# 1   Introduction

Raft is a consensus algorithm designed for understandability and ease of implementation compared to the Paxos algorithm, and thus, it is widely utilized in various products to resolve consensus issues.

In distributed systems, Raft requires odd number of servers to always maintain a quorum that remains functional and communicative in case of network partition faults. Practically, this means a Raft cluster size should be at least three to tolerate a single point of failure.

While the requirement for an odd cluster size isn't problematic for large-scale distributed systems, it can pose a challenge for budget-constrained customers who need to reach consensus with fewer servers. One solution to this issue is to substitute one server in the cluster with a cost-effective entity known as a witness. The witness acts as a tie-breaker, ensuring the maintenance of a quorum in case of an even split of servers due to network issues. Witness typically operates on low-configuration hardware to save costs, seldom participates in critical data paths, and only persists a minimal amount of data. In the Paxos algorithm, a witness can be implemented as shared storage, which is not only inexpensive but also widely available to most customers.

Efforts have been undertaken to incorporate witness into Raft algorithm. For example, Diego Ongaro's Ph.D. dissertation explored the reduction of server numbers by using a witness to store log entries for any failed server until recovery or replacement. Jehan-François Pâris and Darrell D. E. Long proposed a type of follower-only server witness that doesn't persist user data. TiKV implements a witness in a similar manner as described by Ongaro: a log-only entity that persists only raft logs without applying them, and supports switching between witness and non-witness servers.

However, all existing research and implementations necessitate a standalone server as a witness, adding to deployment complexity and reducing cost-efficiency compared to shared storage in the Paxos algorithm. Additionally, the witness must participate in log replication when it needs to conform to a quorum with other servers (when some servers are down), implying that its low-configuration hardware could potentially become a performance bottleneck in such scenarios.

# 2 Algorithm

In this section, we present the extended Raft algorithm, which is a variant of the Raft algorithm. It is designed for cluster with regular servers and at most one witness. The witness functions strictly as a follower server, implying that its role never transitions to candidate or leader. Furthermore, a witness only maintains a minimal set of metadata.

Our newly proposed algorithm minimizes data traffic and reduces the frequency of witness visits while preserving all key properties of the Raft algorithm. This implies that the extended Raft algorithm ensures that each of the following properties is consistently upheld.

- Election Safety

- Leader Append-only

- Log Matching

- Leader Completeness

- State Machine Safety

In the rest of this document, unless otherwise specified, the term 'server' will denote either a regular server or a witness. Consequently, the server set, denoted as $Server$, that constitutes a cluster could either be $RegularServer \cup witness$ or simply $RegularServer$, where $RegularServer$ is the set of all regular servers.

## 2.1 Concepts and Definitions

**Definition 1.** *A leader's **replication set** is a subset of the cluster server set, with its cardinality matching that of the regular server set.*

Let $ReplicationSets$ represent the set of all replication sets within a cluster. We have:

$$ReplicationSets \triangleq \begin{cases} \{Server\}, & witness \notin Server \\ \{Server \setminus \{x\} : x \in Server\}, & witness \in Server \end{cases}$$

A leader maintains its replication set based on its view of the servers and modifies it under specific conditions. A leader may alter its replication set multiple times within a single term. To uniquely identify each replication set, we introduce the notion of a 'subterm' as detailed below.

**Definition 2.** *A term is segmented into* SUBTERMS*, each begins with a replication set.*

Subterms are sequentially numbered using consecutive integers, beginning from 0 for each term. The leader retains its current subterm and increments it when the replication set is altered. Consequently, a leader's replication set remains static throughout its subterm.

## 2.2 States

Figure 2.1 illustrates the additional states introduced by the Extended Raft algorithm to accommodate the witness.

Each leader maintains three additional volatile variables: *replicationSet*, *currentSubterm*, and *witnessSubterm*. The *replicationSet* represents the leader's current replication set, initialized to *RegularServer*, and is modified by *AdjustReplicationSet* action. The *currentSubterm* represents the latest subterm of the leader, initialized to 0 and incremented upon a change in the leader's replication set. The *witnessSubterm* denotes the latest subterm during which the leader replicated its log entry to the witness.

