i] \in \{Nil, j\}$

IN  $\wedge m.mterm \leq currentTerm[i]$

$\wedge \vee grant \wedge votedFor' = [votedFor \text{ EXCEPT } ![i] = j]$

$\vee \neg grant \wedge \text{UNCHANGED } votedFor$

$\wedge Reply([mtype \mapsto RequestVoteResponse,$

$mterm \mapsto currentTerm[i],$

$mvoteGranted \mapsto grant,$

$mlog$  is used just for the *elections* history variable for the proof. It would not exist in a real implementation.

$mlog \mapsto log[i],$

$msource \mapsto i,$

$mdest \mapsto j],$

$m)$

$\wedge \text{UNCHANGED } \langle state, currentTerm, candidateVars, leaderVars, logVars, maxc, leaderCount \rangle$

Server  $i$  receives a *RequestVote* response from server  $j$  with

$m.mterm = currentTerm[i]$ .

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

This tallies votes even when the current state is not *Candidate*, but they won't be looked at, so it doesn't matter.

$\wedge m.mterm = currentTerm[i]$

$\wedge votesResponded' = [votesResponded \text{ EXCEPT } ![i] =$   
 $votesResponded[i] \cup \{j\}]$

$\wedge \vee \wedge m.mvoteGranted$

$\wedge votesGranted' = [votesGranted \text{ EXCEPT } ![i] =$   
 $votesGranted[i] \cup \{j\}]$

$\wedge voterLog' = [voterLog \text{ EXCEPT } ![i] =$   
 $voterLog[i] @@ (j :> m.mlog)]$

$\vee \wedge \neg m.mvoteGranted$

$\wedge \text{UNCHANGED } \langle votesGranted, voterLog \rangle$



$\wedge \text{Discard}(m)$   
 $\wedge \text{UNCHANGED } \langle \text{serverVars}, \text{votedFor}, \text{leaderVars}, \text{logVars}, \text{maxc}, \text{leaderCount} \rangle$

Server  $i$  receives an *AppendEntries* request from server  $j$  with  
 $m.mterm \leq \text{currentTerm}[i]$ . 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.

$\text{HandleAppendEntriesRequest}(i, j, m) \triangleq$   
 $\text{LET } \text{logOk} \triangleq \vee m.mprevLogIndex = 0$   
 $\vee \wedge m.mprevLogIndex > 0$   
 $\wedge m.mprevLogIndex \leq \text{Len}(\text{log}[i])$   
 $\wedge m.mprevLogTerm = \text{log}[i][m.mprevLogIndex].term$   
 $\text{IN } \wedge m.mterm \leq \text{currentTerm}[i]$   
 $\wedge \vee \wedge \text{reject request}$   
 $\vee m.mterm < \text{currentTerm}[i]$   
 $\vee \wedge m.mterm = \text{currentTerm}[i]$   
 $\wedge \text{state}[i] = \text{Follower}$   
 $\wedge \neg \text{logOk}$   
 $\wedge \text{Reply}([mtype \mapsto \text{AppendEntriesResponse},$   
 $mterm \mapsto \text{currentTerm}[i],$   
 $msuccess \mapsto \text{FALSE},$   
 $mmatchIndex \mapsto 0,$   
 $msource \mapsto i,$   
 $mdest \mapsto j],$   
 $m)$   
 $\wedge \text{UNCHANGED } \langle \text{serverVars}, \text{logVars} \rangle$   
 $\vee \text{return to follower state}$   
 $\wedge m.mterm = \text{currentTerm}[i]$   
 $\wedge \text{state}[i] = \text{Candidate}$   
 $\wedge \text{state}' = [\text{state} \text{ EXCEPT } ![i] = \text{Follower}]$   
 $\wedge \text{UNCHANGED } \langle \text{currentTerm}, \text{votedFor}, \text{logVars}, \text{messages} \rangle$   
 $\vee \text{accept request}$   
 $\wedge m.mterm = \text{currentTerm}[i]$   
 $\wedge \text{state}[i] = \text{Follower}$   
 $\wedge \text{logOk}$   
 $\wedge \text{LET } \text{index} \triangleq m.mprevLogIndex + 1$   
 $\text{IN } \vee \text{already done with request}$   
 $\wedge \vee m.mentries = \langle \rangle$   
 $\vee \wedge m.mentries \neq \langle \rangle$   
 $\wedge \text{Len}(\text{log}[i]) \geq \text{index}$   
 $\wedge \text{log}[i][\text{index}].term = m.mentries[1].term$   
 This could make our *commitIndex* decrease (for  
 example if we process an old, duplicated request),  
 but that doesn't really affect anything.  
 $\wedge \text{commitIndex}' = [\text{commitIndex} \text{ EXCEPT } ![i] =$

