# 15-812 Term Paper:
# Specifying and proving cluster membership for the Raft distributed consensus algorithm

Brandon Amos and Huanchen Zhang[*]

2015/05/10

[*]Equal contribution, alphabetical order.

# Contents

# 1   Introduction

Distributed consensus is popular in today's world as many large-scale production systems rely on reaching consensus among a set of decentralized servers. Consensus algorithms are notoriously difficult to correctly implement and formal verification methods are helpful in proving properties of the algorithms.

Raft is a newly released consensus algorithm that is beginning to be adopted in large-scale systems [OO14]. A partial formal specification of Raft is presented in [Ong14] and used in hand proofs for a subset of properties.

In this report, we add new functionality to the formal specification in §3. We prove (by hand) a safety property of there being at most one leader per term under our modifications in §4.1. We describe a proof sketch in §4.2 showing that at any point, a leader can be elected in the future.

# 2   Background

## 2.1   The Raft Consensus Algorithm

Raft is consensus algorithm that allows a collection of machines to work as a coherent group that can survive failures of some members and is presented at USENIX ATC'14 [OO14] and further expanded on in Diego Ongaro's thesis [Ong14]. The Secret Lives of Data [sec] provides a visual walkthrough and introduction to Raft.

Raft has moved beyond academia and is being implemented and deployed in large-scale production systems, as described on the website [raf].

Some important concepts and terms for understanding Raft are:

- **Replicated Log.** Each node maintains a log that contains values and configuration entries. Because the system is distributed, the logs aren't guaranteed to be consistent on every server. Log entries can be **committed**, which means that a majority of the nodes agree on the entry. A majority of nodes is also called a **quorum**.

- **Server states.** Servers in the cluster exist in the following three states.

    - **Leader.** The leader receives requests from external entities to append values to the replicated log.
    - **Follower.** Followers receive commands from the leader to add new entries to their logs.
    - **Candidate.** If a follower doesn't hear from a leader within a specified interval, it times out and becomes a candidate.

- **Configuration.** The configuration is the set of servers in the Raft system. The protocol allows servers to be added and removed from the system.

In this report, we study adding and removing servers from the cluster. Adding and removing servers is done by operating on one server at a time and keeping track of the configuration with the normal log replication mechanisms. The RPC's for adding and removing servers are fully described in Figure 1.

| **AddServer RPC** | **RemoveServer RPC** |
|---|---|
| Invoked by admin to add a server to the cluster configuration. | Invoked by admin to remove a server from the cluster configuration. |
| **Arguments:** | **Arguments:** |
| **newServer**    address of server to add to configuration | **oldServer**    address of server to remove from configuration |
| **Results:** | **Results:** |
| **status**    OK if server was added successfully | **status**    OK if server was removed successfully |
| **leaderHint**    address of recent leader, if known | **leaderHint**    address of recent leader, if known |
| **Receiver implementation:** | **Receiver implementation:** |
| 1. Reply NOT_LEADER if not leader (§6.2) | 1. Reply NOT_LEADER if not leader (§6.2) |
| 2. Catch up new server for fixed number of rounds. Reply TIMEOUT if new server does not make progress for an election timeout or if the last round takes longer than the election timeout. (§4.2.1) | 2. Wait until previous configuration in log is committed (§4.1) |
| 3. Wait until previous configuration in log is committed (§4.1) | 3. Append new configuration entry to log (old configuration without oldServer), commit it using majority of new configuration (§4.1) |
| 4. Append new configuration entry to log (old configuration plus newServer), commit it using majority of new configuration (§4.1) | 4. Reply OK and, if this server was removed, step down (§4.2.2) |
| 5. Reply OK | |

Figure 1: Implementation of Raft's Add and Remove RPC's. Copied from Figure 4.1 of [Ong14] and included here for completeness.

### 2.1.1 Safety and Availability

Safety and availability (or liveness) are fundamental properties systems that are important to formally verify [AS87]. The safety property of Raft we focus on is that two leaders can never be elected in the same term. An availability property of Raft is that a leader can be elected at some point in the future.

## 2.2 Temporal Logic of Actions (TLA)

Lamport's temporal logic of actions (TLA) [Lam94] is a logic for specifying and reasoning about concurrent systems. Figure 2 summarizes minimal syntax and semantics of TLA. TLA+ is formal specification language that describes system behavior using TLA [Lam02]. TLA+ breaks distributed algorithms into state transition functions that specify all possible behaviors of the system. The TLA+ Model Checker (TLC) [YML99] exhaustively checks whether a property or invariant holds. The TLA+ Proof System [CDLM08] mechanically checks TLA+ proofs. [Lam00] provides a helpful summary and description for reading and writing TLA+.

    Appendix B of Ongaro's thesis [Ong14] provides a TLA+ specification and hand-written proofs of a subset of Raft's properties and features. §8 of the thesis provides informal arguments about correctness. For completeness (and convenience), this report includes the original TLA+ specification in Appendix A.

**Syntax**

$$\langle formula \rangle \quad \triangleq \quad \langle predicate \rangle \mid \Box[\langle action \rangle]_{\langle state\ function \rangle} \mid \neg \langle formula \rangle$$
$$\mid \langle formula \rangle \wedge \langle formula \rangle \mid \Box \langle formula \rangle$$

$$\langle action \rangle \quad \triangleq \quad \text{boolean-valued expression containing constant symbols,}$$
$$\text{variables, and primed variables}$$

$$\langle predicate \rangle \quad \triangleq \quad \langle action \rangle \text{ with no primed variables} \mid Enabled\ \langle action \rangle$$

$$\langle state\ function \rangle \triangleq \text{nonboolean expression containing constant symbols and variables}$$

**Semantics**

$$s[\![f]\!] \quad \triangleq \quad f(\forall\ `v\text{'} : s[\![v]\!]/v) \qquad\qquad \sigma[\![F \wedge G]\!] \quad \triangleq \quad \sigma[\![F]\!] \wedge \sigma[\![G]\!]$$

$$s[\![\mathcal{A}]\!]t \quad \triangleq \quad \mathcal{A}(\forall\ `v\text{'} : s[\![v]\!]/v,\ t[\![v]\!]/v') \qquad \sigma[\![\neg F]\!] \quad \triangleq \quad \neg\sigma[\![F]\!]$$

$$\models \mathcal{A} \quad \triangleq \quad \forall s, t \in \mathbf{St} : s[\![\mathcal{A}]\!]t \qquad\qquad\quad \models F \qquad \triangleq \quad \forall \sigma \in \mathbf{St}^{\infty} : \sigma[\![F]\!]$$

$$s[\![Enabled\ \mathcal{A}]\!] \quad \triangleq \quad \exists t \in \mathbf{St} : s[\![\mathcal{A}]\!]t$$

$$\langle s_0, s_1, \ldots \rangle[\![\Box F]\!] \quad \triangleq \quad \forall n \in \mathsf{Nat} : \langle s_n, s_{n+1}, \ldots \rangle[\![F]\!]$$

$$\langle s_0, s_1, \ldots \rangle[\![\mathcal{A}]\!] \quad \triangleq \quad s_0[\![\mathcal{A}]\!]s_1$$

**Additional notation**

$$p' \quad \triangleq \quad p(\forall\ `v\text{'} : v'/v) \qquad\qquad \Diamond F \qquad \triangleq \quad \neg\Box\neg F$$

$$[\mathcal{A}]_f \quad \triangleq \quad \mathcal{A} \vee (f' = f) \qquad\qquad F \rightsquigarrow G \quad \triangleq \quad \Box(F \Rightarrow \Diamond G)$$

$$\langle \mathcal{A} \rangle_f \quad \triangleq \quad \mathcal{A} \wedge (f' \neq f) \qquad\qquad \mathrm{WF}_f(\mathcal{A}) \quad \triangleq \quad \Box\Diamond\langle \mathcal{A} \rangle_f \vee \Box\Diamond\neg Enabled\ \langle \mathcal{A} \rangle_f$$

$$Unchanged\ f \quad \triangleq \quad f' = f \qquad\qquad \mathrm{SF}_f(\mathcal{A}) \quad \triangleq \quad \Box\Diamond\langle \mathcal{A} \rangle_f \vee \Diamond\Box\neg Enabled\ \langle \mathcal{A} \rangle_f$$

**where**    $f$ is a $\langle state\ function \rangle$        $s, s_0, s_1, \ldots$ are states
           $\mathcal{A}$ is an $\langle action \rangle$            $\sigma$ is a behavior
           $F$ and $G$ are $\langle formula \rangle$s     $(\forall\ `v\text{'} : \ldots/v,\ \ldots/v')$ denotes substitution
           $p$ is a $\langle state\ function \rangle$ or $\langle predicate \rangle$       for all variables $v$

Figure 2: Summary of TLA's simple syntax and semantics. Copied from Figure 4 of [Lam94] and included here for completeness.

# 3    Adding cluster membership changes to Raft's formal specification

We have extended Raft's formal TLA+ specification to allow server configuration changes. For completeness, Appendix B provides our modified specification.

## 3.1    Modeling Network Messages

We utilize the existing specification for messaging between Raft nodes by using the `Send`, `Discard`, and `Reply` helper functions. Messages in the system are represented as a bag in `messages` that maps a message's content to an integer. This integer counts the number of active messages in the system and is initialized to one, incremented by one when a message is duplicated or sent again, and decremented by one when a message is discarded or replied to.

Network packets can be duplicated or dropped, which the TLA+ specification models with `Duplicate` and `Drop` in the state transition function.

## 3.2  New Variables

We have added the following new variables and constants to the specification. Our original modifications included other variables that introduced a new state for detached servers and kept track of additional indexes. However, we realized these could be deduced mathematically from other variables in the system.

- `NumRounds`. The number of rounds to catch each server up by.

- `InitServer` and `Server`. Previously, there was only a single constant describing the set of servers in the system. We have modified this to describe both an initial and global set of servers that can be added and removed.

- `ValueEntry` and `ConfigEntry`. Previously, the log only contained homogeneous entries. Now, configuration changes are also stored in the log and each entry is now identified as either a value or config with a `type` metadata.

- `CatchupRequest`, `CatchupResponse`, and `CheckOldConfig`. New message types in the system to catch up servers and check if the old config have been committed.

## 3.3  Initial state of the system

We have only slightly modified the system initialization in `Init` to correctly handle the changed set of servers. Every variable is initialized to contain information for the global set of servers, even if they aren't in the initial configuration, so that the lists do not have to be resized every time a server receives a configuration change. This prevents some corner cases when server receives a configuration change in it's log that doesn't get committed that is then overwritten by another log entry.

## 3.4  State Transitions

In the `Next` state transition definitions, we modify the existential operators to operate on the global set of servers. Some servers might not be in any configurations, so we add restrictions to the state transition functions.

- `Timeout`. A server can only timeout, become a candidate, and start a new election if it is in its own configuration.

- `RequestVote`. Candidates only request votes from servers in their configuration.

- `AppendEntries`. Leaders only send new log entries to servers in their configuration.