In addition to *currentTerm* and *votedFor*, the witness also includes *witnessReplicationSet*, *witnessLastLogTerm*, and *witnessLastLogSubterm*. These represent the most recent replication set, log entry term, and log entry subterm that the leader sent to the witness. We will further discuss this in Section 2.3.

Beyond the additional states in the leader and witness, each log entry is also associated with the leader's current subterm number when it is appended to the leader's log (Figure 2.2). This subterm number is replicated and persisted alongside the log entry on all regular servers. We now denote a log entry as $\langle index, term, subterm \rangle$, which is uniquely identified by *index* and *term*, and associated with *subterm*.

## 2.3 Log Replication

The leader replicates its log to regular servers in the exact same way as that in the Raft algorithm, and the regular servers handle the received log entries in the same way as well.

However, for the witness, log replication is performed differently, as shown in Figure **??**. In the Extended Raft algorithm, the leader sends an *AppendEntriesToWitnessRequest*

3

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 $matchIndex$

Leader's replication set

VARIABLE $replicationSet$

Leader's subterm number

VARIABLE $currentSubterm$

The latest subterm that this leader has written to the witness

VARIABLE $witnessSubterm$

$$leaderVars \triangleq \langle nextIndex,\ matchIndex,\ elections,$$
$$replicationSet,\ currentSubterm,\ witnessSubterm \rangle$$

The following variables are persisted only on witness:

The latest replication set sent from leader

VARIABLE $witnessReplicationSet$

The latest term of replicated log entry

VARIABLE $witnessLastLogTerm$

The latest subterm of replicated log entry

VARIABLE $witnessLastLogSubterm$

$$witnessVars \triangleq \langle witnessReplicationSet,$$
$$witnessLastLogTerm,\ witnessLastLogSubterm \rangle$$

Figure 2.1: State

```
  Leader i receives a client request to add v to the log.
ClientRequest(i, v)  ≜
    ∧ state[i] = Leader
    ∧ LET  entry    ≜ [term      ↦ currentTerm[i],
                        subterm   ↦ currentSubterm[i],
                        value     ↦ v]
           newLog   ≜ Append(log[i], entry)
       IN   log′ = [log EXCEPT ![i] = newLog]
    ∧ UNCHANGED ⟨messages, serverVars, candidateVars,
                   leaderVars, commitIndex, witnessVars⟩
```
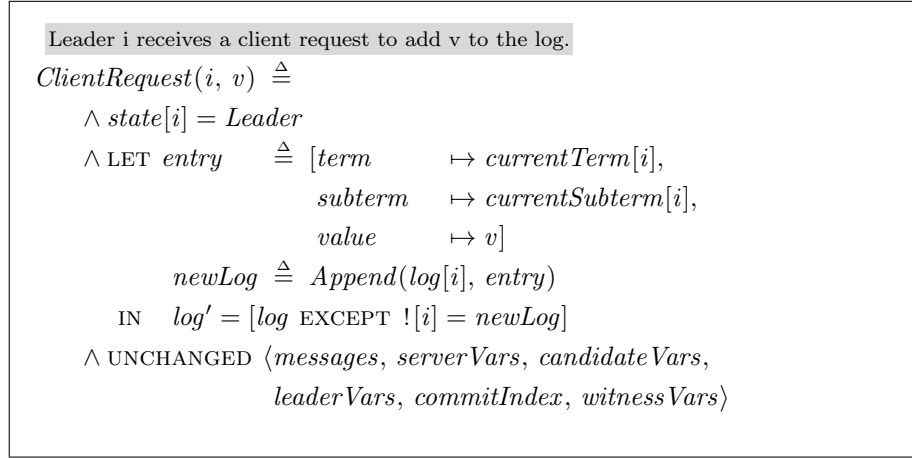
Figure 2.2: Client request

to the witness, which includes the current replication and metadata of the log entry that satisfies the following conditions:

1. The log entry's term and subterm are equal to *currentTerm* and *currentSubterm*, respectively.

2. The leader's current replication set includes the witness.

3. The leader has received acknowledgments for the log entry from at least a subquorum (one server away from forming a quorum) in the leader's current replication set.