IF  $commitIndex[i] < m.mcommitIndex$  THEN  
      $Min(\{m.mcommitIndex, Len(log[i])\})$   
 ELSE  
      $commitIndex[i]$   
 $\wedge Reply([mtype \mapsto AppendEntriesResponse,$   
            $mterm \mapsto currentTerm[i],$   
            $msuccess \mapsto TRUE,$   
            $mmatchIndex \mapsto m.mprevLogIndex +$   
                            $Len(m.mentries),$   
            $msource \mapsto i,$   
            $mdest \mapsto j],$   
            $m)$   
 $\wedge UNCHANGED \langle serverVars, log \rangle$   
 $\vee$  **conflict: remove 1 entry**  
      $\wedge m.mentries \neq \langle \rangle$   
      $\wedge Len(log[i]) \geq index$   
      $\wedge log[i][index].term \neq m.mentries[1].term$   
      $\wedge LET new \triangleq [index2 \in 1 .. (Len(log[i]) - 1) \mapsto$   
                            $log[i][index2]]$   
         IN  $log' = [log \text{ EXCEPT } ![i] = new]$   
      $\wedge UNCHANGED \langle serverVars, commitIndex, messages \rangle$   
 $\vee$  **no conflict: append entry**  
      $\wedge m.mentries \neq \langle \rangle$   
      $\wedge Len(log[i]) = m.mprevLogIndex$   
      $\wedge log' = [log \text{ EXCEPT } ![i] =$   
                            $Append(log[i], m.mentries[1])]$   
      $\wedge UNCHANGED \langle serverVars, commitIndex, messages \rangle$   
 $\wedge UNCHANGED \langle candidateVars, leaderVars, maxc, leaderCount \rangle$

Server  $i$  receives an *AppendEntries* response from server  $j$  with  
 $m.mterm = currentTerm[i]$ .

$HandleAppendEntriesResponse(i, j, m) \triangleq$   
 $\wedge m.mterm = currentTerm[i]$   
 $\wedge \vee \wedge m.msuccess$  **successful**  
      $\wedge nextIndex' = [nextIndex \text{ EXCEPT } ![i][j] = \text{IF } matchIndex[i][j] > m.mmatchIndex \text{ THEN } matchIndex[i][j] - 1]$   
      $\wedge matchIndex' = [matchIndex \text{ EXCEPT } ![i][j] = \text{IF } matchIndex[i][j] > m.mmatchIndex \text{ THEN } matchIndex[i][j] - 1]$   
 $\vee \wedge \neg m.msuccess$  **not successful**  
      $\wedge nextIndex' = [nextIndex \text{ EXCEPT } ![i][j] =$   
                            $Max(\{nextIndex[i][j] - 1, 1\})]$   
      $\wedge UNCHANGED \langle matchIndex \rangle$   
 $\wedge Discard(m)$   
 $\wedge UNCHANGED \langle serverVars, candidateVars, logVars, elections, maxc, leaderCount \rangle$

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

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

$\wedge m.mterm > currentTerm[i]$   
 $\wedge m.mterm \leq MaxTerm$   
 $\wedge currentTerm' = [currentTerm \text{ EXCEPT } ![i] = m.mterm]$   
 $\wedge state' = [state \text{ EXCEPT } ![i] = Follower]$   
 $\wedge votedFor' = [votedFor \text{ EXCEPT } ![i] = Nil]$   
 messages is unchanged so  $m$  can be processed further.  
 $\wedge \text{UNCHANGED } \langle messages, candidateVars, leaderVars, logVars, maxc, leaderCount \rangle$

Responses with stale terms are ignored.