- `BecomeLeader`. A candidate can only become a leader if they receive votes from a majority of their quorum.

- `ClientRequest`. Unmodified, only leaders receive requests from clients to add new values to the replicated state machine.

- `AdvanceCommitIndex`. Leaders can advance the commit index if all servers in their config agree.

### 3.4.1 AddNewServer

We have added a new state transition function to add a new server to the system. This can be called when some server $i$ is the leader and adds a new server that's not in it's configuration. This sends a `CatchupRequest` message to the server to be added with log entries to append.

The first time this is called, `nextIndex[i][j]` will be 0 and the entire committed log will be sent. However, this can be called multiple times before a server is added when $i$ is still a leader, since $j$ will not be added to it's configuration until the server is sufficiently caught up. Therefore, the leader uses `nextIndex[i][j]` to keep track of the new server's state so that duplicate requests are not harmful.

---

 Leader $i$ adds a new server $j$ to the cluster.
$AddNewServer(i, j) \triangleq$
$\quad \wedge state[i] = Leader$
$\quad \wedge j \notin GetConfig(i)$
$\quad \wedge currentTerm' = [currentTerm \text{ EXCEPT } ![j] = 1]$
$\quad \wedge votedFor' = [votedFor \text{ EXCEPT } ![j] = Nil]$
$\quad \wedge Send([mtype \mapsto CatchupRequest,$
$\qquad\qquad mterm \mapsto currentTerm[i],$
$\qquad\qquad mlogLen \mapsto matchIndex[i][j],$
$\qquad\qquad mentries \mapsto SubSeq(log[i], nextIndex[i][j], commitIndex[i]),$
$\qquad\qquad mcommitIndex \mapsto commitIndex[i],$
$\qquad\qquad msource \mapsto i,$
$\qquad\qquad mdest \mapsto j,$
$\qquad\qquad mrounds \mapsto NumRounds])$
$\quad \wedge \text{UNCHANGED } \langle state, leaderVars, logVars, candidateVars \rangle$

---

### 3.4.2 DeleteServer

Deleting a server is simpler than adding a server because no catching up needs to be done. The system needs to wait until a previous configuration change has been committed. One edge case that we haven't specified is when a leader is asked to delete itself.

---

 Leader $i$ removes a server $j$ (possibly itself) from the cluster.
$DeleteServer(i, j) \triangleq$
$\quad \wedge state[i] = Leader$
$\quad \wedge state[j] \in \{Follower, Candidate\}$
$\quad \wedge j \in GetConfig(i)$

$\wedge\, j \neq i$  *TODO*: A leader cannot remove itself.
$\wedge\, Send([mtype \mapsto CheckOldConfig,$
$\qquad mterm \mapsto currentTerm[i],$
$\qquad madd \mapsto \text{FALSE},$
$\qquad mserver \mapsto j,$
$\qquad msource \mapsto i,$
$\qquad mdest \mapsto i])$
$\wedge\, \text{UNCHANGED } \langle serverVars,\ candidateVars,\ leaderVars,\ logVars\rangle$

## 3.5  Modifying helper functions

### 3.5.1  Quorum

With static configurations, the quorum remains constant throughout execution. However, with dynamically changing configurations, a quorum is specific to each server's current view of the system, so we have added a parameter to the `Quorum` helper function definition so each server can compute a quorum for it's current configuration.

The set of all quorums for a server configuration.

This just calculates simple majorities, but the only

important property is that every quorum overlaps with every other.
$$Quorum(config) \triangleq \{i \in \text{SUBSET } (config) : Cardinality(i) * 2 > Cardinality(config)\}$$

### 3.5.2  Getting a server's configuration

Servers immediately start using configuration entries as they are appended to their logs, before they're committed. If a server's log has no configuration entries, the initial set of servers is used. We introduce the following helper functions `GetMaxConfigIndex` and `GetConfig` because many portions of the handlers and state transition functions require the server's configuration.

Return the index of the latest configuration in server $i$'s *log*.
$GetMaxConfigIndex(i) \triangleq$
$\quad \text{LET } configIndexes \triangleq \{index \in 1 .. Len(log[i]) : log[i][index].type = ConfigEntry\}$
$\quad \text{IN } \quad \text{IF } configIndexes = \{\} \text{ THEN } 0$
$\qquad\qquad \text{ELSE } \quad Max(configIndexes)$

Return the configuration of teh latest configuration in server $i$'s *log*.
$GetConfig(i) \triangleq$
$\quad \text{IF } GetMaxConfigIndex(i) = 0 \text{ THEN } InitServer$
$\quad \text{ELSE } \quad log[i][GetMaxConfigIndex(i)].value$

## 3.6 Handlers for configuration changes

We have introduced the following handlers for the new messages in the system.

### 3.6.1 Handling `CatchupRequest` messages

When a detached server receives this message, it should first check if the message is still valid, by checking `mterm` in the message. If this agrees, the server will appropriately overwrite and/or append the new entries (`mentries`) to it's log and respond to the leader indicating the current log position and that it has one less round to complete.

---

Detached server $i$ receives a *CatchupRequest* from leader $j$.

$HandleCatchupRequest(i, j, m) \triangleq$
  $\lor \ \land m.mterm < currentTerm[i]$
    $\land Reply([mtype \ \mapsto CatchupResponse,$
            $mterm \mapsto currentTerm[i],$
            $msuccess \mapsto \text{FALSE},$
            $mmatchIndex \mapsto 0,$
            $msource \mapsto i,$
            $mdest \mapsto j,$
            $mroundsLeft \mapsto 0],$
            $m)$
    $\land \text{UNCHANGED} \ \langle serverVars, candidateVars,$
          $leaderVars, logVars\rangle$
  $\lor \ \land m.mterm \geq currentTerm[i]$
    $\land currentTerm' = [currentTerm \ \text{EXCEPT} \ ![i] = m.mterm]$
    $\land log' = [log \ \text{EXCEPT} \ ![i] = SubSeq(log[i], 1, m.mlogLen) \circ m.mentries]$
    $\land Reply([mtype \ \mapsto CatchupResponse,$
            $mterm \mapsto currentTerm[i],$
            $msuccess \mapsto \text{TRUE},$
            $mmatchIndex \mapsto Len(log[i]),$
            $msource \mapsto i,$
            $mdest \mapsto j,$
            $mroundsLeft \mapsto m.mrounds - 1],$
            $m)$
    $\land \text{UNCHANGED} \ \langle state, votedFor, candidateVars, leaderVars,$
          $commitIndex\rangle$

---

### 3.6.2 Handling `CatchupResponse` messages

When a leader receives the `CatchupResponse` message, it checks if the server indicated it was successful in `msuccess`, then makes sure the `mmatchIndex` is correctly set. If so, it will send another request to the server with new log entries to catch up if there are still rounds

remaining. Otherwise, it will send a message to itself to wait until any uncommitted entries in it's log have been committed.

---

Leader $i$ receives a *CatchupResponse* from detached server $j$.
$HandleCatchupResponse(i, j, m)$ $\triangleq$

    A real system checks for progress every timeout interval.

    Assume that if this response is called, the new server

    has made progress.

    $\wedge$ $\vee$ $\wedge$ $m.msuccess$

         $\wedge$ $\vee$ $\wedge$ $m.mmatchIndex \neq commitIndex[i]$

             $\wedge$ $m.mmatchIndex \neq matchIndex[i][j]$

           $\vee$ $m.mmatchIndex = commitIndex[i]$

        $\wedge$ $state[i] = Leader$

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

        $\wedge$ $j \notin GetConfig(i)$

        $\wedge$ $nextIndex' = [nextIndex$ EXCEPT $![i][j] = m.mmatchIndex + 1]$

        $\wedge$ $matchIndex' = [matchIndex$ EXCEPT $![i][j] = m.mmatchIndex]$

        $\wedge$ $\vee$ $\wedge$ $m.mroundsLeft \neq 0$

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

                    $mterm \mapsto currentTerm[i],$

                    $mentries \mapsto SubSeq(log[i],$

                                  $nextIndex[i][j],$

                                  $commitIndex[i]),$

                    $mLogLen \mapsto nextIndex[i][j] - 1,$

                    $msource \mapsto i,$

                    $mdest \mapsto j,$

                    $mrounds \mapsto m.mroundsLeft],$

                    $m)$

      $\vee$ $\wedge$ $m.mroundsLeft = 0$

           A real system makes sure the final call to this handler is

           received after a timeout interval.

           We assume that if a timeout happened, the message

           has already been dropped.

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

                  $mterm \mapsto currentTerm[i],$

                  $madd \mapsto$ TRUE,

                  $mserver \mapsto j,$

                  $msource \mapsto i,$

                  $mdest \mapsto i], m)$

      $\wedge$ UNCHANGED $\langle elections \rangle$

    $\vee$ $\wedge$ $\vee$ $\neg m.msuccess$

         $\vee$ $\wedge$ $\vee$ $m.mmatchIndex = commitIndex[i]$

             $\vee$ $m.mmatchIndex = matchIndex[i][j]$

           $\wedge$ $m.mmatchIndex \neq commitIndex[i]$

         $\vee$ $state[i] \neq Leader$

$$\vee \ m.mterm \neq currentTerm[i]$$
$$\vee \ j \in GetConfig(i)$$
$$\wedge \ Discard(m)$$
$$\wedge \ \text{UNCHANGED} \ \langle leaderVars \rangle$$
$$\wedge \ \text{UNCHANGED} \ \langle serverVars, \ candidateVars, \ logVars \rangle$$

---

### 3.6.3  Handling `CheckOldConfig` messages

This handler causes the leader to wait until an uncommitted configuration is committed before adding a new entry. This is used both for adding and removing servers. If there is still an uncommitted entry, the leader will send itself another message to check again in the future. In a real system, this could be implemented by using a background thread on the server that sleeps and periodically checks, but this is nontrivial to model in the TLA+ spec and is equivalent to sending itself a message, even though the message can be duplicated or dropped.