4. The leader has not received any acknowledgment from the witness during the current subterm.

The *AppendEntriesToWitness* can be considered as a special implementation of the *AppendEntries* action with batching of log prefix up to a specific log entry. This log entry belongs to the current subterm and has received subquorum acknowledgements from the current replication set. In addition to the message fields that also exist in *AppendEntriesRequest*, the dispatched message *AppendEntriesToWitnessReques* also contains the current replication set and the metadata of the last batched entries.

Upon receiving the request, the witness persists the replication set, as well as the term and subterm of the log entry, then responds to the leader with an *AppendEntriesResponse* message (Figure 2.4).

$AppendEntriesToWitness(i) \triangleq$

    $\wedge state[i] = Leader$

    $\wedge WitnessID \in replicationSet[i]$

    $\wedge$ LET $Agree(index) \triangleq$

                $\{i\} \cup \{k \in replicationSet[i] : matchIndex[i][k] \geq index\}$

          $IsAgreed(k) \triangleq$

              $\wedge log[i][k].term \quad\; = currentTerm[i]$

              $\wedge log[i][k].subterm = currentSubterm[i]$

              $\wedge (\{WitnessID\} \cup Agree(k)) \in Quorum$

          $agreeIndexes \triangleq \{k \in 1 .. Len(log[i]) : IsAgreed(k)\}$

          $lastEntry \quad\;\; \triangleq Max(agreedIndex)$

      IN    $\wedge agreeIndexes \neq \{\}$

          $\wedge \vee$

              $\wedge currentSubterm[i] > witnessSubterm$

              $\wedge Send([mtype \qquad\qquad \mapsto AppendEntriesToWitnessRequest,$

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

                       $mlogTerm \qquad\;\; \mapsto log[i][lastEntry].term,$

                       $mlogSubterm \;\; \mapsto log[i][lastEntry].subterm,$

                       $mreplicationSet \mapsto replicaitonSet[i],$

                       $mindex \qquad\qquad \mapsto lastEntry,$

                    

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

                       $mentries \qquad\;\; \mapsto SubSeq(log[i], 1, lastEntry),$

                       $msource \qquad\qquad \mapsto i,$

                       $mdest \qquad\qquad\;\; \mapsto WitnessID])$

             $\wedge$ UNCHANGED $\langle leaderVars \rangle$

           $\vee$

              $\wedge currentSubterm[i] = witnessSubterm[i]$

              $\wedge Send([mtype \qquad\qquad\;\; \mapsto AppendEntriesResponse,$

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

                       $msuccess \qquad\qquad \mapsto$ TRUE,

                       $mmatchIndex \quad\;\; \mapsto lastEntry,$

                       $msource \qquad\qquad \mapsto WitnessID,$

                       $mdest \qquad\qquad\;\;\; \mapsto i])$

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

6

Figure 2.3: AppendEntriesToWitness

$HandleAppendEntriesToWitnessRequest(j, m) \triangleq$
    $\land\ m.mterm \leq currentTerm[WitnessID]$
    $\land\ \lor$  reject request
        $\land\ m.mterm < currentTerm[WitnessID]$
        $\land\ Reply([mtype \quad\quad\quad \mapsto AppendEntriesResponse,$
                 $mterm \quad\quad\quad\quad \mapsto currentTerm[WitnessID],$
                 $msuccess \quad\quad\quad \mapsto \text{FALSE},$
                 $mmatchIndex \quad\quad \mapsto 0,$
                 $msource \quad\quad\quad \mapsto WitnessID,$
                 $mdest \quad\quad\quad\quad \mapsto j],$
                 $m)$
        $\land\ \text{UNCHANGED}\ \langle serverVars,\ logVars,\ witnessVars \rangle$
      $\lor$  accept request, always no conflict.
        $\land\ m.mterm = currentTerm[WitnessID]$
        $\land\ \lor\ m.mlastLogTerm > witnessLastLogTerm[WitnessID]$
          $\lor\ \land\ m.mlastLogTerm = witnessLastLogTerm[WitnessID]$
            $\land\ m.mlastLogSubterm > witnessLastLogSubterm[WitnessID]$
        $\land\ witnessReplicationSet' =$
          $[witnessReplicationSet\ \text{EXCEPT}\ ![WitnessID] = m.mreplicationSet]$
        $\land\ witnessLastLogTerm' =$
          $[witnessLastLogTerm\ \text{EXCEPT}\ ![WitnessID] = m.mlastLogTerm]$
        $\land\ witnessLastLogSubterm' =$
          $[witnessLastLogSubterm\ \text{EXCEPT}\ ![WitnessID] = m.mlastLogSubterm]$
        $\land\ log' = [log\ \text{EXCEPT}\ ![WitnessID] = m.mentries]$
        $\land\ Reply([mtype \quad\quad\quad\quad \mapsto AppendEntriesResponse,$
                 $mterm \quad\quad\quad\quad \mapsto currentTerm[WitnessID],$
                 $msuccess \quad\quad\quad \mapsto \text{TRUE},$
                 $mmatchIndex \quad\quad \mapsto m.mindex,$
                 $msource \quad\quad\quad \mapsto WitnessID,$
                 $mdest \quad\quad\quad\quad \mapsto j],$
                 $m)$
        $\land\ \text{UNCHANGED}\ \langle serverVars,\ log \rangle$
    $\land\ \text{UNCHANGED}\ \langle candidateVars,\ leaderVars \rangle$

