# Safety Proof and Formal Specification for Raft

authorship removed for blind review

last updated March 28, 2013

## 1 Introduction

This document is a proof of safety for the Raft replication algorithm. Liveness is not addressed, nor is reconfiguration.

## 2 System model

This document uses the usual asynchronous system assumptions.

- Messages may take an arbitrary number of steps to arrive at a server; once they arrive, they are modeled as being processed in one atomic step.

- Servers fail by stopping and may later restart from stable storage on disk.

- The network may reorder, drop, and duplicate messages.

## 3 Conventions

The specification uses the syntax and semantics of the TLA+ language version 2. The remainder of this document uses the same syntax and semantics but with the following minor allowances for convenience.

- As in TLA+, $foo'$ has a specific meaning: the value of variable $foo$ in the next state of the system.

- We say $\langle index, term \rangle \in log$ iff
  $Len(log) \geq index \ \wedge \ log[index].term = term$.

- We use the symbol $\parallel$ for concatenation of logs and entries.

- We ignore values in log entries, since it's easy to see that a value is attached to a particular $\langle index, term \rangle$, and those uniquely identify a log entry.

## 4 Specification

This section provides a complete, formal description of the Raft algorithm.

1 ———————————————————————— MODULE $raft$ ————————————————————————
2    This is the formal specification for the Raft replication algorithm.

4 EXTENDS $Naturals$, $FiniteSets$, $Sequences$, $TLC$, $TLAPS$

6    The set of server $IDs$
7 CONSTANTS $Replica$

9   The set of requests that can go into the *log*
10  CONSTANTS *Value*

12  Server states.
13  CONSTANTS *Follower*, *Candidate*, *Leader*

15  A reserved value.
16  CONSTANTS *Nil*

18  Message types:
19  CONSTANTS *RequestVoteRequest*, *RequestVoteResponse*,
20            *AppendEntriesRequest*, *AppendEntriesResponse*

22  A bag of records representing requests and responses set from one server to another.
23  *TLAPS* doesn't support the Bags module, so this is a function mapping Message to *PNat*.
24  VARIABLE *messages*

26  A history variable used in the proof. This would not be present in an implementation.
27  Keeps track of successful elections, including the initial logs of the leader and voters' logs.
28  VARIABLE *elections*

30  A history variable used in the proof. This would not be present in an implementation.
31  Keeps track of every *log* ever in the system.
32  VARIABLE *allLogs*

34  ├────────────────────────────────────────────────────────────────────────────────┤
35  The following variables are all per server (functions with the domain *Replica*).

37  The server's term number.
38  VARIABLE *currentTerm*
39  The server's state (Follower, *Candidate*, or *Leader*).
40  VARIABLE *state*
41  The candidate the server voted for in its current term, or *Nil* if it hasn't voted for any.
42  VARIABLE *votedFor*
43  *replicaVars* $\triangleq$ $\langle currentTerm, state, votedFor \rangle$

45  A Sequence of *log* entries. The index into this sequence is the index of the *log* entry.
46  Unfortunately, the Sequence module defines *Head(s)* as the entry with index 1, so be careful not to use that!
47  VARIABLE *log*
48  The latest entry the state machine may apply is *commitIndex*.
49  VARIABLE *commitIndex*
50  *logVars* $\triangleq$ $\langle log, commitIndex \rangle$

52  The following variables are used only on candidates:
53  The set of servers from which the candidate has received a *RequestVote* response in this term.
54  VARIABLE *votesResponded*
55  The set of servers from which the candidate has received a vote in this term.
56  VARIABLE *votesGranted*
57  A history variable used in the proof. This would not be present in an implementation.
58  Keeps track of the *log* of each voter.
59  VARIABLE *voterLog*
60  *candidateVars* $\triangleq$ $\langle votesResponded, votesGranted, voterLog \rangle$

62  The following variables are used only on leaders:
63  The next entry to send to each follower.

64    VARIABLE $nextIndex$

65      The latest entry that each follower has acknowledged is the same as the leader's.

66      This is used to calculate $commitIndex$ on the leader.

67    VARIABLE $lastAgreeIndex$

68    $leaderVars \triangleq \langle nextIndex, \ lastAgreeIndex, \ elections \rangle$

70      End of per server variables.

71 ⊢──────────────────────────────────────────────────────────────────────

73      All variables; used for stuttering.

74    $vars \triangleq \langle messages, \ replicaVars, \ candidateVars, \ leaderVars, \ logVars \rangle$

76 ⊢──────────────────────────────────────────────────────────────────────

77      Helpers

79      The set of all quorums. This just calculates simple majorities, but the only

80      important property is that every quorum overlaps with every other.

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

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

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

86      Helper for $Send$ and $Reply$

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

88        IF $m \in \text{DOMAIN} \ msgs$ THEN

89           $[msgs \ \text{EXCEPT} \ ![m] = msgs[m] + 1]$

90         ELSE

91           $msgs \ @@ \ (m :> 1)$

93      Helper for $Discard$ and $Reply$

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

95        IF $m \in \text{DOMAIN} \ msgs$ THEN

96           $[msgs \ \text{EXCEPT} \ ![m] = msgs[m] - 1]$

97         ELSE

98           $msgs$

100      Add a message to the set of messages.

101    $Send(m) \triangleq messages' = WithMessage(m, \ messages)$

103      The recipient is done processing the message.

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

106      Combination of $Send$ and $Discard$

107    $Reply(response, \ request) \triangleq$

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

110      Return the minimum value from a set, or undefined if the set is empty.

111    $Min(s) \triangleq \text{CHOOSE} \ x \in s : \forall \, y \in s : x \leq y$

112      Return the maximum value from a set, or undefined if the set is empty.

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

115 ⊢──────────────────────────────────────────────────────────────────────

117      Define initial values for all variables

118    $InitHistoryVars \triangleq \ \land elections = \{\}$

119                       $\land allLogs \quad = \{\}$

$$120 \qquad\qquad\qquad\qquad \wedge\ voterLog = [i \in Replica \mapsto [j \in \{\} \mapsto \langle\rangle]]$$

$$121 \quad InitReplicaVars \ \triangleq\ \wedge\ currentTerm = [i\ \in Replica \mapsto 1]$$

$$122 \qquad\qquad\qquad\qquad \wedge\ state \qquad\ = [i\ \in Replica \mapsto Follower]$$

$$123 \qquad\qquad\qquad\qquad \wedge\ votedFor \quad\ = [i\ \in Replica \mapsto Nil]$$

$$124 \quad InitCandidateVars \ \triangleq\ \wedge\ votesResponded = [i \in Replica \mapsto \{\}]$$

$$125 \qquad\qquad\qquad\qquad\quad \wedge\ votesGranted \quad = [i \in Replica \mapsto \{\}]$$