---

Leader $i$ receives a *CheckOldConfig* message.
$HandleCheckOldConfig(i, m) \triangleq$
$\quad \vee \ \wedge state[i] \neq Leader \vee m.mterm = currentTerm[i]$
$\quad\quad \wedge Discard(m)$
$\quad\quad \wedge \text{UNCHANGED} \ \langle serverVars, \ candidateVars, \ leaderVars, \ logVars \rangle$
$\quad \vee \ \wedge state[i] = Leader \wedge m.mterm = currentTerm[i]$
$\quad\quad \wedge \ \vee \ \wedge GetMaxConfigIndex(i) \leq commitIndex[i]$
$\quad\quad\quad \wedge \text{LET} \ newConfig \ \triangleq \ \text{IF} \ m.madd \ \text{THEN} \ \text{UNION} \ \{GetConfig(i), \{m.mserver\}\}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \text{ELSE} \ \ GetConfig(i) \setminus \{m.mserver\}$
$\quad\quad\quad\quad\quad newEntry \ \triangleq \ [term \mapsto currentTerm[i], \ type \mapsto ConfigEntry, \ value \mapsto newConfig]$
$\quad\quad\quad\quad\quad newLog \ \triangleq \ Append(log[i], \ newEntry)$
$\quad\quad\quad\quad \text{IN} \ \ \ log' = [log \ \text{EXCEPT} \ ![i] = newLog]$
$\quad\quad\quad \wedge Discard(m)$
$\quad\quad\quad \wedge \text{UNCHANGED} \ \langle commitIndex \rangle$
$\quad\quad \vee \ \wedge GetMaxConfigIndex(i) > commitIndex[i]$
$\quad\quad\quad \wedge Reply([mtype \ \mapsto CheckOldConfig,$
$\quad\quad\quad\quad\quad\quad mterm \mapsto currentTerm[i],$
$\quad\quad\quad\quad\quad\quad madd \mapsto m.madd,$
$\quad\quad\quad\quad\quad\quad mserver \mapsto m.mserver,$
$\quad\quad\quad\quad\quad\quad msource \mapsto i,$
$\quad\quad\quad\quad\quad\quad mdest \mapsto i],$
$\quad\quad\quad\quad\quad\quad m)$
$\quad\quad\quad \wedge \text{UNCHANGED} \ \langle logVars \rangle$
$\quad\quad \wedge \text{UNCHANGED} \ \langle serverVars, \ candidateVars, \ leaderVars \rangle$

---

## 3.7 Mitigating effects of disruptive servers

Configuration changes can servers that have been removed to cause suboptimal (but still correct) system performance, as illustrated in Figure 3.

By studying our new specification, we have added a slight modification to the Raft algorithm to lessen the impacts disruptive servers can have: Servers can only timeout if they are in their own configuration.
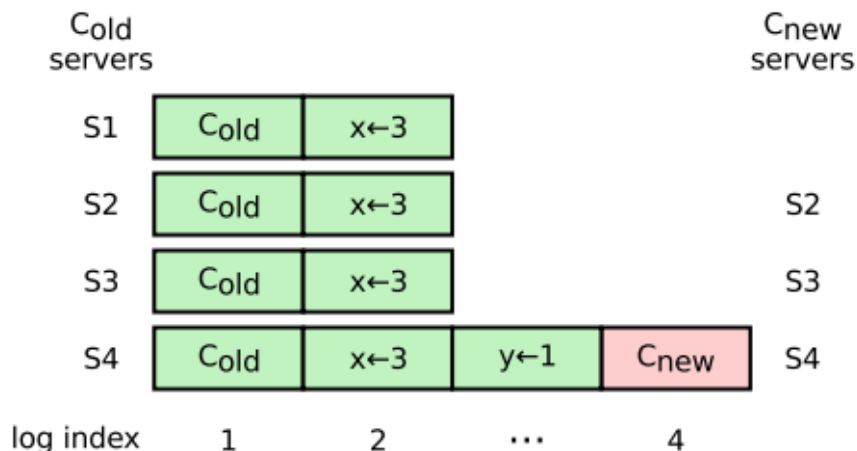


Figure 3: An example of how a server can be disruptive even before the $C_{\text{new}}$ log entry has been committed. The figure shows the removal of S1 from a four-server cluster. S4 is leader of the new cluster and has created the $C_{\text{new}}$ entry in its log, but it hasn't yet replicated that entry. Even before $C_{\text{new}}$ is committed, S1 can time out, increment its term, and send this larger term number to the new cluster, forcing S4 to step down. **Figure and description copied from Figure 4.7 of [Ong14] and included here for completeness.**

## 3.8 Model checking the specification

We have used the TLC model checker to validate simple cases of our modified specification. We created invariants that we knew would be broken so that we could obtain a traceback of the operations and messages that caused the point to be reached. One example is that a server that's not in the initial configuration eventually receives log entries because it has been added to the cluster.

# 4 Proofs

## 4.1 Safety: There is never more than one leader.

**Lemma 1.** *Let $n \geq 2$, $c_1 = \{1, \ldots, n\}$, $c_2 = \{1, \ldots, n-1\}$. If $s \in Quorum(c_1)$, $t \in Quorum(c_2)$, then $s \cap t \neq \emptyset$.*

*Proof.*

$$|s| \geq \left\lfloor \frac{n}{2} \right\rfloor + 1$$

$$|t| \geq \left\lfloor \frac{n-1}{2} \right\rfloor + 1$$

$$|s| + |t| \geq n + 1$$

Since there are only $n$ unique elements in $c_1 \cup c_2$, $s \cap t \neq \emptyset$. □

**Lemma 2.** *Let $n \geq 1$, $c_1 = \{1, \ldots, n\}$, $c_2 = \{1, \ldots, n{+}1\}$. If $s \in Quorum(c_1)$, $t \in Quorum(c_2)$, then $s \cap t \neq \emptyset$.*

*Proof.*

$$|s| \geq \left\lfloor \frac{n}{2} \right\rfloor + 1$$

$$|t| \geq \left\lfloor \frac{n+1}{2} \right\rfloor + 1$$

$$|s| + |t| \geq n + 2$$

Since there are only $n + 1$ unique elements in $c_1 \cup c_2$, $s \cap t \neq \emptyset$. □

**Lemma 3.** *Let $n \geq 1$, $c_1 = \{1, \ldots, n{-}1\}$, $c_2 = \{1, \ldots, n{+}1\}$. If $s \in Quorum(c_1)$, $t \in Quorum(c_2)$, then $s \cap t \neq \emptyset$.*

*Proof.*

$$|s| \geq \left\lfloor \frac{n-1}{2} \right\rfloor + 1$$

$$|t| \geq \left\lfloor \frac{n+1}{2} \right\rfloor + 1$$

$$|s| + |t| \geq n + 2$$

Since there are only $n + 1$ unique elements in $c_1 \cup c_2$, $s \cap t \neq \emptyset$. □

**Lemma 4.** *A quorum cannot be formed based on a stale config (i.e. a config that is before the latest committed config)*

*Proof.* Let $C_{\text{latest}}$ be the latest committed config and $C_{\text{latest}-1}$ be the config that is committed right before $C_{\text{latest}}$.

Suppose $C_{\text{latest}} = \{1, \ldots, n\}$. Then, $C_{\text{latest}-1}$ can either be $\{1, \ldots, n{-}1\}$ or $\{1, \ldots, n{+}1\}$. For simplicity, assume the last server is the one that changes.

Since $C_{\text{latest}}$ is committed, at least $\lfloor n/2 \rfloor + 1$ servers have $C_{\text{latest}}$ in their logs.

- **Case 1.** $C_{\text{latest}-1} = \{1, \ldots, n-1\}$. In order to form a quorum based on $C_{\text{latest}-1}$, it requires at least $\left\lfloor \frac{n-1}{2} \right\rfloor + 1$ votes.

  However, any server with $C_{\text{latest}}$ in its log won't vote yes because of the "Election Restriction" (§3.6.1 in [Ong14]) that "the voter denies its vote if its own log is more up-to-date than that of the candidate."

13

Therefore, it can only get at most $n - \lfloor n/2 \rfloor - 1$ votes.

Since

$$\left(n - \left\lfloor \frac{n}{2} \right\rfloor - 1\right) - \left(\left\lfloor \frac{n-1}{2} \right\rfloor - 1\right) = -1 < 0,$$

it can never get enough votes to form a quorum based on $C_{\text{latest}-1}$.

- **Case 2.** $C_{\text{latest}-1} = \{1, \ldots, n+1\}$. Similar argument as in **Case 1**.

Therefore, as long as $C_{\text{latest}}$ is committed, a quorum cannot be formed based on $C_{\text{latest}-1}$. Induction can show that any config prior to $C_{\text{latest}}$ cannot be the basis to form a quorum.

$\square$

**Lemma 5.** *Let $C_{\text{latest}}$ be the latest committed config. Let $C_{\text{new}}$ be any uncommitted config in the system, suppose $C_{\text{latest}} = \{1, \ldots, n\}$. Then, $C_{\text{new}}$ is either $\{1, \ldots, n-1\}$ or $\{1, \ldots, n+1\}$. For simplicity, assume the last server is the one that changes.*

*Proof.* By Lemma 4, since any stale config cannot be the basis of a quorum, any leader before a newer config gets committed in the system must have $C_{\text{latest}}$ in its log. Since in HandleCheckOldConfig, we require $GetMaxConfigIndex(i) \le commitIndex(i)$ to hold before the leader can append any newer config to its log, $C_{\text{new}}$ can only be "one step" away from $C_{\text{latest}}$. $\square$

**Theorem 1.** *There is at most one leader per term. This is the "Election Safety" property in Figure 3.2 and is proved for statically sized configurations in Lemma 2 of B.3 of [Ong14].*

$$\forall e, f \in elections$$

$$e.eterm = f.eterm \Rightarrow e.eleader = f.eleader$$

*Proof.* By Lemma 4 and Lemma 5, there can only be 3 possible configurations in the system at a time to form quorums:

$$C_{\text{latest}} = \{1, \ldots, n\}$$

$$C_{\text{new}+} = \{1, \ldots, n+1\}$$

$$C_{\text{new}-} = \{1, \ldots, n-1\}$$

For simplicity, assume the last server is the one that changes. Also note that if $n \ge 2$, all 3 are possible. If $n = 1$, only $C_{\text{latest}}$ and $C_{\text{new}+}$ are possible.

- **Case 1.** $e.evotes, f.evotes \in Quorum(C_{\text{latest}})$.

  Because any two quorums of a config overlap, $e.evotes \cap f.evotes \ne \emptyset$. Suppose $s \in (e.evotes \cap f.evotes)$. In HandleRequestVoteRequest,

  $$grant \triangleq \land m.mterm = currentTerm[i] \ (1)$$
  $$\land logOk$$
  $$\land votedFor[i] \in \{Nil, j\} \ (3)$$

  Properties (1) and (3) guarantee that a server can only vote for at most one server per term.

  Since $s \in e.evotes$ and $s \in f.evotes$, $e.eleader = f.eleader$.

14

- **Case 2.** $e.evotes, f.evotes \in Quorum(C_{new+})$. Similar proof to **Case 1**.

- **Case 3.** $e.evotes, f.evotes \in Quorum(C_{new-})$. Similar proof to **Case 1**.

- **Case 4.** $e.evotes \in Quorum(C_{latest})$, $f.evotes \in Quorum(C_{new+})$.

  By Lemma 2, $e.evotes \cap f.evotes \neq \emptyset$. Afterwards, similar proof to **Case 1**.

- **Case 5.** $e.evotes \in Quorum(C_{latest})$, $f.evotes \in Quorum(C_{new-})$.

  By Lemma 1, $e.evotes \cap f.evotes \neq \emptyset$. Afterwards, similar proof to **Case 1**.

- **Case 6.** $e.evotes \in Quorum(C_{new+})$, $f.evotes \in Quorum(C_{new-})$.

  By Lemma 3, $e.evotes \cap f.evotes \neq \emptyset$. Afterwards, similar proof to **Case 1**.