7

Figure 2.4: Handle AppendEntriesToWitness

Note that the *mlog* and *mentries* fields in the message are solely used for proof. They do not exist in the actual implementation, nor will the witness persist the full log prefix in its durable storage. The persistence of *m.mentries* to *log* in *HandleAppendEntriesToWitnessRequest* is exclusively used for proof purposes.

When the leader receives an acknowledgment from the witness, it also updates its *witnessSubterm* to the current subterm if the acknowledged log entry is in the current subterm. Then, according to condition **??**, the leader no longer sends *AppendEntriesToWitnessRequest* in the current subterm. Instead, the leader treats subsequent log entries as already acknowledged by the witness if they fulfill all conditions except condition **??**. This is referred to as **Shortcut Replication**, where the leader sends an *AppendEntriesResponse* to itself on behalf of the witness.

Note that the log entry in *AppendEntriesToWitness* must be within the current term and subterm. So, when the leader advances its subterm (a consequence of replication set adjustment), a new empty log entry must be appended to the leader's log. This is to ensure the leader won't be impeded by condition 1 in committing log entries.

Leader immediately commits a log entry during its term when the log entry is acknowledged by a quorum, in the same way as in Raft. The acknowledgment either comes from a regular server or the witness. Since the leader sends *AppendEntriesToWitnessRequest* to the witness only after it has received subquorum acknowledgments from regular servers, any log entry acknowledged by the witness is immediately committed.

## 2.4 Replication Set Adjustment

Leader maintains the replication set and modifies it whenever necessary. This process is referred to as **Replication Set Adjustment**. Figure 2.5 illustrates a straightforward method to adjust the replication set. Although there could be multiple ways to adjust a replication set in practical implementations, such a simple method can be used in the algorithm to minimize the atomic regions without loss of generality.

To commit, the leader's replication set must contain at least a subquorum of healthy regular servers if it includes a witness. Therefore, in a practical implementation, replication set adjustment should be performed in a way such that the resulting replication set includes as many healthy regular servers as

Adjust replication set on leader i. This action updates replication set
by swapping items inside and outside. While implementations may change
replication set in various ways, the spec uses this simple swapping to
minimize atomic regions without loss of generality.

$AdjustReplicationSet(i) \triangleq$

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

$\quad \wedge replicationSet[i] \neq Server$

$\quad$ Swap server outside replication set with some server inside

$\quad \wedge \text{LET } in \quad \triangleq \text{ CHOOSE } x \in replicationSet[i] : x \neq i$

$\qquad\qquad out \triangleq \text{ CHOOSE } x \in Server \setminus replicationSet[i] : \text{TRUE}$

$\qquad \text{IN} \quad replicationSet' \quad =$

$\qquad\qquad\qquad [replicationSet \text{ EXCEPT } ![i] = (@ \setminus \{in\}) \cup \{out\}]$