126    The values $nextIndex[i][i]$ and $lastAgreeIndex[i][i]$ are never read, since

127    the leader does not send itself messages. It's still easier to include these

128    in the functions.

$$129 \quad InitLeaderVars \ \triangleq\ \wedge\ nextIndex \qquad\ = [i \in Replica \mapsto [j \in Replica \mapsto 1]]$$

$$130 \qquad\qquad\qquad\qquad \wedge\ lastAgreeIndex\ = [i \in Replica \mapsto [j \in Replica \mapsto 0]]$$

$$131 \quad InitLogVars \ \triangleq\ \wedge\ log \qquad\qquad = [i \in Replica \mapsto \langle\rangle]$$

$$132 \qquad\qquad\qquad \wedge\ commitIndex\ = [i \in Replica \mapsto 0]$$

$$133 \quad Init \ \triangleq\ \wedge\ messages = [m \in \{\} \mapsto 0]$$

$$134 \qquad\qquad \wedge\ InitHistoryVars$$

$$135 \qquad\qquad \wedge\ InitReplicaVars$$

$$136 \qquad\qquad \wedge\ InitCandidateVars$$

$$137 \qquad\qquad \wedge\ InitLeaderVars$$

$$138 \qquad\qquad \wedge\ InitLogVars$$

140    The network duplicates a message

$$141 \quad DuplicateMessage \ \triangleq\ \wedge\ \exists\, m \in \textsc{domain}\ messages :$$

$$142 \qquad\qquad\qquad\qquad\qquad Send(m)$$

$$143 \qquad\qquad\qquad\qquad \wedge\ \textsc{unchanged}\ \langle replicaVars,\ candidateVars,\ leaderVars,\ logVars\rangle$$

145    The network drops a message

$$146 \quad DropMessage \ \triangleq\ \wedge\ \exists\, m \in \textsc{domain}\ (messages) :$$

$$147 \qquad\qquad\qquad\qquad Discard(m)$$

$$148 \qquad\qquad\qquad \wedge\ \textsc{unchanged}\ \langle replicaVars,\ candidateVars,\ leaderVars,\ logVars\rangle$$

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

$$151 \quad Timeout(i) \ \triangleq\ \wedge\ state[i] \in \{Follower,\ Candidate\}$$

$$152 \qquad\qquad\qquad \wedge\ state' = [state\ \textsc{except}\ ![i] = Candidate]$$

$$153 \qquad\qquad\qquad \wedge\ currentTerm' = [currentTerm\ \textsc{except}\ ![i] = currentTerm[i] + 1]$$

154      Most implementations would probably just set the local vote

155      atomically, but messaging localhost for it is weaker.

$$156 \qquad\qquad\qquad \wedge\ votedFor' = [votedFor\ \textsc{except}\ ![i] = Nil]$$

$$157 \qquad\qquad\qquad \wedge\ votesResponded' = [votesResponded\ \textsc{except}\ ![i] = \{\}]$$

$$158 \qquad\qquad\qquad \wedge\ votesGranted' \quad = [votesGranted\ \textsc{except}\ ![i] = \{\}]$$

$$159 \qquad\qquad\qquad \wedge\ voterLog' \qquad\ = [voterLog\ \textsc{except}\ ![i] = [j \in \{\} \mapsto \langle\rangle]]$$

$$160 \qquad\qquad\qquad \wedge\ \textsc{unchanged}\ \langle messages,\ leaderVars,\ logVars\rangle$$

162    Server $i$ restarts from stable storage.

163    It loses everything but its $currentTerm$, $votedFor$, and $log$.

$$164 \quad Restart(i) \ \triangleq\ \wedge\ \textsc{unchanged}\ \langle messages,\ currentTerm,\ votedFor,\ log,\ elections\rangle$$

$$165 \qquad\qquad\qquad \wedge\ state' \qquad\qquad = [state\ \textsc{except}\ ![i] = Follower]$$

$$166 \qquad\qquad\qquad \wedge\ votesResponded' = [votesResponded\ \textsc{except}\ ![i] = \{\}]$$

$$167 \qquad\qquad\qquad \wedge\ votesGranted' \quad = [votesGranted\ \textsc{except}\ ![i] = \{\}]$$

$$168 \qquad\qquad\qquad \wedge\ voterLog' \qquad\ = [voterLog\ \textsc{except}\ ![i] = [j \in \{\} \mapsto \langle\rangle]]$$

$$169 \qquad\qquad\qquad \wedge\ nextIndex' \qquad = [nextIndex\ \textsc{except}\ ![i] = [j \in Replica \mapsto 1]]$$

$$170 \qquad\qquad\qquad \wedge\ lastAgreeIndex' \ = [lastAgreeIndex\ \textsc{except}\ ![i] = [j \in Replica \mapsto 0]]$$

$$171 \qquad\qquad\qquad \wedge\ commitIndex' \quad = [commitIndex\ \textsc{except}\ ![i] = 0]$$

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

174   $RequestVote(i, j) \triangleq \land state[i] = Candidate$

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

176                          $\land Send([mtype \quad\quad \mapsto RequestVoteRequest,$

177                                     $mterm \quad\quad\quad \mapsto currentTerm[i],$

178                                     $mlastLogTerm \mapsto LastTerm(log[i]),$

179                                     $mlastLogIndex \mapsto Len(log[i]),$

180                                     $msource \quad\quad\quad \mapsto i,$

181                                     $mdest \quad\quad\quad \mapsto j])$

182                          $\land$ UNCHANGED $\langle replicaVars, candidateVars, leaderVars, logVars \rangle$

184     Leader $i$ sends $j$ an *AppendEntries* request containing up to 1 entry.

185     While implementations may want to send more than 1 at a time, this spec uses

186     just 1 because it minimizes atomic regions without loss of generality.