Therefore, there is at most one leader per term. □

## 4.2 Proof Sketch for Availability: A leader can be elected in the future

One availability property of the system is that a leader is able to be elected in some future state from any state. Our proof sketch is to choose a server that has the most updated log. Then, this server can time out and cause a quorum of it's configuration to vote for it, which will always be able to happen because servers will vote if a candidate's log is up-to-date and the term is greater than theirs.

Other servers can also time out while this server times out. It is not harmful for another server to receive a majority of the votes and become leader, nor is a split vote harmful, since the randomized timeouts will not collide in future elections in practice.

# 5 Broken Raft?

## 5.1 Cluster membership changes

We present two possible edge cases during cluster membership changes that illustrate a possible area where Raft's description might be inconsistent. We could be misinterpreting the wording in [OO14, Ong14] and plan to send these cases to the author.

### 5.1.1 New servers need to vote for availability

Consider the following initial cluster, where $s_1$ is the leader, represented with the $^*$ and the log of each server is shown on the right. Note the log is a 3-tuple of the term it was appended, the type (configuration or value), and the contents.

$$
\begin{array}{ll}
s_1^* & (1, \text{config}, \{1,2,3\}) \\
s_2 & (1, \text{config}, \{1,2,3\}) \\
s_3 & (1, \text{config}, \{1,2,3\}) \\
s_4 &
\end{array}
$$

$s_1$ gets a request to add $s_4$, so catches up $s_4$ with the config entry.

$$
\begin{array}{ll}
s_1^* & (1,\ \text{config},\ \{1,2,3\}) \\
s_2 & (1,\ \text{config},\ \{1,2,3\}) \\
s_3 & (1,\ \text{config},\ \{1,2,3\}) \\
s_4 & (1,\ \text{config},\ \{1,2,3\})
\end{array}
$$

$s_1$ then appends a new config to its log to add $s_4$.

$$
\begin{array}{ll}
s_1^* & (1,\ \text{config},\ \{1,2,3\}),\ (1,\ \text{config},\ \{1,2,3,4\}) \\
s_2 & (1,\ \text{config},\ \{1,2,3\}) \\
s_3 & (1,\ \text{config},\ \{1,2,3\}) \\
s_4 & (1,\ \text{config},\ \{1,2,3\})
\end{array}
$$

$s_3$ dies and $s_1$ replicates the new config to $s_2$.

$$
\begin{array}{ll}
s_1^* & (1,\ \text{config},\ \{1,2,3\}),\ (1,\ \text{config},\ \{1,2,3,4\}) \\
s_2 & (1,\ \text{config},\ \{1,2,3\}),\ (1,\ \text{config},\ \{1,2,3,4\}) \\
s_3 & (1,\ \text{config},\ \{1,2,3\}) \\
s_4 & (1,\ \text{config},\ \{1,2,3\})
\end{array}
$$

$s_2$ times out and starts an election and $s_1$ steps down. In this case, both $s_1$ and $s_2$ need $s_4$'s vote to become the leader. Otherwise the system won't have a leader and is thus non-available.

### 5.1.2 New members voting causes inconsistencies

Consider the following situation with 4 initial servers and $s_5$ is added.

Use $s^{*n}$ to denote a server being leader and $s^{Tn}$ to denote a server timing out, both in term $n$.

$$
\begin{array}{ll}
s_1^{*1} & (1,\ \text{config},\ \{1,2,3,4\}) \\
s_2 & (1,\ \text{config},\ \{1,2,3,4\}) \\
s_3 & (1,\ \text{config},\ \{1,2,3,4\}) \\
s_4 & (1,\ \text{config},\ \{1,2,3,4\}) \\
s_5 &
\end{array}
$$

$s_1$ catches up $s_5$.

$$
\begin{array}{ll}
s_1^{*1} & (1,\ \text{config},\ \{1,2,3,4\}) \\
s_2 & (1,\ \text{config},\ \{1,2,3,4\}) \\
s_3 & (1,\ \text{config},\ \{1,2,3,4\}) \\
s_4 & (1,\ \text{config},\ \{1,2,3,4\}) \\
s_5 & (1,\ \text{config},\ \{1,2,3,4\})
\end{array}
$$

$s_1$ appends new config.

$$s_1^{*1} \quad (1, \text{config}, \{1,2,3,4\}), (1,\text{config},\{1,2,3,4,5\})$$
$$s_2 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_3 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_4 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_5 \quad (1, \text{config}, \{1,2,3,4\})$$

$s_1$ replicates new config to $s_5$.

$$s_1^{*1} \quad (1, \text{config}, \{1,2,3,4\}), (1,\text{config},\{1,2,3,4,5\})$$
$$s_2 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_3 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_4 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_5 \quad (1, \text{config}, \{1,2,3,4\}), (1,\text{config},\{1,2,3,4,5\})$$

$s_1$ dies temporarily.

$$s_1^{D*1} \quad (1, \text{config}, \{1,2,3,4\}), (1,\text{config},\{1,2,3,4,5\})$$
$$s_2 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_3 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_4 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_5 \quad (1, \text{config}, \{1,2,3,4\}), (1,\text{config},\{1,2,3,4,5\})$$

$s_2$ times out and starts an election.

$$s_1^{D*1} \quad (1, \text{config}, \{1,2,3,4\}), (1,\text{config},\{1,2,3,4,5\})$$
$$s_2^{T2} \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_3 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_4 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_5 \quad (1, \text{config}, \{1,2,3,4\}), (1,\text{config},\{1,2,3,4,5\})$$

$s_2$, $s_3$, $s_4$ vote for $s_2$. $s_5$ rejects. $s_2$ becomes leader.

$$s_1^{D*1} \quad (1, \text{config}, \{1,2,3,4\}), (1,\text{config},\{1,2,3,4,5\})$$
$$s_2^{*2} \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_3 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_4 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_5 \quad (1, \text{config}, \{1,2,3,4\}), (1,\text{config},\{1,2,3,4,5\})$$

$s_2$ appends a new config to its log.

$$s_1^{D*1} \quad (1, \text{config}, \{1,2,3,4\}), (1,\text{config},\{1,2,3,4,5\})$$
$$s_2^{*2} \quad (1, \text{config}, \{1,2,3,4\}), (2, \text{config}, \{2,3,4\})$$
$$s_3 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_4 \quad (1, \text{config}, \{1,2,3,4\})$$
$$s_5 \quad (1, \text{config}, \{1,2,3,4\}), (1,\text{config},\{1,2,3,4,5\})$$

$s_2$ replicates new config to $s_3$ and **is committed**!

$$
\begin{array}{ll}
s_1^{D*1} & (1,\text{ config, }\{1,2,3,4\}),\ (1,\text{config},\{1,2,3,4,5\}) \\
s_2^{*2} & (1,\text{ config, }\{1,2,3,4\}),\ (2,\text{ config, }\{2,3,4\}) \\
s_3 & (1,\text{ config, }\{1,2,3,4\}),\ (2,\text{ config, }\{2,3,4\}) \\
s_4 & (1,\text{ config, }\{1,2,3,4\}) \\
s_5 & (1,\text{ config, }\{1,2,3,4\}),\ (1,\text{config},\{1,2,3,4,5\})
\end{array}
$$

$s_1$ comes backs alive and times out and starts an election.

$$
\begin{array}{ll}
s_1^{T3} & (1,\text{ config, }\{1,2,3,4\}),\ (1,\text{config},\{1,2,3,4,5\}) \\
s_2^{*2} & (1,\text{ config, }\{1,2,3,4\}),\ (2,\text{ config, }\{2,3,4\}) \\
s_3 & (1,\text{ config, }\{1,2,3,4\}),\ (2,\text{ config, }\{2,3,4\}) \\
s_4 & (1,\text{ config, }\{1,2,3,4\}) \\
s_5 & (1,\text{ config, }\{1,2,3,4\}),\ (1,\text{config},\{1,2,3,4,5\})
\end{array}
$$

If $s_5$ can vote, then $s_1$ can receive $s_1$, $s_4$, and $s_5$'s votes and become the new leader. Then $s_1$ will try to replicate its log to everyone, including $s_2$ and $s_3$, which will conflict and overwrite the already committed entry $(2,\text{ config, }\{2,3,4\})$ with an older uncommitted entry. This breaks the leader completeness property presented in Figure 3.2 of [Ong14]: "If a log entry is committed in a given term, then that entry will be present in the logs of the leaders for all higher-numbered terms."

# 6    Conclusion and Future Work

We have presented a formal specification for Raft cluster membership changes and have proved that properties of the cluster are preserved during these changes. Future work involves further validating our modifications to the specification and modeling more invariants and properties of Raft. An interesting direction could be to study other formal verifications of Raft, such as Verdi's case study of Raft in PLDI 2015 [WWP+15].

# References

[AS87]      Bowen Alpern and Fred B Schneider. Recognizing safety and liveness. *Distributed computing*, 2(3):117–126, 1987.

[CDLM08]  Kaustuv C Chaudhuri, Damien Doligez, Leslie Lamport, and Stephan Merz. A tla+ proof system. *arXiv preprint arXiv:0811.1914*, 2008.

[Lam94]    Leslie Lamport. The temporal logic of actions. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 16(3):872–923, 1994.

[Lam00]    Leslie Lamport. A summary of TLA+. 2000.

[Lam02]    Leslie Lamport. *Specifying systems: The TLA+ language and tools for hardware and software engineers*. Addison-Wesley Longman Publishing Co., Inc., 2002.

[Ong14]    Diego Ongaro. *Consensus: Bridging theory and practice*. PhD thesis, Stanford University, 2014.

[OO14]     Diego Ongaro and John Ousterhout. In search of an understandable consensus algorithm. In *Proc. USENIX Annual Technical Conference*, pages 305–320, 2014.

[raf]       `https://raftconsensus.github.io/`.

[sec]       `http://thesecretlivesofdata.com/raft/`.

[WWP+15] James R Wilcox, Doug Woos, Pavel Panchekha, Zachary Tatlock, Xi Wang, Michael D Ernst, and Thomas Anderson. Verdi: A framework for implementing and formally verifying distributed systems. *PLDI*, 2015.

[YML99]    Yuan Yu, Panagiotis Manolios, and Leslie Lamport. Model checking tla+ specifications. In *Correct Hardware Design and Verification Methods*, pages 54–66. Springer, 1999.

# A    Original TLA+ Specification

Starts on next page.

———— MODULE *raft_orig* ————

This is the formal specification for the Raft consensus algorithm.

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*

⊢───────────────────────────────────────────

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*

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*

⊢───────────────────────────────────────────

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

The server's term number.

1

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).

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

---

  Helpers

2

The set of all quorums. This just calculates simple majorities, but the only
important property is that every quorum overlaps with every other.
$$Quorum \triangleq \{i \in \text{SUBSET } (Server) : Cardinality(i) * 2 > Cardinality(Server)\}$$

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

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

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

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

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

3

$InitCandidateVars \;\stackrel{\Delta}{=}\; \wedge\; votesResponded = [i \in Server \mapsto \{\}]$
$\qquad\qquad\qquad\qquad\;\; \wedge\; votesGranted \;\;\;\; = [i \in Server \mapsto \{\}]$

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.