$\quad \wedge currentSubTerm' = [currentSubTerm \text{ EXCEPT } ![i] = @ + 1]$

$\quad \wedge \text{UNCHANGED } \langle messages, serverVars, candidateVars, nextIndex,$

$\qquad\qquad\qquad\qquad matchIndex, witnessSubterm, log, commitIndex,$

$\qquad\qquad\qquad\qquad witnessVars \rangle$

Figure 2.5: Replicaiton set adjustment

possible. To achieve this, leader can track the status of each peer regular server by monitoring their responses. If leader receives no response from a regular peer server over an *electiontimeout*, it assumes the regular server is unreachable and initiates a replication set adjustment. Let *Reachable* and *Unreachable* represent the set of reachable regular servers and unreachable regular servers from the leader's perspective. A practical replication set adjustment can be implemented as below:

1. When a regular server becomes a leader, its replication set is initialized to all regular servers *RegularServer* in the leader's configuration.

2. If all servers in *RegularServer* are reachable from the leader, leader changes its replication set to *RegularServer*, i.e.,

$$\forall s \in RegularServer : s \in Reachable \implies ReplicationSet' = RegularServer$$

3. Leader swaps one unreachable regular server inside its replication set with the reachable regular server or witness outside, i.e.,

$$\exists x \in ReplicationSet, y \in Server \setminus ReplicationSet : \wedge\, x \in Unreachable$$
$$\wedge \vee\, y \in Reachable$$
$$\vee\, y = witness$$
$$\implies ReplicationSet' = Server \setminus \{x\}$$

4. Leader increments its *currentSubterm* and appends a new empty log entry if its replication set changes.

## 2.5   Leader Election

When a regular server transitions into candidate, it requests votes from all peer servers including witness. The extended Raft algorithm adheres to the exact same rules as those in Raft when it comes to requesting and granting votes between a candidate and a regular server. However, the candidate does not request a vote from the witness until it has received subquorum votes from regular servers. Figure 2.6 describes the new *RequestWitnessVote* action for a candidate to request vote from a witness.

Considering that the witness does not persist the full log prefix, the extended Raft algorithm employs a new set of rules for the witness in leader elections.

```
Candidate i sends witness a RequestWitnessVote request.
RequestWitnessVote(i) ≜
    ∧ state[i] = Candidate
    ∧ WitnessID ∉ votesResponded[i]
    ∧ (votesGranted[i] ∪ {WitnessID}) ∈ Quorum
    ∧ Send([mtype              ↦ RequestWitnessVoteRequest,
            mterm              ↦ currentTerm[i],
            mlastLogTerm       ↦ LastTerm(log[i]),
            mlastLogSubterm    ↦ LastSubterm(log[i]),
            mvotesGranted      ↦ votesGranted[i],
            msource            ↦ i,
            mdest              ↦ WitnessID])
    ∧ UNCHANGED ⟨serverVars, candidateVars, leaderVars, logVars,
                 witnessVars⟩
```
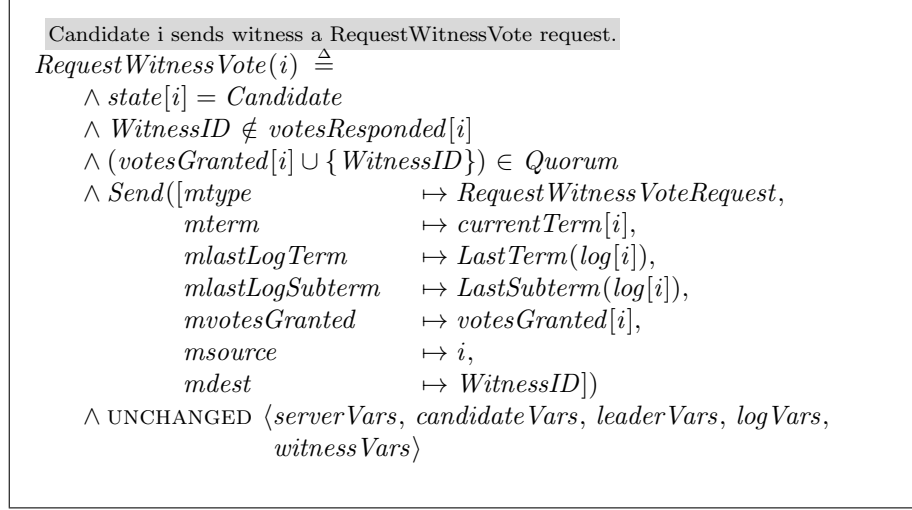
Figure 2.6: RequestWitnessVote