187   $AppendEntries(i, j) \triangleq$

188         $\land i \neq j$

189         $\land state[i] = Leader$

190         $\land$ LET $prevLogIndex \triangleq nextIndex[i][j] - 1$

191                  $prevLogTerm \triangleq$ IF $prevLogIndex > 0$ THEN

192                                     $log[i][prevLogIndex].term$

193                            ELSE

194                                     0

195                 Send up to 1 entry, constrained by the end of the *log*.

196                 $lastEntry = Min(\{Len(log[i]), nextIndex[i][j] + 1\})$

197                 $entries \triangleq SubSeq(log[i], nextIndex[i][j], upper)$

198       IN     $Send([mtype \quad\quad\quad \mapsto AppendEntriesRequest,$

199                       $mterm \quad\quad\quad\quad \mapsto currentTerm[i],$

200                       $mprevLogIndex \mapsto prevLogIndex,$

201                       $mprevLogTerm \quad \mapsto prevLogTerm,$

202                       $mentries \quad\quad\quad \mapsto entries,$

203                       $mcommitIndex \quad \mapsto Min(\{commitIndex[i], lastEntry\}),$

204                       $msource \quad\quad\quad \mapsto i,$

205                       $mdest \quad\quad\quad\quad \mapsto j])$

206       $\land$ UNCHANGED $\langle replicaVars, candidateVars, leaderVars, logVars \rangle$

208     Candidate $i$ transitions to leader.

209   $BecomeLeader(i) \triangleq \land state[i] = Candidate$

210                           $\land votesGranted[i] \in Quorum$

211                           $\land state' = [state$ EXCEPT $![i] = Leader]$

212                           $\land nextIndex' \quad\quad = [nextIndex$ EXCEPT $![i] =$

213                                       $[j \in Replica \mapsto Len(log[i]) + 1]]$

214                           $\land lastAgreeIndex' = [lastAgreeIndex$ EXCEPT $![i] =$

215                                       $[j \in Replica \mapsto 0]]$

216                           $\land elections' = elections \cup$

217                                     $\{[e.eterm \quad\quad \mapsto currentTerm[i],$

218                                       $e.eleader \quad \mapsto i,$

219                                       $e.elog \quad\quad\quad \mapsto log[i],$

220                                       $e.evotes \quad\quad \mapsto votesGranted[i],$

221                                       $e.evoterLog \mapsto voterLog[i]]\}$

222                           $\land$ UNCHANGED $\langle messages, currentTerm, votedFor, candidateVars, logVars \rangle$

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

$225 \quad ClientRequest(i, v) \triangleq \land state[i] = Leader$

$226 \qquad\qquad\qquad\qquad \land \text{LET } entry \triangleq [term \mapsto currentTerm[i],$

$227 \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad value \mapsto v]$

$228 \qquad\qquad\qquad\qquad\qquad\quad newIndex \triangleq Len(log[i]) + 1$

$229 \qquad\qquad\qquad\qquad\qquad\quad\; newLog \triangleq Append(log[i], entry)$

$230 \qquad\qquad\qquad\qquad\quad \text{IN} \quad log' = [log \text{ EXCEPT } ![i] = newLog]$

$231 \qquad\qquad\qquad\qquad \land \text{UNCHANGED } \langle messages,\, replicaVars,\, candidateVars,\, leaderVars,$

$232 \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad commitIndex\rangle$

$234 \vdash$ ————————————————————————————————————————

238    Server $i$ receives a *RequestVote* request from server $j$ with $m.mterm \leq currentTerm[i]$.

$239 \quad HandleRequestVoteRequest(i, j, m) \triangleq$

$240 \qquad \text{LET } logIsCurrent \triangleq \lor m.mlastLogTerm > LastTerm(log[i])$

$241 \qquad\qquad\qquad\qquad\qquad\quad \lor \land m.mlastLogTerm = LastTerm(log[i])$

$242 \qquad\qquad\qquad\qquad\qquad\qquad\;\; \land m.mlastLogIndex \geq Len(log[i])$

$243 \qquad\qquad grant \triangleq \land m.mterm = currentTerm[i]$

$244 \qquad\qquad\qquad\qquad\; \land logIsCurrent$

$245 \qquad\qquad\qquad\qquad\; \land votedFor[i] \in \{Nil, j\}$

$246 \qquad \text{IN} \quad \land m.mterm \leq currentTerm[i]$

$247 \qquad\qquad \land \lor grant \quad \land votedFor' = [votedFor \text{ EXCEPT } ![i] = j]$

$248 \qquad\qquad\quad \lor \neg grant \land \text{UNCHANGED } votedFor$

$249 \qquad\qquad \land Reply([mtype \qquad\qquad \mapsto RequestVoteResponse,$

$250 \qquad\qquad\qquad\qquad\; mterm \qquad\qquad \mapsto currentTerm[i],$

$251 \qquad\qquad\qquad\qquad\; mvoteGranted \mapsto grant,$

252    *mlog* is used just for the 'elections' history variable for

253    the proof. It would not exist in a real implementation.

$254 \qquad\qquad\qquad\qquad\; mlog \qquad\qquad \mapsto log[i],$

$255 \qquad\qquad\qquad\qquad\; msource \qquad\; \mapsto i,$

$256 \qquad\qquad\qquad\qquad\; mdest \qquad\qquad \mapsto j],$

$257 \qquad\qquad\qquad\qquad\; m)$

$258 \qquad\qquad \land \text{UNCHANGED } \langle state,\, currentTerm,\, candidateVars,\, leaderVars,\, logVars\rangle$

260    Server $i$ receives a *RequestVote* response from server $j$ with $m.mterm = currentTerm[i]$.

$261 \quad HandleRequestVoteResponse(i, j, m) \triangleq$

$262 \qquad \land m.mterm = currentTerm[i]$

$263 \qquad \land votesResponded' = [votesResponded \text{ EXCEPT } ![i] = votesResponded[i] \cup \{j\}]$

$264 \qquad \land \lor m.mvoteGranted \quad \land votesGranted' = [votesGranted \text{ EXCEPT } ![i] = votesGranted[i] \cup \{j\}]$

$265 \qquad\qquad\qquad\qquad\qquad\qquad\quad \land voterLog' = [voterLog \text{ EXCEPT } ![i][j] = m.mlog]$

$266 \qquad\quad \lor \neg m.mvoteGranted \land \text{UNCHANGED } \langle votesGranted,\, voterLog\rangle$

$267 \qquad \land Discard(m)$

$268 \qquad \land \text{UNCHANGED } \langle replicaVars,\, votedFor,\, leaderVars,\, logVars\rangle$

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

271    This just handles $m.entries$ of length 0 or 1, but implementations could safely accept

272    more by treating them the same as multiple independent requests of 1 entry.