$InitLeaderVars \;\stackrel{\Delta}{=}\; \wedge\; nextIndex \;\;\; = [i \in Server \mapsto [j \in Server \mapsto 1]]$
$\qquad\qquad\qquad\qquad \wedge\; matchIndex = [i \in Server \mapsto [j \in Server \mapsto 0]]$
$InitLogVars \;\stackrel{\Delta}{=}\; \wedge\; log \qquad\qquad = [i \in Server \mapsto \langle\rangle]$
$\qquad\qquad\qquad\quad \wedge\; commitIndex \;\; = [i \in Server \mapsto 0]$
$Init \;\stackrel{\Delta}{=}\; \wedge\; messages = [m \in \{\} \mapsto 0]$
$\qquad\quad \wedge\; InitHistoryVars$
$\qquad\quad \wedge\; InitServerVars$
$\qquad\quad \wedge\; InitCandidateVars$
$\qquad\quad \wedge\; InitLeaderVars$
$\qquad\quad \wedge\; InitLogVars$

---

Define state transitions

Server $i$ restarts from stable storage.

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

$Restart(i) \;\stackrel{\Delta}{=}$
$\quad \wedge\;\; state' \qquad\qquad\quad = [state \text{ EXCEPT } ![i] = Follower]$
$\quad \wedge\;\; votesResponded' = [votesResponded \text{ EXCEPT } ![i] = \{\}]$
$\quad \wedge\;\; votesGranted' \;\;\; = [votesGranted \text{ EXCEPT } ![i] = \{\}]$
$\quad \wedge\;\; voterLog' \qquad\;\; = [voterLog \text{ EXCEPT } ![i] = [j \in \{\} \mapsto \langle\rangle]]$
$\quad \wedge\;\; nextIndex' \qquad\;\; = [nextIndex \text{ EXCEPT } ![i] = [j \in Server \mapsto 1]]$
$\quad \wedge\;\; matchIndex' \qquad = [matchIndex \text{ EXCEPT } ![i] = [j \in Server \mapsto 0]]$
$\quad \wedge\;\; commitIndex' \;\;\; = [commitIndex \text{ EXCEPT } ![i] = 0]$
$\quad \wedge\;\; \text{UNCHANGED } \langle messages,\, currentTerm,\, votedFor,\, log,\, elections\rangle$

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

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

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

$RequestVote(i,\, j) \;\stackrel{\Delta}{=}$
$\quad \wedge\; state[i] = Candidate$

4

$\quad \wedge j \notin votesResponded[i]$

$\quad \wedge Send([mtype \qquad\qquad \mapsto RequestVoteRequest,$

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

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

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

$\qquad\qquad msource \qquad\quad \mapsto i,$

$\qquad\qquad mdest \qquad\qquad \mapsto j])$

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

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.

$AppendEntries(i, j) \triangleq$

$\quad \wedge i \neq j$

$\quad \wedge state[i] = Leader$

$\quad \wedge \text{LET } prevLogIndex \triangleq nextIndex[i][j] - 1$

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

$\qquad\qquad\qquad\qquad\qquad\qquad log[i][prevLogIndex].term$

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

$\qquad\qquad\qquad\qquad\qquad 0$

$\qquad\quad$ Send up to 1 entry, constrained by the end of the *log*.

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

$\qquad\quad entries \triangleq SubSeq(log[i], nextIndex[i][j], lastEntry)$

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

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

$\qquad\qquad\quad mprevLogIndex \mapsto prevLogIndex,$

$\qquad\qquad\quad mprevLogTerm \mapsto prevLogTerm,$

$\qquad\qquad\quad mentries \qquad\qquad \mapsto entries,$

$\qquad\qquad\qquad$ *mlog* is used as a history variable for the proof.

$\qquad\qquad\qquad$ It would not exist in a real implementation.

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

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

$\qquad\qquad\quad msource \qquad\qquad \mapsto i,$

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

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

Candidate $i$ transitions to leader.

$BecomeLeader(i) \triangleq$

$\quad \wedge state[i] = Candidate$

$\quad \wedge votesGranted[i] \in Quorum$

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

$\quad \wedge nextIndex' \quad = [nextIndex \text{ EXCEPT } ![i] =$

$\qquad\qquad\qquad\qquad\qquad [j \in Server \mapsto Len(log[i]) + 1]]$

$\quad \wedge matchIndex' = [matchIndex \text{ EXCEPT } ![i] =$

$\qquad\qquad\qquad\qquad\qquad [j \in Server \mapsto 0]]$

5

$$\land\ elections' \quad = elections\ \cup$$
$$\{[eterm \qquad \mapsto currentTerm[i],$$
$$eleader \quad \mapsto i,$$
$$elog \qquad \mapsto log[i],$$
$$evotes \quad \mapsto votesGranted[i],$$
$$evoterLog \mapsto voterLog[i]]\}$$
$$\land\ \text{UNCHANGED}\ \langle messages,\ currentTerm,\ votedFor,\ candidateVars,\ logVars\rangle$$

Leader $i$ receives a client request to add $v$ to the *log*.
$$ClientRequest(i,\ v)\ \triangleq$$
$$\land\ state[i] = Leader$$
$$\land\ \text{LET}\ entry\ \triangleq\ [term\ \mapsto currentTerm[i],$$
$$value\ \mapsto v]$$
$$newLog\ \triangleq\ Append(log[i],\ entry)$$
$$\text{IN}\quad log' = [log\ \text{EXCEPT}\ ![i] = newLog]$$
$$\land\ \text{UNCHANGED}\ \langle messages,\ serverVars,\ candidateVars,$$
$$leaderVars,\ commitIndex\rangle$$

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.
$$AdvanceCommitIndex(i)\ \triangleq$$
$$\land\ state[i] = Leader$$
$$\land\ \text{LET}\quad \text{The set of servers that agree up through index.}$$
$$Agree(index)\ \triangleq\ \{i\} \cup \{k \in Server :$$
$$matchIndex[i][k] \geq index\}$$
$$\text{The maximum indexes for which a quorum agrees}$$
$$agreeIndexes\ \triangleq\ \{index \in 1\ ..\ Len(log[i]) :$$
$$Agree(index) \in Quorum\}$$
$$\text{New value for } commitIndex'[i]$$
$$newCommitIndex\ \triangleq$$
$$\text{IF}\ \land\ agreeIndexes \neq \{\}$$
$$\land\ log[i][Max(agreeIndexes)].term = currentTerm[i]$$
$$\text{THEN}$$
$$Max(agreeIndexes)$$
$$\text{ELSE}$$
$$commitIndex[i]$$
$$\text{IN}\quad commitIndex' = [commitIndex\ \text{EXCEPT}\ ![i] = newCommitIndex]$$
$$\land\ \text{UNCHANGED}\ \langle messages,\ serverVars,\ candidateVars,\ leaderVars,\ log\rangle$$

---

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

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

6

$m.mterm \leq currentTerm[i]$.
$HandleRequestVoteRequest(i, j, m) \triangleq$
  LET $logOk \triangleq \lor m.mlastLogTerm > LastTerm(log[i])$
           $\lor \land m.mlastLogTerm = LastTerm(log[i])$
              $\land m.mlastLogIndex \geq Len(log[i])$
    $grant \triangleq \land m.mterm = currentTerm[i]$
            $\land logOk$
            $\land votedFor[i] \in \{Nil, j\}$
  IN  $\land m.mterm \leq currentTerm[i]$
      $\land \lor grant \quad \land votedFor' = [votedFor \text{ EXCEPT } ![i] = j]$
        $\lor \neg grant \land \text{UNCHANGED } votedFor$
      $\land Reply([mtype \quad\quad \mapsto RequestVoteResponse,$
              $mterm \quad\quad \mapsto currentTerm[i],$
              $mvoteGranted \mapsto grant,$
                  $mlog$ is used just for the *elections* history variable for
                  the proof. It would not exist in a real implementation.
              $mlog \quad\quad\quad \mapsto log[i],$
              $msource \quad\quad \mapsto i,$
              $mdest \quad\quad \mapsto j],$
              $m)$
      $\land \text{UNCHANGED } \langle state, currentTerm, candidateVars, leaderVars, logVars \rangle$

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.
    $\land m.mterm = currentTerm[i]$
    $\land votesResponded' = [votesResponded \text{ EXCEPT } ![i] =$
                    $votesResponded[i] \cup \{j\}]$
    $\land \lor \land m.mvoteGranted$
        $\land votesGranted' = [votesGranted \text{ EXCEPT } ![i] =$
                      $votesGranted[i] \cup \{j\}]$
        $\land voterLog' = [voterLog \text{ EXCEPT } ![i] =$
                    $voterLog[i] @@ (j :> m.mlog)]$
      $\lor \land \neg m.mvoteGranted$
        $\land \text{UNCHANGED } \langle votesGranted, voterLog \rangle$
    $\land Discard(m)$
    $\land \text{UNCHANGED } \langle serverVars, votedFor, leaderVars, logVars \rangle$

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.
$HandleAppendEntriesRequest(i, j, m) \triangleq$

7

26

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

8

$$\land \text{UNCHANGED } \langle serverVars,\ logVars \rangle$$
$\lor$   conflict: remove 1 entry
$$\land m.mentries \neq \langle \rangle$$
$$\land Len(log[i]) \geq index$$
$$\land log[i][index].term \quad \neq m.mentries[1].term$$
$$\land \text{LET } new \stackrel{\Delta}{=} [index2 \in 1 \mathinner{.\,.} (Len(log[i]) - 1) \mapsto$$
$$log[i][index2]]$$
$$\text{IN} \quad log' = [log \text{ EXCEPT } ![i] = new]$$
$$\land \text{UNCHANGED } \langle serverVars,\ commitIndex,\ messages \rangle$$
$\lor$   no conflict: append entry
$$\land m.mentries \neq \langle \rangle$$
$$\land Len(log[i]) = m.mprevLogIndex$$
$$\land log' = [log \text{ EXCEPT } ![i] =$$
$$Append(log[i],\ m.mentries[1])]$$
$$\land \text{UNCHANGED } \langle serverVars,\ commitIndex,\ messages \rangle$$
$$\land \text{UNCHANGED } \langle candidateVars,\ leaderVars \rangle$$

Server $i$ receives an *AppendEntries* response from server $j$ with
$m.mterm = currentTerm[i]$.
$HandleAppendEntriesResponse(i,\ j,\ m) \stackrel{\Delta}{=}$
$$\land m.mterm = currentTerm[i]$$
$$\land \lor \land m.msuccess \text{ successful}$$
$$\land nextIndex' \quad = [nextIndex \text{ EXCEPT } ![i][j] = m.mmatchIndex + 1]$$
$$\land matchIndex' = [matchIndex \text{ EXCEPT } ![i][j] = m.mmatchIndex]$$
$$\lor \land \neg m.msuccess \text{ not successful}$$
$$\land nextIndex' = [nextIndex \text{ EXCEPT } ![i][j] =$$
$$Max(\{nextIndex[i][j] - 1,\ 1\})]$$
$$\land \text{UNCHANGED } \langle matchIndex \rangle$$
$$\land Discard(m)$$
$$\land \text{UNCHANGED } \langle serverVars,\ candidateVars,\ logVars,\ elections \rangle$$