$DropStaleResponse(i, j, m) \triangleq$   
 $\wedge m.mterm < currentTerm[i]$   
 $\wedge Discard(m)$   
 $\wedge \text{UNCHANGED } \langle serverVars, candidateVars, leaderVars, logVars, maxc, leaderCount \rangle$

Receive a message.

$Receive(m) \triangleq$   
 LET  $i \triangleq m.mdest$   
 $j \triangleq m.msource$   
 IN Any  $RPC$  with a newer term causes the recipient to advance its term first. Responses with stale terms are ignored.  
 $\vee UpdateTerm(i, j, m)$   
 $\vee \wedge m.mtype = RequestVoteRequest$   
 $\wedge HandleRequestVoteRequest(i, j, m)$   
 $\vee \wedge m.mtype = RequestVoteResponse$   
 $\wedge \vee DropStaleResponse(i, j, m)$   
 $\vee HandleRequestVoteResponse(i, j, m)$   
 $\vee \wedge m.mtype = AppendEntriesRequest$   
 $\wedge HandleAppendEntriesRequest(i, j, m)$   
 $\vee \wedge m.mtype = AppendEntriesResponse$   
 $\wedge \vee DropStaleResponse(i, j, m)$   
 $\vee HandleAppendEntriesResponse(i, j, m)$

End of message handlers.

---

Network state transitions

The network duplicates a message

$DuplicateMessage(m) \triangleq$   
 $\wedge Send(m)$   
 $\wedge \text{UNCHANGED } \langle serverVars, candidateVars, leaderVars, logVars, maxc, leaderCount \rangle$

The network drops a message

$DropMessage(m) \triangleq$   
 $\wedge Discard(m)$   
 $\wedge \text{UNCHANGED } \langle serverVars, candidateVars, leaderVars, logVars, maxc, leaderCount \rangle$

---

Defines how the variables may transition.

$$\begin{aligned}
Next \triangleq & \bigwedge \bigvee \exists i \in Server : Timeout(i) \\
& \bigvee \exists i \in Server : Restart(i) \\
& \bigvee \exists i, j \in Server : i \neq j \wedge RequestVote(i, j) \\
& \bigvee \exists i \in Server : BecomeLeader(i) \\
& \bigvee \exists i \in Server, v \in Value : ClientRequest(i, v) \\
& \bigvee \exists i \in Server : AdvanceCommitIndex(i) \\
& \bigvee \exists i, j \in Server : i \neq j \wedge AppendEntries(i, j) \\
& \bigvee \exists m \in \{msg \in ValidMessage(messages) : \\
& \quad msg.mtype \in \{RequestVoteRequest, RequestVoteResponse, AppendEntriesRequest, AppendEntriesResponse\} : \\
& \quad \bigvee \exists m \in \{msg \in ValidMessage(messages) : \\
& \quad \quad msg.mtype \in \{AppendEntriesRequest\} : DuplicateMessage(m) \\
& \quad \bigvee \exists m \in \{msg \in ValidMessage(messages) : \\
& \quad \quad msg.mtype \in \{RequestVoteRequest\} : DropMessage(m) \\
& \quad \text{History variable that tracks every log ever:} \\
& \quad \bigwedge allLogs' = allLogs \cup \{log[i] : i \in Server\}
\end{aligned}$$

The specification must start with the initial state and transition according to *Next*.

$$Spec \triangleq Init \wedge \Box [Next]_{vars}$$

$$\begin{aligned}
MoreThanOneLeaderInv \triangleq & \\
& \forall i, j \in Server : \\
& \quad (\wedge currentTerm[i] = currentTerm[j] \\
& \quad \wedge state[i] = Leader \\
& \quad \wedge state[j] = Leader) \\
& \Rightarrow i = j
\end{aligned}$$