$273 \quad HandleAppendEntriesRequest(i, j, m) \triangleq$

$274 \qquad \text{LET } accept \triangleq \land m.mterm = currentTerm[i]$

$275 \qquad\qquad\qquad\qquad\;\; \land \lor m.mprevLogIndex = 0$

$276 \qquad\qquad\qquad\qquad\qquad\; \lor \land m.mprevLogIndex > 0$

$277 \qquad\qquad\qquad\qquad\qquad\qquad\;\; \land m.mprevLogIndex \leq Len(log[i])$

$$278 \qquad\qquad\qquad\qquad \land\ m.mprevLogTerm = log[i][m.mprevLogIndex].term$$

$$279 \quad \text{IN} \quad \land\ m.mterm \le currentTerm[i]$$

$$280 \qquad \land\ \lor\ \boxed{\text{reject request}}$$

$$281 \qquad\qquad\qquad \land\ \neg accept$$

$$282 \qquad\qquad\qquad \land\ Reply([mtype \qquad\qquad \mapsto AppendEntriesResponse,$$

$$283 \qquad\qquad\qquad\qquad\qquad mterm \qquad\qquad \mapsto currentTerm[i],$$

$$284 \qquad\qquad\qquad\qquad\qquad mlastAgreeIndex \mapsto 0,$$

$$285 \qquad\qquad\qquad\qquad\qquad msource \qquad\qquad \mapsto i,$$

$$286 \qquad\qquad\qquad\qquad\qquad mdest \qquad\qquad \mapsto j],$$

$$287 \qquad\qquad\qquad\qquad\qquad m)$$

$$288 \qquad\qquad\qquad \land\ \text{UNCHANGED}\ \langle replicaVars,\ candidateVars,\ leaderVars,\ logVars \rangle$$

$$289 \qquad\quad \lor\ \land\ accept$$

$$290 \qquad\qquad\qquad \land\ \text{LET}\ index\ \triangleq\ m.mprevLogIndex + 1$$

$$291 \qquad\qquad\qquad\quad \text{IN} \quad \lor\ \boxed{\text{already done with request}}$$

$$292 \qquad\qquad\qquad\qquad\qquad \land\ \lor\ m.mentries = \langle\rangle$$

$$293 \qquad\qquad\qquad\qquad\qquad\qquad \lor\ \land\ Len(log[i]) \ge index$$

$$294 \qquad\qquad\qquad\qquad\qquad\qquad\qquad \land\ log[i][index].term = m.mentries[1].term$$

$$295 \qquad\qquad\qquad\qquad\qquad \boxed{\text{This could make our } commitIndex \text{ decrease (for}}$$

$$296 \qquad\qquad\qquad\qquad\qquad \boxed{\text{example if we process an old, duplicated request),}}$$

$$297 \qquad\qquad\qquad\qquad\qquad \boxed{\text{but that doesn't really affect anything.}}$$

$$298 \qquad\qquad\qquad\qquad\qquad \land\ commitIndex' = [commitIndex\ \text{EXCEPT}\ ![i] = m.mcommitIndex]$$

$$299 \qquad\qquad\qquad\qquad\qquad \land\ Reply([mtype \qquad\qquad \mapsto AppendEntriesResponse,$$

$$300 \qquad\qquad\qquad\qquad\qquad\qquad\qquad mterm \qquad\qquad \mapsto currentTerm[i],$$

$$301 \qquad\qquad\qquad\qquad\qquad\qquad\qquad mlastAgreeIndex \mapsto m.mprevLogIndex + Len(m.mentries),$$

$$302 \qquad\qquad\qquad\qquad\qquad\qquad\qquad msource \qquad\qquad \mapsto i,$$

$$303 \qquad\qquad\qquad\qquad\qquad\qquad\qquad mdest \qquad\qquad \mapsto j],$$

$$304 \qquad\qquad\qquad\qquad\qquad\qquad\qquad m)$$

$$305 \qquad\qquad\qquad\qquad\qquad \land\ \text{UNCHANGED}\ \langle logVars \rangle$$

$$306 \qquad\qquad\qquad\qquad \lor\ \boxed{\text{conflict: remove 1 entry}}$$

$$307 \qquad\qquad\qquad\qquad\qquad \land\ m.mentries \ne \langle\rangle$$

$$308 \qquad\qquad\qquad\qquad\qquad \land\ Len(log[i]) \ge index$$

$$309 \qquad\qquad\qquad\qquad\qquad \land\ log[i][index].term \quad \ne m.mentries[1].term$$

$$310 \qquad\qquad\qquad\qquad\qquad \land\ \text{LET}\ new\ \triangleq\ [index2 \in 1\ ..\ (Len(log[i]) - 1) \mapsto log[i][index2]]$$

$$311 \qquad\qquad\qquad\qquad\qquad\quad \text{IN} \quad log' = [log\ \text{EXCEPT}\ ![i] = new]$$

$$312 \qquad\qquad\qquad\qquad\qquad \land\ \text{UNCHANGED}\ \langle commitIndex,\ messages \rangle$$

$$313 \qquad\qquad\qquad\qquad \lor\ \boxed{\text{no conflict: add entry}}$$

$$314 \qquad\qquad\qquad\qquad\qquad \land\ m.mentries \ne \langle\rangle$$

$$315 \qquad\qquad\qquad\qquad\qquad \land\ Len(log[i]) = m.mprevLogIndex$$

$$316 \qquad\qquad\qquad\qquad\qquad \land\ log' = [log\ \text{EXCEPT}\ ![i] = Append(log[i],\ m.mentries[1])]$$

$$317 \qquad\qquad\qquad\qquad\qquad \land\ \text{UNCHANGED}\ \langle commitIndex,\ messages \rangle$$

$$318 \qquad\qquad \land\ \text{UNCHANGED}\ \langle replicaVars,\ candidateVars,\ leaderVars \rangle$$

$$320 \quad \boxed{\text{Server } i \text{ receives an } AppendEntries \text{ response from server } j \text{ with } m.mterm = currentTerm[i].}$$

$$321 \quad HandleAppendEntriesResponse(i,\ j,\ m)\ \triangleq$$

$$322 \qquad \land\ m.mterm = currentTerm[i]$$

$$323 \qquad \land\ \lor\ \land\ m.mlastAgreeIndex > 0$$

$$324 \qquad\qquad\quad \land\ nextIndex' \qquad = [nextIndex \qquad\qquad \text{EXCEPT}\ ![i][j] = m.mlastAgreeIndex + 1]$$

$$325 \qquad\qquad\quad \land\ lastAgreeIndex' = [lastAgreeIndex\ \text{EXCEPT}\ ![i][j] = m.mlastAgreeIndex]$$

$$326 \qquad\qquad\quad \land\ \text{LET} \quad \boxed{Agree(index) \text{ is the set of replicas that agree up through } index.}$$

$$327 \qquad\qquad\qquad\qquad Agree(index)\ \triangleq\ \{i\} \cup \{k \in Replica : lastAgreeIndex'[i][k] \ge index\}$$

$$328 \qquad\qquad\qquad\qquad \boxed{\text{The set of indexes for which a quorum agrees}}$$

$$329 \qquad\qquad\qquad\qquad agreeIndexes\ \triangleq\ \{index \in 1\ ..\ Len(log[i]) : Agree(index) \in Quorum\}$$

330                   New value for $commitIndex'[i]$

331            $newCommitIndex \triangleq$ IF $\wedge\ agreeIndexes \neq \{\}$

332                            $\wedge\ log[i][Max(agreeIndexes)].term = currentTerm[i]$