Any *RPC* with a newer term causes the recipient to advance its term first.
$UpdateTerm(i,\ j,\ m) \stackrel{\Delta}{=}$
$$\land m.mterm > currentTerm[i]$$
$$\land currentTerm' \quad = [currentTerm \text{ EXCEPT } ![i] = m.mterm]$$
$$\land state' \qquad\quad = [state \qquad\quad \text{ EXCEPT } ![i] \quad = Follower]$$
$$\land votedFor' \qquad = [votedFor \qquad \text{ EXCEPT } ![i] \quad = Nil]$$
    messages is unchanged so $m$ can be processed further.
$$\land \text{UNCHANGED } \langle messages,\ candidateVars,\ leaderVars,\ logVars \rangle$$

Responses with stale terms are ignored.
$DropStaleResponse(i,\ j,\ m) \stackrel{\Delta}{=}$
$$\land m.mterm < currentTerm[i]$$
$$\land Discard(m)$$
$$\land \text{UNCHANGED } \langle serverVars,\ candidateVars,\ leaderVars,\ logVars \rangle$$

9

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

$DuplicateMessage(m) \triangleq$

 $\wedge Send(m)$

 $\wedge$ UNCHANGED $\langle serverVars, candidateVars, leaderVars, logVars \rangle$

The network drops a message

$DropMessage(m) \triangleq$

 $\wedge Discard(m)$

 $\wedge$ UNCHANGED $\langle serverVars, candidateVars, leaderVars, logVars \rangle$

---

Defines how the variables may transition.

$Next \triangleq \wedge \vee \exists i \in Server : Restart(i)$

     $\vee \exists i \in Server : Timeout(i)$

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

     $\vee \exists i \in Server : BecomeLeader(i)$

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

     $\vee \exists i \in Server : AdvanceCommitIndex(i)$

     $\vee \exists i, j \in Server : AppendEntries(i, j)$

     $\vee \exists m \in$ DOMAIN $messages : Receive(m)$

     $\vee \exists m \in$ DOMAIN $messages : DuplicateMessage(m)$

     $\vee \exists m \in$ DOMAIN $messages : DropMessage(m)$

    History variable that tracks every $log$ ever:

   $\wedge allLogs' = allLogs \cup \{log[i] : i \in Server\}$

10

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

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

\ ∗ *Changelog*:
\ ∗
\ ∗ $2014 - 12 - 02$:
\ ∗ $-$ Fix *AppendEntries* to only send one entry at a time, as originally
\ ∗ intended. Since *SubSeq* is inclusive, the upper bound of the range should
\ ∗ have been *nextIndex*, not $nextIndex + 1$. Thanks to *Igor Kovalenko* for
\ ∗ reporting the issue.
\ ∗ $-$ Change $matchIndex'$ to *matchIndex* (without the apostrophe) in
\ ∗ *AdvanceCommitIndex*. This apostrophe was not intentional and perhaps
\ ∗ confusing, though it makes no practical difference ($matchIndex'$ equals
\ ∗ *matchIndex*). Thanks to *Hugues Evrard* for reporting the issue.
\ ∗
\ ∗ $2014 - 07 - 06$:
\ ∗ $-$ Version from *PhD* dissertation

11

# B Our Modified TLA+ Specification

Starts on next page.

—— MODULE *raft* ——

This is the formal specification for the Raft consensus algorithm.

Original Copyright 2014 *Diego Ongaro*
Modifications for cluster membership changes by
  Brandon *Amos* and *Huanchen Zhang*, 2015

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 number of rounds to catch new servers up for.
Must be $\geq 1$.
CONSTANTS *NumRounds*

The initial and global set of servers.
CONSTANTS *InitServer*, *Server*

Log metadata to distinguish values from configuration changes.
CONSTANT *ValueEntry*, *ConfigEntry*

The set of values 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*,
              *CatchupRequest*, *CatchupResponse*,
              *CheckOldConfig*

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*

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

1

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*

---

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 *Server* 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*

2

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).

$vars \triangleq \langle messages,\ allLogs,\ serverVars,\ candidateVars,$
$\qquad\qquad leaderVars,\ logVars \rangle$

---

Helpers

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

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

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

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

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.

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

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.

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

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

Combination of *Send* and *Discard*

$Reply(response,\ request) \triangleq$

3

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

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

Return the index of the latest configuration in server $i$'s *log*.
$GetMaxConfigIndex(i) \triangleq$
 LET $configIndexes \triangleq \{index \in 1\,..\,Len(log[i]) : log[i][index].type = ConfigEntry\}$
 IN IF $configIndexes = \{\}$ THEN $0$
   ELSE $Max(configIndexes)$

Return the configuration of teh latest configuration in server $i$'s *log*.
$GetConfig(i) \triangleq$
 IF $GetMaxConfigIndex(i) = 0$ THEN $InitServer$
 ELSE $log[i][GetMaxConfigIndex(i)].value$

---

Define initial values for all variables

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

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

---

Define state transitions

Server $i$ restarts from stable storage.

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

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

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

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

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

$RequestVote(i, j) \triangleq$
$\quad \wedge\ state[i] = Candidate$
$\quad \wedge\ j \in (GetConfig(i) \setminus votesResponded[i])$
$\quad \wedge\ Send([mtype\qquad\quad \mapsto RequestVoteRequest,$
$\qquad\qquad mterm\qquad\qquad \mapsto currentTerm[i],$
$\qquad\qquad mlastLogTerm \mapsto LastTerm(log[i]),$
$\qquad\qquad mlastLogIndex \mapsto Len(log[i]),$
$\qquad\qquad msource\qquad\quad \mapsto i,$
$\qquad\qquad mdest\qquad\qquad \mapsto j])$
$\quad \wedge\ \text{UNCHANGED } \langle serverVars, candidateVars, leaderVars, logVars \rangle$

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.

$AppendEntries(i, j) \triangleq$
$\quad \wedge\ i \neq j$
$\quad \wedge\ state[i] = Leader$
$\quad \wedge\ j \in GetConfig(i)$
$\quad \wedge\ \text{LET } prevLogIndex \triangleq\ nextIndex[i][j] - 1$

5

$$prevLogTerm \triangleq \text{IF } prevLogIndex > 0 \text{ THEN}$$
$$log[i][prevLogIndex].term$$
$$\text{ELSE}$$
$$0$$

Send up to 1 entry, constrained by the end of the *log*.
$$lastEntry \triangleq Min(\{Len(log[i]), nextIndex[i][j]\})$$
$$entries \triangleq SubSeq(log[i], nextIndex[i][j], lastEntry)$$

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

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

$\qquad\qquad mprevLogIndex \ \mapsto prevLogIndex,$

$\qquad\qquad mprevLogTerm \ \ \mapsto prevLogTerm,$

$\qquad\qquad mentries \qquad\quad \mapsto entries,$

    *mlog* is used as a history variable for the proof.

    It would not exist in a real implementation.

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

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

$\qquad\qquad msource \qquad\quad \mapsto i,$

$\qquad\qquad mdest \qquad\qquad \mapsto j])$

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

Candidate $i$ transitions to leader.
$BecomeLeader(i) \triangleq$

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

$\quad \land votesGranted[i] \in Quorum(GetConfig(i))$

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

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

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

$\quad \land elections' \qquad = elections \cup$
$$\{[eterm \qquad \mapsto currentTerm[i],$$
$$eleader \quad \mapsto i,$$
$$elog \qquad \mapsto log[i],$$
$$evotes \qquad \mapsto votesGranted[i],$$
$$evoterLog \mapsto voterLog[i]]\}$$

$\quad \land \text{UNCHANGED } \langle messages, currentTerm, votedFor, candidateVars, logVars \rangle$

Leader $i$ receives a client request to add $v$ to the *log*.
$ClientRequest(i, v) \triangleq$

$\quad \land state[i] = Leader$

$\quad \land \text{LET } entry \triangleq [term \ \mapsto currentTerm[i],$
$$type \quad \mapsto ValueEntry,$$
$$value \mapsto v]$$
$$newLog \triangleq Append(log[i], entry)$$

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

6

$\wedge$ UNCHANGED $\langle messages, serverVars, candidateVars,$
$\qquad\qquad\qquad leaderVars, commitIndex\rangle$

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.
$AdvanceCommitIndex(i) \triangleq$
$\quad \wedge state[i] = Leader$
$\quad \wedge$ LET $\quad$ The set of servers that agree up through index.
$\qquad\qquad Agree(index) \triangleq \{i\} \cup \{k \in GetConfig(i) :$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad matchIndex[i][k] \geq index\}$
$\qquad\qquad$ The maximum indexes for which a quorum agrees
$\qquad\qquad agreeIndexes \triangleq \{index \in 1 .. Len(log[i]) :$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad Agree(index) \in Quorum(GetConfig(i))\}$
$\qquad\qquad$ New value for $commitIndex'[i]$
$\qquad\qquad newCommitIndex \triangleq$
$\qquad\qquad\quad$ IF $\wedge agreeIndexes \neq \{\}$
$\qquad\qquad\qquad\quad \wedge log[i][Max(agreeIndexes)].term = currentTerm[i]$
$\qquad\qquad\quad$ THEN
$\qquad\qquad\qquad\quad Max(agreeIndexes)$
$\qquad\qquad\quad$ ELSE
$\qquad\qquad\qquad\quad commitIndex[i]$
$\quad\quad$ IN $\quad commitIndex' = [commitIndex$ EXCEPT $![i] = newCommitIndex]$
$\quad \wedge$ UNCHANGED $\langle messages, serverVars, candidateVars, leaderVars, log\rangle$

Leader $i$ adds a new server $j$ to the cluster.
$AddNewServer(i, j) \triangleq$
$\quad \wedge state[i] = Leader$
$\quad \wedge j \notin GetConfig(i)$
$\quad \wedge currentTerm' = [currentTerm$ EXCEPT $![j] = 1]$
$\quad \wedge votedFor' = [votedFor$ EXCEPT $![j] = Nil]$
$\quad \wedge Send([mtype \mapsto CatchupRequest,$
$\qquad\qquad mterm \mapsto currentTerm[i],$
$\qquad\qquad mlogLen \mapsto matchIndex[i][j],$
$\qquad\qquad mentries \mapsto SubSeq(log[i], nextIndex[i][j], commitIndex[i]),$
$\qquad\qquad mcommitIndex \mapsto commitIndex[i],$
$\qquad\qquad msource \mapsto i,$
$\qquad\qquad mdest \mapsto j,$
$\qquad\qquad mrounds \mapsto NumRounds])$
$\quad \wedge$ UNCHANGED $\langle state, leaderVars, logVars, candidateVars\rangle$