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

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 prefixes that contain that entry.”

$$\begin{aligned}
LogMatchingInv \triangleq & \\
& \forall i, j \in Server : i \neq j \Rightarrow \\
& \quad \forall n \in 1 \dots min(Len(log[i]), Len(log[j])) : \\
& \quad \quad log[i][n].term = log[j][n].term \Rightarrow \\
& \quad \quad SubSeq(log[i], 1, n) = SubSeq(log[j], 1, n)
\end{aligned}$$

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*

$$\text{IF } Len(log[i]) \neq 0 \wedge \exists y \in \text{DOMAIN } log[i] : log[i][y].term \leq x$$

THEN

then, we use the subsequence up to the maximum committed term of the leader

$$\text{LET } maxTermIndex \triangleq$$

$$\text{CHOOSE } y \in \text{DOMAIN } log[i] :$$

$\wedge \log[i][y].term \leq x$   
 $\wedge \forall z \in \text{DOMAIN } \log[i] : \log[i][z].term \leq x \Rightarrow y \geq z$   
 IN  $\text{SubSeq}(\log[i], 1, \min(\text{maxTermIndex}, \text{commitIndex}[i]))$   
 Otherwise the prefix is the empty tuple  
 ELSE  $\langle \rangle$

$\text{CheckIsPrefix}(\text{seq1}, \text{seq2}) \triangleq$   
 $\wedge \text{Len}(\text{seq1}) \leq \text{Len}(\text{seq2})$   
 $\wedge \forall i \in 1 \dots \text{Len}(\text{seq1}) : \text{seq1}[i] = \text{seq2}[i]$

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)

$\text{LeaderCompletenessInv} \triangleq$   
 $\forall i \in \text{Server} :$   
 $\text{state}[i] = \text{Leader} \Rightarrow$   
 $\forall j \in \text{Server} : i \neq j \Rightarrow$   
 $\text{CheckIsPrefix}(\text{CommittedTermPrefix}(j, \text{currentTerm}[i]), \log[i])$

The prefix of the *log* of server *i* that has been committed

$\text{Committed}(i) \triangleq$   
 IF  $\text{commitIndex}[i] = 0$   
 THEN  $\langle \rangle$   
 ELSE  $\text{SubSeq}(\log[i], 1, \text{commitIndex}[i])$

Committed *log* entries should never conflict between servers

$\text{LogInv} \triangleq$   
 $\forall i, j \in \text{Server} :$   
 $\vee \text{CheckIsPrefix}(\text{Committed}(i), \text{Committed}(j))$   
 $\vee \text{CheckIsPrefix}(\text{Committed}(j), \text{Committed}(i))$

Note that *LogInv* checks for safety violations across space

This is a key safety invariant and should always be checked

THEOREM  $\text{Spec} \Rightarrow \Box \text{LogInv}$

A leader's *maxc* should remain under *MaxClientRequests*

$\text{MaxCInv} \triangleq (\exists i \in \text{Server} : \text{state}[i] = \text{Leader}) \Rightarrow \text{maxc} \leq \text{MaxClientRequests}$

No server can become leader more than *MaxBecomeLeader* times

$\text{LeaderCountInv} \triangleq \exists i \in \text{Server} : (\text{state}[i] = \text{Leader} \Rightarrow \text{leaderCount}[i] \leq \text{MaxBecomeLeader})$

No server can have a term exceeding *MaxTerm*

$\text{MaxTermInv} \triangleq \forall i \in \text{Server} : \text{currentTerm}[i] \leq \text{MaxTerm}$

$\text{MyConstraint} \triangleq (\forall i \in \text{Server} : \text{currentTerm}[i] \leq 3 \wedge \text{Len}(\log[i]) \leq 3) \wedge (\forall m \in \text{DOMAIN } \text{messages} : \text{messages}$

$\text{Symmetry} \triangleq \text{Permutations}(\text{Server})$