333                   THEN

334                       $Max(agreeIndexes)$

335                   ELSE

336                       $commitIndex[i]$

337         IN     $commitIndex' = [commitIndex$ EXCEPT $![i] = newCommitIndex]$

338      $\vee\ \wedge\ m.mlastAgreeIndex = 0$

339         $\wedge\ nextIndex' = [nextIndex$ EXCEPT $![i][j] = Max(\{nextIndex[i][j] - 1, 1\})]$

340         $\wedge$ UNCHANGED $\langle lastAgreeIndex, commitIndex\rangle$

341    $\wedge\ Discard(m)$

342    $\wedge$ UNCHANGED $\langle replicaVars, candidateVars, log, elections\rangle$

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

345 $UpdateTerm(i, j, m) \triangleq\ \wedge\ m.mterm > currentTerm[i]$

346                  $\wedge\ currentTerm'\quad = [currentTerm$ EXCEPT $![i] = m.mterm]$

347                  $\wedge\ state'\qquad\quad = [state\qquad$ EXCEPT $![i]\quad = Follower]$

348                  $\wedge\ votedFor'\qquad = [votedFor\quad$ EXCEPT $![i]\quad = Nil]$

349                 messages is unchanged so $m$ can be processed further.

350                  $\wedge$ UNCHANGED $\langle messages, candidateVars, leaderVars, logVars\rangle$

352    Responses with stale terms are ignored.

353 $DropStaleResponse(i, j, m) \triangleq\ \wedge\ m.mterm < currentTerm[i]$

354                       $\wedge\ Discard(m)$

355                       $\wedge$ UNCHANGED $\langle replicaVars, candidateVars, leaderVars, logVars\rangle$

357    Receive a message.

358 $Receive \triangleq \exists\, m \in$ DOMAIN $messages :$

359            LET $i \triangleq m.mdest$

360                $j \triangleq m.msource$

361          IN      Any $RPC$ with a newer term causes the recipient to advance its term first.

362                  Responses with stale terms are ignored.

363                $\vee\ UpdateTerm(i, j, m)$

364                $\vee\ m.mtype = RequestVoteRequest\qquad \wedge\ HandleRequestVoteRequest(i, j, m)$

365                $\vee\ m.mtype = RequestVoteResponse\qquad \wedge\ \vee\ DropStaleResponse(i, j, m)$

366                                            $\vee\ HandleRequestVoteResponse(i, j, m)$

367                $\vee\ m.mtype = AppendEntriesRequest\ \wedge\ HandleAppendEntriesRequest(i, j, m)$

368                $\vee\ m.mtype = AppendEntriesResponse\ \wedge\ \vee\ DropStaleResponse(i, j, m)$

369                                            $\vee\ HandleAppendEntriesResponse(i, j, m)$

371    End of message handlers.

372 ├───────────────────────────────────────────────────────

374    Defines how the variables may transition.

375 $Next \triangleq\ \wedge\ \vee\ DuplicateMessage$

376               $\vee\ DropMessage$

377               $\vee\ Receive$

378               $\vee\ \exists\, i \in Replica : Timeout(i)$

379               $\vee\ \exists\, i \in Replica : Restart(i)$

380               $\vee\ \exists\, i \in Replica : BecomeLeader(i)$

381               $\vee\ \exists\, i \in Replica, v \in Value : ClientRequest(i, v)$

382               $\vee\ \exists\, i, j \in Replica : RequestVote(i, j)$

383     $\lor \exists\, i, j \in Replica : AppendEntries(i, j)$
384            History variable that tracks every *log* ever:
385     $\land\ allLogs' = allLogs \cup \{\forall\, i \in Replica : log[i]\}$

387     The specification must start with the initial state and transition according to *Next*.
388     $Spec \triangleq Init \land \Box[Next]_{vars}$

390

# 5    Definitions

**Definition 1.** An entry $\langle index, term \rangle$ is **committed at term** $t$ if it is present in every leader's log following $t$:

$$committed(t) \triangleq \{\langle index, term \rangle :$$
$$\forall\ election \in elections :$$
$$election.term > t \Rightarrow$$
$$\langle index, term \rangle \in election.log\}$$

**Definition 2.** An entry $\langle index, term \rangle$ is **immediately committed** if it is acknowledged by a quorum (including the leader) during *term*. Theorem 9 shows that these entries are *committed at term*.

$$immediatelyCommitted \triangleq \{\langle index, term \rangle :$$
$$\exists\ leader, subquorum :$$
$$\land\ subquorum \cup \{leader\} \in Quorum$$
$$\land\ \forall\, i \in subquorum :$$
$$\exists\, m \in messages :$$
$$\land\ m.msource = i$$
$$\land\ m.mdest = leader$$
$$\land\ m.mtype = AppendEntriesResponse$$
$$\land\ m.mterm = term$$
$$\land\ m.mlastAgreeIndex \geq index\}$$

**Definition 3.** An entry $\langle index, term \rangle$ is **prefix committed at term** $t$ if there is another entry that is *committed at term t* following it in some log. Theorem 10 shows that these entries are *committed at term t*.

$$prefixCommitted(t) \triangleq \{\langle index, term \rangle :$$
$$\exists\, i \in Server :$$
$$\land\ \langle index, term \rangle \in log[i]$$
$$\land\ \exists\, \langle rindex, rterm \rangle \in log[i] :$$
$$\land\ index < rindex$$
$$\land\ \langle rindex, rterm \rangle \in committed(t)\}$$

# 6    Proof

**Lemma 1.** Each server's *currentTerm* monotonically increases:

$$\forall\, i \in Server :$$
$$currentTerm[i] \leq currentTerm'[i]$$

*Proof.* This follows immediately from the specification. $\qquad\square$

**Lemma 2.** There is at most one leader per term:

$$\forall \ e, f \in elections :$$
$$e.term = f.term \Rightarrow e = f$$

*Sketch.* It takes votes from a quorum to become leader, voters may only vote once per term, and any two quorums overlap. Moreover, no single server starts multiple elections in the same term.

*Proof.*

1. Consider two elections, $e$ and $f$, both members of *elections*, where $e.term = f.term$.

2. $e.votes \in Quroum$ and $f.votes \in Quorum$, since this is a necessary condition for elements of *elections*.

3. Let *voter* be an arbitrary member of $e.votes \cap f.votes$. Such a member must exist since any two quorums overlap.

4. Once *voter* casts a vote for $e.leader$ in $e.term$, it can not cast a vote for a different server in $e.term$ (the specification ensures this: once it increments its *currentTerm*, it can never vote again for the same server (Lemma 1); and until then, it safely retains its vote information).