Leader $i$ removes a server $j$ (possibly itself) from the cluster.
$DeleteServer(i, j) \triangleq$
$\quad \wedge state[i] = Leader$

7

$\land state[j] \in \{Follower,\ Candidate\}$
$\land j \in GetConfig(i)$
$\land j \neq i$  TODO: A leader cannot remove itself.
$\land Send([mtype\ \mapsto CheckOldConfig,$
$\qquad\qquad mterm \mapsto currentTerm[i],$
$\qquad\qquad madd \mapsto \text{FALSE},$
$\qquad\qquad mserver \mapsto j,$
$\qquad\qquad msource \mapsto i,$
$\qquad\qquad mdest \mapsto i])$
$\land \text{UNCHANGED}\ \langle serverVars,\ candidateVars,\ leaderVars,\ logVars \rangle$

---

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

Server $i$ receives a *RequestVote* request from server $j$ with
$m.mterm \leq currentTerm[i]$.
$HandleRequestVoteRequest(i,\ j,\ m)\ \triangleq$
$\quad \text{LET}\ logOk\ \triangleq\ \lor\ m.mlastLogTerm > LastTerm(log[i])$
$\qquad\qquad\qquad\qquad \lor\ \land\ m.mlastLogTerm = LastTerm(log[i])$
$\qquad\qquad\qquad\qquad\quad\ \land\ m.mlastLogIndex \geq Len(log[i])$
$\qquad\quad grant\ \triangleq\ \land\ m.mterm = currentTerm[i]$
$\qquad\qquad\qquad\quad\ \land\ logOk$
$\qquad\qquad\qquad\quad\ \land\ votedFor[i] \in \{Nil,\ j\}$
$\quad \text{IN}\quad \land\ m.mterm \leq currentTerm[i]$
$\qquad\qquad \land\ \lor\ grant\quad \land\ votedFor' = [votedFor\ \text{EXCEPT}\ ![i] = j]$
$\qquad\qquad\quad\ \ \lor\ \neg grant \land \text{UNCHANGED}\ votedFor$
$\qquad\qquad \land\ Reply([mtype\qquad\quad \mapsto RequestVoteResponse,$
$\qquad\qquad\qquad\qquad mterm\qquad\quad\ \mapsto currentTerm[i],$
$\qquad\qquad\qquad\qquad mvoteGranted \mapsto grant,$
$\qquad\qquad\qquad\qquad\quad mlog$ is used just for the *elections* history variable for
$\qquad\qquad\qquad\qquad\quad$ the proof. It would not exist in a real implementation.
$\qquad\qquad\qquad\qquad mlog\qquad\qquad \mapsto log[i],$
$\qquad\qquad\qquad\qquad msource\qquad\ \mapsto i,$
$\qquad\qquad\qquad\qquad mdest\qquad\quad\ \mapsto j],$
$\qquad\qquad\qquad\qquad m)$
$\qquad\qquad \land \text{UNCHANGED}\ \langle state,\ currentTerm,\ candidateVars,\ leaderVars,\ logVars \rangle$

Server $i$ receives a *RequestVote* response from server $j$ with
$m.mterm = currentTerm[i]$.
$HandleRequestVoteResponse(i,\ j,\ m)\ \triangleq$
$\qquad$ This tallies votes even when the current state is not *Candidate*, but
$\qquad$ they won't be looked at, so it doesn't matter.
$\qquad \land\ m.mterm = currentTerm[i]$
$\qquad \land\ votesResponded' = [votesResponded\ \text{EXCEPT}\ ![i] =$

8

$$votesResponded[i] \cup \{j\}]$$
$\wedge \ \vee \ \wedge \ m.mvoteGranted$
  $\wedge votesGranted' = [votesGranted \text{ EXCEPT } ![i] =$
  $$votesGranted[i] \cup \{j\}]$$
  $\wedge voterLog' = [voterLog \text{ EXCEPT } ![i] =$
  $$voterLog[i] @@ (j :> m.mlog)]$$
 $\vee \ \wedge \ \neg m.mvoteGranted$
  $\wedge \text{UNCHANGED } \langle votesGranted, \ voterLog \rangle$
$\wedge Discard(m)$
$\wedge \text{UNCHANGED } \langle serverVars, \ votedFor, \ leaderVars, \ logVars \rangle$

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.
$HandleAppendEntriesRequest(i, j, m) \ \triangleq$
  LET $logOk \ \triangleq \ \vee \ m.mprevLogIndex = 0$
    $\vee \ \wedge \ m.mprevLogIndex > 0$
     $\wedge \ m.mprevLogIndex \leq Len(log[i])$
     $\wedge \ m.mprevLogTerm = log[i][m.mprevLogIndex].term$
  IN  $\wedge \ m.mterm \leq currentTerm[i]$
    $\wedge \ \vee \ \wedge \ $ reject request
       $\vee \ m.mterm < currentTerm[i]$
       $\vee \ \wedge \ m.mterm = currentTerm[i]$
        $\wedge \ state[i] = Follower$
        $\wedge \ \neg logOk$
      $\wedge Reply([mtype \qquad\qquad \mapsto AppendEntriesResponse,$
        $mterm \qquad\qquad \mapsto currentTerm[i],$
        $msuccess \qquad\quad \mapsto \text{FALSE},$
        $mmatchIndex \quad\ \ \mapsto 0,$
        $msource \qquad\quad\ \mapsto i,$
        $mdest \qquad\qquad \mapsto j],$
        $m)$
      $\wedge \text{UNCHANGED } \langle serverVars, \ logVars \rangle$
     $\vee \ $ return to follower state
      $\wedge \ m.mterm = currentTerm[i]$
      $\wedge \ state[i] = Candidate$
      $\wedge \ state' = [state \text{ EXCEPT } ![i] = Follower]$
      $\wedge \text{UNCHANGED } \langle currentTerm, \ votedFor, \ logVars, \ messages \rangle$
     $\vee \ $ accept request
      $\wedge \ m.mterm = currentTerm[i]$
      $\wedge \ state[i] = Follower$
      $\wedge \ logOk$
      $\wedge \text{LET } index \ \triangleq \ m.mprevLogIndex + 1$
        IN  $\vee \ $ already done with request

9

40

$$\land \lor m.mentries = \langle\rangle$$
$$\lor \land Len(log[i]) \geq index$$
$$\land log[i][index].term = m.mentries[1].term$$

This could make our *commitIndex* decrease (for
example if we process an old, duplicated request),
but that doesn't really affect anything.

$$\land commitIndex' = [commitIndex \text{ EXCEPT } ![i] = m.mcommitIndex]$$
$$\land Reply([mtype \qquad\qquad \mapsto AppendEntriesResponse,$$
$$mterm \qquad\qquad \mapsto currentTerm[i],$$
$$msuccess \qquad\qquad \mapsto \text{TRUE},$$
$$mmatchIndex \qquad \mapsto m.mprevLogIndex + Len(m.mentries),$$
$$msource \qquad\qquad \mapsto i,$$
$$mdest \qquad\qquad \mapsto j],$$
$$m)$$
$$\land \text{UNCHANGED } \langle votedFor, currentTerm, log, state\rangle$$
$$\lor \quad \text{conflict: remove 1 entry}$$
$$\land m.mentries \neq \langle\rangle$$
$$\land Len(log[i]) \geq index$$
$$\land log[i][index].term \quad \neq m.mentries[1].term$$
$$\land \text{LET } new \triangleq [index2 \in 1 .. (Len(log[i]) - 1) \mapsto log[i][index2]]$$
$$\text{IN } \quad log' = [log \text{ EXCEPT } ![i] = new]$$
$$\land \text{UNCHANGED } \langle serverVars, commitIndex, messages\rangle$$
$$\lor \quad \text{no conflict: append entry}$$
$$\land m.mentries \neq \langle\rangle$$
$$\land Len(log[i]) = m.mprevLogIndex$$
$$\land log' = [log \text{ EXCEPT } ![i] = Append(log[i], m.mentries[1])]$$
$$\land \text{UNCHANGED } \langle serverVars, commitIndex, messages\rangle$$
$$\land \text{UNCHANGED } \langle candidateVars, leaderVars\rangle$$

Server $i$ receives an *AppendEntries* response from server $j$ with
$m.mterm = currentTerm[i]$.
$$HandleAppendEntriesResponse(i, j, m) \triangleq$$
$$\land m.mterm = currentTerm[i]$$
$$\land \lor \land m.msuccess \quad \text{successful}$$
$$\land nextIndex' \quad = [nextIndex \text{ EXCEPT } ![i][j] = m.mmatchIndex + 1]$$
$$\land matchIndex' = [matchIndex \text{ EXCEPT } ![i][j] = m.mmatchIndex]$$
$$\lor \land \neg m.msuccess \quad \text{not successful}$$
$$\land nextIndex' = [nextIndex \text{ EXCEPT } ![i][j] = Max(\{nextIndex[i][j] - 1, 1\})]$$
$$\land \text{UNCHANGED } \langle matchIndex\rangle$$
$$\land Discard(m)$$

10

41

$\land$ UNCHANGED $\langle serverVars,\ candidateVars,\ logVars,\ elections\rangle$

Detached server $i$ receives a *CatchupRequest* from leader $j$.
$HandleCatchupRequest(i,\ j,\ m) \triangleq$
  $\lor\ \land\ m.mterm < currentTerm[i]$
    $\land\ Reply([mtype\ \mapsto\ CatchupResponse,$
            $mterm \mapsto currentTerm[i],$
            $msuccess \mapsto \text{FALSE},$
            $mmatchIndex \mapsto 0,$
            $msource \mapsto i,$
            $mdest \mapsto j,$
            $mroundsLeft \mapsto 0],$
            $m)$
    $\land$ UNCHANGED $\langle serverVars,\ candidateVars,$
        $leaderVars,\ logVars\rangle$
  $\lor\ \land\ m.mterm \geq currentTerm[i]$
    $\land\ currentTerm' = [currentTerm \text{ EXCEPT } ![i] = m.mterm]$
    $\land\ log' = [log \text{ EXCEPT } ![i] = SubSeq(log[i],\ 1,\ m.mlogLen) \circ m.mentries]$
    $\land\ Reply([mtype\ \mapsto\ CatchupResponse,$
            $mterm \mapsto currentTerm[i],$
            $msuccess \mapsto \text{TRUE},$
            $mmatchIndex \mapsto Len(log[i]),$
            $msource \mapsto i,$
            $mdest \mapsto j,$
            $mroundsLeft \mapsto m.mrounds - 1],$
            $m)$
    $\land$ UNCHANGED $\langle state,\ votedFor,\ candidateVars,\ leaderVars,$
        $commitIndex\rangle$

Leader $i$ receives a *CatchupResponse* from detached server $j$.
$HandleCatchupResponse(i,\ j,\ m) \triangleq$
  A real system checks for progress every timeout interval.
  Assume that if this response is called, the new server
  has made progress.
  $\land\ \lor\ \land\ m.msuccess$
      $\land\ \lor\ \land\ m.mmatchIndex \neq commitIndex[i]$
          $\land\ m.mmatchIndex \neq matchIndex[i][j]$
        $\lor\ m.mmatchIndex = commitIndex[i]$
      $\land\ state[i] = Leader$
      $\land\ m.mterm = currentTerm[i]$
      $\land\ j \notin GetConfig(i)$
      $\land\ nextIndex' = [nextIndex \text{ EXCEPT } ![i][j] = m.mmatchIndex + 1]$
      $\land\ matchIndex' = [matchIndex \text{ EXCEPT } ![i][j] = m.mmatchIndex]$
      $\land\ \lor\ \land\ m.mroundsLeft \neq 0$
          $\land\ Reply([mtype\ \mapsto\ CatchupRequest,$

11

$$
\begin{aligned}
&\quad\quad\quad mterm \mapsto currentTerm[i], \\
&\quad\quad\quad mentries \mapsto SubSeq(log[i], \\
&\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad nextIndex[i][j], \\
&\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad commitIndex[i]), \\
&\quad\quad\quad mLogLen \mapsto nextIndex[i][j] - 1, \\
&\quad\quad\quad msource \ \mapsto i, \\
&\quad\quad\quad mdest \mapsto j, \\
&\quad\quad\quad mrounds \mapsto m.mroundsLeft], \\
&\quad\quad\quad m)
\end{aligned}
$$

$\vee \ \wedge m.mroundsLeft = 0$

A real system makes sure the final call to this handler is
received after a timeout interval.
We assume that if a timeout happened, the message
has already been dropped.