Figure 2.7 describes the new action for the witness to handle a vote request from
a candidate. The *RequestWitnessVoteRequest* message includes the *term* and
*subterm* of the last log entry in the candidate, as well as the votes the candidate
has already received in the current election. The witness votes for a candidate
if it has not casted vote in the current term and any of the following conditions
is true:

- Candidate's last log entry has larger term, i.e. $m.mlastLogTerm > witnessLastLogTerm$.

- Candidate's last log entry has same term but larger subterm, i.e. $m.mlastLogTerm = witnessLastLogTerm \land m.lastLogSubterm > witnessLastLogSubterm$.

- Candidate's last log entry has same term and subterm comparing to witness, and all subqorum votes it got are from regular servers in witness's replication set, i.e.

$$m.mlastLogTerm = witnessLastLogTerm$$
$$m.mlastLogSubterm = witnessLastLogSubterm$$
$$m.mvotesGranted \subseteq witnessReplicationSet$$

11

$\fbox{\begin{minipage}{0.9\textwidth}\end{minipage}}$

Witness receives a RequestWitnessVote request from server j with
m.mterm ¡= currentTerm[WitnessID].

$HandleRequestWitnessVoteRequest(j, m) \triangleq$

    LET $logOk \triangleq \lor m.mlastLogTerm > witnessLastLogTerm$

                        $\lor \land m.mlastLogTerm = witnessLastLogTerm$

                           $\land m.mlastLogSubterm > witnessLastLogSubterm$

                        $\lor \land m.mlastLogTerm = witnessLastLogTerm$

                           $\land m.mlastLogSubterm = witnessLastLogSubterm$

                           $\land m.mvotesGranted \subseteq witnessReplicationSet$

        $grant \triangleq \land m.mterm = currentTerm[WitnessID]$

                  $\land logOk$

                  $\land votedFor[WitnessID] \in \{Nil, j\}$

    IN    $\land m.mterm \leq currentTerm[WitnessID]$

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

          $\lor \neg grant \land \text{UNCHANGED } votedFor$

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

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

                $mvoteGranted \mapsto grant,$

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

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

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

                $msource \qquad\quad \mapsto WitnessID,$

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

                $m)$

        $\land \text{UNCHANGED } \langle state, currentTerm, candidateVars, leaderVars,$

                            $logVars, witnessVars\rangle$

Figure 2.7: HandleRequestVoteToWitness

## 2.6 Membership Reconfiguration

Cluster configuration change can be done in both simple and complex approach in Raft algorithm. The simple approach changes one server at a time. And the complex approach allows arbitrary configuration changes at one time.

The key to the safety of configuration change in Raft is preserving overlap between any majority of the old cluster and any majority of the new cluster. This overlap prevents the cluster from splitting into two independent majorities where two leaders being elected in same term, since otherwise the overlapped server would be able to vote two servers in a term, which is a a contradiction of the specification. This applies to the extended Raft algorithm too because a new leader also needs quorum voters to succeed in an election.Therefore, it is safe to apply same membership configuration methods in Raft algorithm to the extended Raft algorithm either with a sequence of single-server membership change or with an intermediate joint configuration that overlaps both old and new membership configurations.

# 3   Witness Implementation

With the extended Raft algorithm, the witness possesses the following properties that make it particularly suitable for implementation as a storage object:

- Witness does not have volatile variables, implying that all states within the witness are persisted in durable storage.

- Witness only persists a small amount of metadata. Neither user data nor the full log prefix metadata are stored in witness. This results in a trivial and predictable capacity requirement.

- Witness is rarely accessed. The shortcut replication ensures that leader only needs to write to witness when its replication set changes. Therefore, witness may only need to be accessed when faults change (e.g. a server shuts down or recovers) or a candidate requests its vote. In a relatively stable cluster, neither of these events occur frequently.