5. $e.leader = f.leader$, since *voter* voted for $e.leader$ and *voter* voted for $f.leader$ in $e.term = f.term$.

6. A server does not start more than one election in the same term, since it increments its term upon starting a new election, and its *currentTerm* monotonically increases.

7. Therefore, $e = f$.

$\square$

**Lemma 3.** A leader's log monotonically grows during its term:

$$\forall \ e \in elections :$$
$$currentTerm[e.leader] = e.term \Rightarrow$$
$$\forall \ index \in Len(log[e.leader]) :$$
$$log'[e.leader][index] = log[e.leader][index]$$

*Sketch.* While $state[i] = Leader$, $i$ only appends to its log. $i$ won't ever get an AppendEntries request from some other server for this term, since there is at most one leader per term. And $i$ rejects AppendEntries requests for other terms until increasing term.

*Proof.*

1. Three variables are involved: *elections*, *currentTerm*, and *log*.

2. When a new election is added to *elections*, the *log* of the leader is not changed in the same step, so the invariant is maintained.

3. Logs can change either from client requests or AppendEntries requests:

   (a) Case: client request: follows directly from the spec
   (b) Case: AppendEntries request
       i. Only servers with $state[i] = Leader$ can send AppendEntries requests for their *currentTerm*.
       ii. By Lemma 2, $e.leader$ is the only leader for $e.term$.
       iii. Servers don't send themselves AppendEntries requests (see spec).
       iv. So $e.leader$ will receive no AppendEntries requests during $e.term$.

4. *currentTerm*[*e.leader*] monotonically increases by Lemma 1, so once *e.leader* moves to a new term, it will trivially satisfy the invariant forever after.

<div align="right">□</div>

**Lemma 4.** An ⟨*index*, *term*⟩ identifies a log prefix:

$$\forall \ l, m \in allLogs :$$
$$\quad \forall \ index, term :$$
$$\quad\quad \langle index, term \rangle \in l \wedge \langle index, term \rangle \in m \Rightarrow$$
$$\quad\quad\quad \forall \ pindex \in 1..index :$$
$$\quad\quad\quad\quad l[pindex] = m[pindex]$$

*Sketch.* Only leaders create entries, and they assign the new entries term numbers that will never be assigned again by other leaders (there's at most one leader per term). The consistency check in AppendEntries guarantees that when followers accept new entries, they do so in a way that's consistent with the leader's log at the time it sent the entries.

*Assertion.* If *log′*[*i*] is a prefix of some log in *allLogs*, then *allLogs′* = *allLogs* ∪ {*log′*[*i*]} maintains the invariant.

*Proof by induction on an execution.*

1. Initial state: all logs are empty, so the invariant holds.

2. Inductive step: logs change by either:

   (a) Case: a leader adds one entry (client request)

      i. This ⟨*index*, *term*⟩ uniquely identifies this entry, since there's only one leader per term (Lemma 2) and leaders only append to their logs.
      ii. By the inductive hypothesis, *log*[*leader*] ∈ *allLogs*.
      iii. Then *allLogs′* = *allLogs* ∪ {*log*[*leader*] ∥ ⟨*index*, *term*⟩} maintains the invariant.

   (b) Case: a follower, *follower*, removes one entry (AppendEntries request *m*)

      i. The invariant still holds, since *log′*[*follower*] is a prefix of *log*[*follower*] (by the Assertion above).

   (c) Case: a follower, *follower*, adds one entry (AppendEntries request *m*)

      i. Let $l \triangleq log[m.msource]$ at the time the leader creates the AppendEntries request.
      ii. *l* ∈ *allLogs* by definition of *allLogs*.
      iii. In the two cases below, we show that *log′*[*follower*] is a prefix of *l*. By the Assertion above, this suffices to show that the invariant is maintained.
      iv. Case: *m.mentries* is a prefix of *l*.
         A. *m.mprevLogIndex* = 0, since nothing precedes *m.mentries* in *l*.
         B. *log*[*follower*] is empty, as a necessary condition for accepting the request.
         C. *log′*[*follower*] = *m.mentries* upon accepting the request, which is a prefix of *l*.
      v. Case: *start* ∥ ⟨*m.mprevLogIndex*, *m.mprevLogTerm*⟩ ∥ *m.mentries* is a prefix of *l*, where *start* is some (possibly empty) log prefix.
         A. The follower accepts the request by assumption, so it contains the entry ⟨*m.mprevLogIndex*, *m.mprevLogTerm*⟩.
         B. By the inductive hypothesis, *log*[*follower*] contains the prefix *start* ∥ ⟨*m.mprevLogIndex*, *m.mprevLogTerm*⟩.
         C. *follower′*[*log*] = *start* ∥ ⟨*m.mprevLogIndex*, *m.mprevLogTerm*⟩ ∥ *m.mentries* upon accepting the request, which is a prefix of *l*.

$\square$

**Lemma 5.** A server's current term is always at least as large as the terms in its log:

$$\forall\ i \in Server :$$
$$\forall\ index \in 1..Len(log[i]) :$$
$$log[i][index].term \leq currentTerm[i]$$

*Sketch.* Servers' current terms monotonically increase. When leaders create new entries, they assign them their current term. And when followers accept new entries from a leader, they agree with the leader's term at the time it sent the entries.

*Proof by induction on an execution.*

1. Initial state: all logs are empty, so the invariant holds.

2. Inductive step: logs change by either:

   (a) Case: a leader adds one entry (client request)

      i. By the inductive hypothesis, all prior entries in $log[leader]$ have term $\leq currentTerm[leader]$.
      ii. The new entry's term is $currentTerm[leader]$.

   (b) Case: a follower removes one entry (AppendEntries request)

      i. The invariant still holds, since only the length of the log decreased.

   (c) Case: a follower adds one entry (AppendEntries request $m$)

      i. By the inductive hypothesis, $currentTerm[leader]$ was at least as large as the term in every entry of $log[leader]$ when $leader$ created the request.
      ii. The leader's log contained $m.mentries$ and $m.mterm = currentTerm[leader]$ at the time $leader$ created the request.
      iii. $follower'[log]$ ends in $m.mentries$ upon accepting the request.
      iv. The leader's log contained every entry present in $follower'[log]$ when it created the request, since they both have $m.mentries$ in common (by Lemma 4).
      v. As a necessary condition for accepting the request, $currentTerm[follower] = m.mterm$.
      vi. Then $currentTerm[follower]$ is at least as large as the term in every entry in $log'[follower]$, and the invariant is maintained.

3. Inductive step: $currentTerm[i]$ changes

   (a) By Lemma 1, $currentTerm'[i] \geq currentTerm[i]$, so the invariant is maintained.