$\wedge Reply([mtype \mapsto CheckOldConfig,$
$\quad\quad\quad mterm \mapsto currentTerm[i],$
$\quad\quad\quad madd \mapsto \text{TRUE},$
$\quad\quad\quad mserver \mapsto j,$
$\quad\quad\quad msource \mapsto i,$
$\quad\quad\quad mdest \mapsto i], m)$

$\wedge \text{UNCHANGED } \langle elections \rangle$

$\vee \ \wedge \ \vee \neg m.msuccess$
$\quad\quad \vee \ \wedge \ \vee m.mmatchIndex = commitIndex[i]$
$\quad\quad\quad\quad \vee m.mmatchIndex = matchIndex[i][j]$
$\quad\quad\quad \wedge m.mmatchIndex \neq commitIndex[i]$
$\quad\quad \vee state[i] \neq Leader$
$\quad\quad \vee m.mterm \neq currentTerm[i]$
$\quad\quad \vee j \in GetConfig(i)$
$\quad \wedge Discard(m)$
$\quad \wedge \text{UNCHANGED } \langle leaderVars \rangle$
$\wedge \text{UNCHANGED } \langle serverVars, candidateVars, logVars \rangle$

Leader $i$ receives a *CheckOldConfig* message.
$HandleCheckOldConfig(i, m) \ \triangleq$
$\quad \vee \ \wedge state[i] \neq Leader \vee m.mterm = currentTerm[i]$
$\quad\quad \wedge Discard(m)$
$\quad\quad \wedge \text{UNCHANGED } \langle serverVars, candidateVars, leaderVars, logVars \rangle$
$\quad \vee \ \wedge state[i] = Leader \wedge m.mterm = currentTerm[i]$
$\quad\quad \wedge \ \vee \ \wedge GetMaxConfigIndex(i) \leq commitIndex[i]$
$\quad\quad\quad\quad \wedge \text{LET } newConfig \ \triangleq \ \text{IF } m.madd \text{ THEN UNION } \{GetConfig(i), \{m.mserver\}\}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \text{ELSE } \ GetConfig(i) \setminus \{m.mserver\}$
$\quad\quad\quad\quad\quad\quad newEntry \ \triangleq \ [term \mapsto currentTerm[i], type \mapsto ConfigEntry, value \mapsto newConfig]$
$\quad\quad\quad\quad\quad\quad newLog \ \triangleq \ Append(log[i], newEntry)$
$\quad\quad\quad\quad \text{IN } \quad log' = [log \text{ EXCEPT } ![i] = newLog]$
$\quad\quad\quad \wedge Discard(m)$

12

43

$\wedge$ UNCHANGED $\langle commitIndex \rangle$
$\vee \wedge GetMaxConfigIndex(i) > commitIndex[i]$
$\wedge Reply([mtype \mapsto CheckOldConfig,$
$\qquad\qquad mterm \mapsto currentTerm[i],$
$\qquad\qquad madd \mapsto m.madd,$
$\qquad\qquad mserver \mapsto m.mserver,$
$\qquad\qquad msource \mapsto i,$
$\qquad\qquad mdest \mapsto i],$
$\qquad\qquad m)$
$\quad\wedge$ UNCHANGED $\langle logVars \rangle$
$\wedge$ UNCHANGED $\langle serverVars,\ candidateVars,\ leaderVars \rangle$

Any *RPC* with a newer term causes the recipient to advance its term first.
$UpdateTerm(i,\ j,\ m) \triangleq$
$\quad \wedge m.mterm > currentTerm[i]$
$\quad \wedge currentTerm' \quad = [currentTerm \text{ EXCEPT } ![i] = m.mterm]$
$\quad \wedge state' \qquad\quad\ = [state \qquad\quad \text{EXCEPT } ![i] \quad = Follower]$
$\quad \wedge votedFor' \qquad = [votedFor \qquad \text{EXCEPT } ![i] \ = Nil]$
$\qquad$ messages is unchanged so $m$ can be processed further.
$\quad \wedge$ UNCHANGED $\langle messages,\ candidateVars,\ leaderVars,\ logVars \rangle$

Responses with stale terms are ignored.
$DropStaleResponse(i,\ j,\ m) \triangleq$
$\quad \wedge m.mterm < currentTerm[i]$
$\quad \wedge Discard(m)$
$\quad \wedge$ UNCHANGED $\langle serverVars,\ candidateVars,\ leaderVars,\ logVars \rangle$

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.
$\qquad\quad \vee UpdateTerm(i,\ j,\ m)$
$\qquad\quad \vee \wedge m.mtype = RequestVoteRequest$
$\qquad\qquad\quad \wedge HandleRequestVoteRequest(i,\ j,\ m)$
$\qquad\quad \vee \wedge m.mtype = RequestVoteResponse$
$\qquad\qquad\quad \wedge \vee DropStaleResponse(i,\ j,\ m)$
$\qquad\qquad\qquad\quad \vee HandleRequestVoteResponse(i,\ j,\ m)$
$\qquad\quad \vee \wedge m.mtype = AppendEntriesRequest$
$\qquad\qquad\quad \wedge HandleAppendEntriesRequest(i,\ j,\ m)$
$\qquad\quad \vee \wedge m.mtype = AppendEntriesResponse$
$\qquad\qquad\quad \wedge \vee DropStaleResponse(i,\ j,\ m)$
$\qquad\qquad\qquad\quad \vee HandleAppendEntriesResponse(i,\ j,\ m)$
$\qquad\quad \vee \wedge m.mtype = CatchupRequest$
$\qquad\qquad\quad \wedge HandleCatchupRequest(i,\ j,\ m)$

13

$\lor \; \land m.mtype = CatchupResponse$
$\quad\;\; \land HandleCatchupResponse(i,\, j,\, m)$
$\lor \; \land m.mtype = CheckOldConfig$
$\quad\;\; \land HandleCheckOldConfig(i,\, m)$

End of message handlers.

---

Network state transitions

The network duplicates a message
$DuplicateMessage(m) \;\triangleq$
$\quad \land Send(m)$
$\quad \land$ UNCHANGED $\langle serverVars,\, candidateVars,\, leaderVars,\, logVars \rangle$

The network drops a message
$DropMessage(m) \;\triangleq$
$\quad \land Discard(m)$
$\quad \land$ UNCHANGED $\langle serverVars,\, candidateVars,\, leaderVars,\, logVars \rangle$

---

Model invariants.

Safety property that only a single leader can be elected at a time.
$OneLeader \;\triangleq\; Cardinality(\{i \in Server : state[i] = Leader\}) \leq 1$

---

Defines how the variables may transition.
$Next \;\triangleq\; \land \lor \exists\, i \in Server : Restart(i)$
$\qquad\qquad\quad \lor \exists\, i \in Server : Timeout(i)$
$\qquad\qquad\quad \lor \exists\, i,\, j \in Server : RequestVote(i,\, j)$
$\qquad\qquad\quad \lor \exists\, i \in Server : BecomeLeader(i)$
$\qquad\qquad\quad \lor \exists\, i \in Server,\, v \in Value : ClientRequest(i,\, v)$
$\qquad\qquad\quad \lor \exists\, i,\, j \in Server : AddNewServer(i,\, j)$
$\qquad\qquad\quad \lor \exists\, i,\, j \in Server : DeleteServer(i,\, j)$
$\qquad\qquad\quad \lor \exists\, i \in Server : AdvanceCommitIndex(i)$
$\qquad\qquad\quad \lor \exists\, i,\, j \in Server : AppendEntries(i,\, j)$
$\qquad\qquad\quad \lor \exists\, m \in$ DOMAIN $messages : Receive(m)$
$\qquad\qquad\quad \lor \exists\, m \in$ DOMAIN $messages : DuplicateMessage(m)$
$\qquad\qquad\quad \lor \exists\, m \in$ DOMAIN $messages : DropMessage(m)$
$\qquad\qquad\quad$ History variable that tracks every $log$ ever:
$\qquad\qquad \land allLogs' = allLogs \cup \{log[i] : i \in Server\}$

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

14

\ * *Changelog*:

\ *

\ * 2015 − 05 − 10:

\ * − Add cluster membership changes as described in Section 4 of

\ *     *Diego Ongaro*. Consensus: Bridging theory and practice.

\ *     *PhD* thesis, *Stanford* University, 2014.

\ * This introduces: *InitServer*, *ValueEntry*, *ConfigEntry*, *CatchupRequest*,

\ *     *CatchupResponse*, *CheckOldConfig*, *GetMaxConfigIndex*,

\ *     *GetConfig* (parameterized), *AddNewServer*, *DeleteServer*,

\ *     *HandleCatchupRequest*, *HandleCatchupResponse*,

\ *     *HandleCheckOldConfig*

\ *

\ * 2014 − 12 − 02:

\ * − Fix *AppendEntries* to only send one entry at a time, as originally

\ * intended. Since *SubSeq* is inclusive, the upper bound of the range should

\ * have been *nextIndex*, not *nextIndex* + 1. Thanks to *Igor Kovalenko* for

\ * reporting the issue.

\ * − Change *matchIndex'* to *matchIndex* (without the apostrophe) in

\ * *AdvanceCommitIndex*. This apostrophe was not intentional and perhaps

\ * confusing, though it makes no practical difference (*matchIndex'* equals

\ * *matchIndex*). Thanks to *Hugues Evrard* for reporting the issue.

\ *

\ * 2014 − 07 − 06:

\ * − Version from *PhD* dissertation

15