$\square$

**Lemma 6.** The terms of indexes grow monotonically in a server's log:

$$\forall\ i \in Server :$$
$$\forall\ index \in 1..(Len(log[i]) - 1) :$$
$$log[i][index].term \leq log[i][index + 1].term$$

*Sketch.* A leader maintains this by assigning new entries its current term, which is always at least as large as the terms in its log. When followers accept new entries, they are consistent with the leader's log at the time it sent the entries.

*Proof by induction on an execution.*

1. Initial state: all logs are empty, so the invariant holds.

2. Inductive step: logs change by either:

    (a) Case: a leader adds one entry (client request)

        i. The new entry's term is $currentTerm[leader]$

        ii. $currentTerm[leader]$ is at least as large as the term of any entry in $log[leader]$, by Lemma 5.

    (b) Case: a follower removes one entry (AppendEntries request)

        i. The invariant still holds, since only the length of the log decreased.

    (c) Case: a follower adds one entry (AppendEntries request)

        i. $follower'[log]$ ends in $m.mentries$ upon accepting the request.

        ii. Then $follower'[log]$ is a prefix of the leader's log at the time the leader created the request, since $m.mentries$ identifies a log prefix (Lemma 4).

        iii. By the inductive hypothesis, since the leader's log contained $m.mentries$ at the time it created the request, those entries and all preceding entries satisfy the invariant (have monotonically increasing terms).

$\square$

**Lemma 7.** If two logs have the same last term, at least one is a prefix of the other:

$$\forall\ l, m \in allLogs :$$
$$\quad LastTerm(m) = LastTerm(l) \land Len(l) \leq Len(m) \Rightarrow$$
$$\quad\quad \forall\ index \in 1..Len(l) :$$
$$\quad\quad\quad l[index].term = m[index].term$$

*Sketch.* When a leader adds its first entry in a term, its term is distinct from all other entries' terms (there is at most one leader per term), so this trivially holds. When it adds subsequent entries under that term, it only appends those, so its log satisfies this property. When followers accept new entries, they are consistent with the leader's log at the time it sent the entries.

*Proof by induction on an execution.*

1. Initial state: trivially holds for empty logs.

2. Inductive step: logs change by either:

    (a) Case: a leader adds one entry (client request)

        i. This $\langle index, term \rangle$ uniquely identifies a log entry not previously found in any of $allLogs$, since there is at most one leader per term (Lemma 2), leaders only append entries with terms set to $currentTerm[i]$, and $currentTerm[i]$ monotonically increases (Lemma 1).

        ii. Case: $term = LastTerm(log[leader])$: The invariant holds since $log[leader]$ is a prefix of $log'[leader]$.

        iii. Case: $LastTerm(log'[leader]) > LastTerm(log[leader])$ The invariant holds, since $\forall\ n \in allLogs :$ $LastTerm(log'[leader]) \neq LastTerm(n)$ (there is at most one leader per term by Lemma 2).

    (b) Case: a follower removes one entry (AppendEntries request)

        i. The invariant still holds, since $log'[follower]$ is a prefix of $log[follower]$.

    (c) Case: a follower adds one entry (AppendEntries request $m$)

        i. The leader's log contained $m.mentries$ at the time it creates the request.

        ii. $log'[follower]$ ends in $m.mentries$ upon applying the request.

        iii. By Lemma 4, $log'[follower]$ is a prefix of the leader's log at the time it created the request, which is in $allLogs$ by the inductive hypothesis.

$\square$

**Lemma 8.** A leader does not send *AppendEntries* requests that conflict with its log.

> $\forall\ i \in Server :$
> > $(\wedge\ state[i] = Leader$
> > $\wedge\ \langle index, term \rangle \in log[i]$
> > $) \Rightarrow \neg\exists\ m \in \text{DOMAIN } messages :$
> > > $\wedge\ m.mtype = AppendEntriesRequest$
> > > $\wedge\ m.mterm = currentTerm[i]$
> > > $\wedge\ \vee\ \langle index, otherTerm \rangle \in m.mentries$
> > > $\quad\ \vee\ m.mprevLogIndex = index\ \wedge\ m.mprevLogTerm = otherTerm$
> > > $\wedge\ otherTerm \neq term$

*Sketch.* Leaders only send AppendEntries using the contents of their log, and a leader's log monotonically grows.

*Proof.*

1. Since there is only one leader per term, the request would have to come from this server.

2. Leaders only send *AppendEntries* requests that are consistent with their logs.

3. By Lemma 3, a leader's log monotonically grows during its term.

4. Therefore, no server could have sent such a conflicting request.

$\square$

**Lemma 9.** *Immediately committed* entries are *committed*:

> $\forall\ \langle index, term \rangle \in immediatelyCommitted :$
> > $\langle index, term \rangle \in committed(term)$

*Sketch.* A quorum contains the entry in the term when it is created, so all leaders in future terms must receive a vote from at least one of these servers. By the log consistency check during leader election, the first leader following term must have the entry: it couldn't have a larger last log term than the voter's, so the voter's log must be a prefix of the leader's log. The next leader after that could have a larger term, but only if it inherited its log from this leader, so the initial log for this leader (which contains the entry) is a prefix of this next leader's log. This argument continues for all future leaders.

*Proof.*

1. Consider an entry $\langle index, term \rangle$ that is *immediately committed*.

2. Define

> $Contradicting \triangleq \{ election \in elections :$
> > $\wedge\ election.term > term$
> > $\wedge\ \langle index, term \rangle \notin election.log \}$

3. Let *election* be an element in *Contradicting* with a minimal *term* field. That is, $\forall\ e \in Contradicting :$ $election.term \leq e.term$.

4. It suffices to show a contradiction, which implies $Contradicting = \phi$.

5. Let *follower* be any server that both votes in *election* and contains $\langle index, term \rangle$ during *term* (either it acknowledges the entry as a follower or it was leader). Such a server must exist since:

14

(a) A quorum votes in *election*.

(b) A quorum contains $\langle index, term \rangle$ during *term*.

(c) Any two quorums overlap.

6. Let $followerLog \triangleq election.voterLog[follower]$.

7. The follower contains the entry when it cast its vote during *election.term*. That is, $\langle index, term \rangle \in followerLog$:

   (a) $\langle index, term \rangle$ was in the follower's log during *term*.

   (b) The follower must have stored the entry in *term* before voting in *election.term*, since *election.term* > *term*, *currentTerm[follower]* monotonically increases (Lemma 1), and the follower rejects requests with terms smaller than *currentTerm[follower]*.

   (c) The follower couldn't have removed the entry before casting its vote:

      i. Case: No *AppendEntriesRequest* with *mterm* < *term* removes the entry from the follower's log, since *currentTerm[follower]* $\geq$ *term* upon storing the entry (by Lemma 5), and the follower rejects requests with terms smaller than *currentTerm[follower]*.

      ii. Case: No *AppendEntriesRequest* with *mterm* = *term* removes the entry from the follower's log, since:

         A. There is only one leader of *term*.

         B. The leader of *term* created and therefore contains the entry (Lemma 3).

         C. The leader would not send any conflicting requests to *follower* during *term* (Lemma 8).

      iii. Case: No *AppendEntriesRequest* with *election.term* > *mterm* > *term* removes the entry from the follower's log. Since all of these leaders also contained the entry by assumption, they did not send any conflicting entries to the follower for this index (Lemma 8). Nor did they send any conflicting entries for prior indexes: that they have this entry implies they have the entire prefix (Lemma 4).

      iv. Case: No *AppendEntriesRequest* with *mterm* = *election.term* removes the entry from the follower's log prior to it voting, since there is at most one leader per term (Lemma 2), so this request would have to come from the leader of *election.term*, but it hasn't been elected yet.

      v. Case: No *AppendEntriesRequest* with *mterm* > *election.term* removes the entry from the follower's log prior to it voting, since then *currentTerm[follower]* > *election.term* would prevent the follower from voting in *term*.

8. The log completeness check during elections states the following, since *follower* granted its vote during *election*:

$$\lor\ LastTerm(election.log) > LastTerm(followerLog)$$
$$\lor\ \land\ LastTerm(election.log) = LastTerm(followerLog)$$
$$\land\ Len(election.log) \geq Len(followerLog)$$

In the following two steps, we take each of these cases in turn and show a contradiction.

9. Case: $LastTerm(election.log) > LastTerm(followerLog)$

   (a) $LastTerm(followerLog) \geq term$, since $\langle index, term \rangle \in followerLog$ and terms in logs grow monotonically (Lemma 6).

   (b) *election.term* > $LastTerm(election.log)$ since servers increment their *currentTerm* when starting an election, and Lemma 5 states that a server's *currentTerm* is at least as large as the terms in its log.

15

(c) Let *prior* be the election in *elections* with $prior.term = LastTerm(election.log)$. Such an election must exist since $LastTerm(election.log) > 0$ and a server must win an election before creating an entry.

(d) By transitivity, we now have the following inequalities:

$$\begin{aligned} term \leq \\ LastTerm(followerLog) < \\ LastTerm(election.log) = prior.term < \\ election.term \end{aligned}$$

(e) $\langle index, term \rangle \notin prior.log$, since:

    i. *prior.log* is a prefix of *election.log* by Lemma 7, since:

        A. $LastTerm(election.log) = prior.term$

        B. The leader of *prior.term* creates entries by appending them to *prior.log*, so any entry with term *prior.term* must have index greater than $Len(prior.log)$.

        C. Then $Len(election.log) > Len(prior.log)$.

    ii. Since $\langle index, term \rangle \notin election.log$ (by assumption), either $Len(election.log) < index$ or $election.log[index].term \neq term$. We consider these cases separately next.

    iii. Case: $Len(election.log) < index$

        A. Then $\langle index, term \rangle \notin prior.log$, since *prior.log* is a prefix of *election.log*.

    iv. Case: $election.log[index].term \neq term$

        A. Case: $Len(prior.log) < index$: immediate.

        B. Case: $Len(prior.log) > index$: $election.log[index].term = prior.log[index].term \neq term$, since *prior.log* is a prefix is of *election.log*.

(f) $prior \in Contradicting$, by definition of *Contradicting*.

(g) *election* was defined as an element in *Contradicting* with a minimal *term* field, but *prior* is also in *Contradicting* and $prior.term < election.term$.

10. Case: $LastTerm(election.log) = LastTerm(followerLog)$ and $Len(election.log) \geq Len(followerLog)$

    (a) *followerLog* is a prefix of *election.log* by Lemma 7.

    (b) Then $\langle index, term \rangle \in election.log$, since $\langle index, term \rangle \in followerLog$.

    (c) But $election \in Contradicting$ implies that $\langle index, term \rangle \notin election.log$.

$\square$

**Lemma 10.** *Prefix committed* entries are *committed*:

$$\forall\, t : prefixCommitted(t) \subseteq committed(t)$$

*Sketch.* If an entry is committed, it identifies a prefix of a log in which every entry is committed, since those entries will also be present in every future leader's log.

*Proof.*

1. Consider an arbitrary entry $\langle index, term \rangle \in prefixCommitted(t)$.

2. There exists an entry $\langle rindex, rterm \rangle \in committed(t)$ following $\langle index, term \rangle$ in some log, by definition of $prefixCommitted(t)$.

3. $\langle rindex, rterm \rangle$ uniquely identifies the log prefix containing $\langle index, term \rangle$ (Lemma 4).

4. Every leader following $t$ contains $\langle index, term \rangle$, since every leader following $t$ contains $\langle rindex, rterm \rangle$.

5. $\langle index, term \rangle \in committed(t)$ by definition of $committed(t)$.

<div align="right">□</div>

**Theorem 1.** Servers only apply entries that are *committed* in their current term:

$\forall\, i \in Server :$
$\quad \{\langle index, log[i][index].term \rangle : index \in 1..committedIndex[i]\} \subseteq$
$\qquad committed(currentTerm[i])$

*Sketch.* A leader only advances its committedIndex to cover entries that are immediately committed or prefix committed. Followers update committedIndex from the leader's only when they have a prefix of a prior version of the leader's log.

*Proof.*

1. The set of committed entries monotonically increases, by its definition.

2. When $committedIndex[i]$ increases, it covers entries present in $i$'s log that are committed:

   (a) Case: follower completes accepting AppendEntries request

      i. Upon processing the request, the follower's log is a prefix of a prior version of the leader's log, $l$.
      ii. Every entry up to $commitIndex'[i]$ in $l$ is committed by the inductive hypothesis.

   (b) Case: leader $i$ processes AppendEntries response

      i. If the leader sets a new $commitIndex$, the conditions in the spec ensure that $commitIndex'[i] \in immediatelyCommitted$.
      ii. Every entry with index below $commitIndex'[i]$ is prefix committed at $currentTerm[i]$.

3. No log entry below $committedIndex[i]$ is ever removed:

   (a) The leader who sent $i$ the *AppendEntries* request to mark this log prefix committed must have the log prefix (see cases above), and it will not remove any of these entries during its term (Lemma 3).

   (b) By the inductive hypothesis, every leader with a larger term than $currentTerm[i]$ must contain the entry. (And requests from leaders of smaller terms are rejected.)

   (c) If a leader has the entry, it wouldn't have the follower remove it (Lemma 8).

<div align="right">□</div>