The absence of volatile variables allows the synchronization of witness states to be distributed across all regular servers. For each log replication between leader and witness, the *HandleAppendEntriesToWitnessRequest* action can be

implemented in the leader as a single operation that loads states from the witness's durable storage, make new states according to *HandleAppendEntriesToWitnessRequest*, and writes back new states to the witness's durable storage. Similarly for voting between candidate and witness, the *HandleRequestWitnessVoteRequest* action can be implemented in the candidate as a single operation that loads states from the witness's durable storage, make new states according to *HandleRequestVoteToWitness*, and writes back the new states to the witness's durable storage. These two implementations can be synchronized using an optimistic concurrency control scheme since they both occur infrequently.

Figure 3.1 illustrates an implementation using a share (NFS or SMB) as a witness. We associate a version number with the states, starting from 0. Each modification to the states increases the version by 1. We can then persist each version of the witness states in a file uniquely named after the version number (e.g. <version>.st). Whenever *HandleAppendEntriesToWitnessRequest* or *HandleRequestVoteToWitness* are invoked, a regular server (leader or candidate) loads the witness states from the file with the largest version number, calculates the new states, and saves the new states back to a distinct temporary file (e.g. <server id>.<version>.st). Then, the regular server tries to create a hard link from the temporary file to the desired versioned state file (i.e. <version+1>.st). As long as the file system guarantees atomicity in creating hardlinked files, the regular server will only succeed if the target versioned state file does not exist. If any other server has finished update states in the same version, the process needs to be repeated from the beginning. The optimistic concurrency control eliminates the need for locking overhead, which is challenging to implement in practical applications. The infrequency of witness visits makes such a scheme very suitable.

In addition to share, the witness functionality can also be incorporated into cloud storage, such as Azure Storage or AWS S3. Many cloud storage vendors provide optimistic concurrency mechanism for updating object data. For instance, Azure Blob Storage allows clients to update a blob object using the original ETag, combined with a conditional header, to ensure that updates only occur if the ETag remains unchanged. This guarantees that no other client has updated the blob object concurrently. The infrequent access requirements and trivial payload size for witness make cloud storage an ideal solution for this function. Consequently, it provides a broader range of implementation options for end-users.

Figure 3.1: Implement witness as a share

# 4 Discussion

## 4.1 Compatibility

In a cluster without witness, the extended Raft algorithm operates identically to the Raft algorithm. A cluste can be formed with mixed servers running with either Raft or the extended Raft algorithms. This compatibility facilitates a seamless upgrade from a cluster based on the Raft algorithm to one using the extended Raft algorithm in practical applications. A binary upgrade can be performed in the existing cluster to incorporate the extended Raft algorithm. Following this, the upgraded cluster functions similarly to a typical Raft-based cluster, without the inclusion of a witness. Once all servers have been upgraded to support the extended Raft algorithm, a membership reconfiguration operation can be performed to integrate a witness into the cluster. Conversely, a cluster functioning with the extended Raft algorithm can be reverted back to Raft algorithm through a binary downgrade, which is possible after removal of witness from the cluster.

The compatibility is crucial for real-world applications. It eliminates the

necessity to establish a new cluster and migrate data from the existing cluster to take advantage of the benefits offered by a witness. As a result, this minimizes disruptions and enhances the operational efficiency of these applications.

## 4.2   Availability

The Raft algorithm ensures the presence of a leader that preserves all committed log entries and continues to commit incoming requests, provided that the cluster has a quorum of connected, healthy servers. This implies that if a follower fails, the leader cannot commit if the remaining followers cannot form a quorum with the leader. Additionally, if a leader fails, a new leader can be elected among the other servers that form a quorum.

However, this is a slightly different in the extended Raft algorithm.

## 4.3   Cascading Failures

The extended Raft algorithm's correctness doesn't guarantee that a leader can always be elected. This predicament often arises when faults occur, leading to a situation where no regular server has a more up-to-date log than witness. In other words, witness would have been the only feasible next leader if it weren't for the lack of a full log prefix. This typically happens as a result of cascading failures where servers go down sequentially.

For instance, consider a cluster composed of two servers ($A$ and $B$) and a witness. If $A$ is the leader and a sequence of events unfold as described below, the cluster will find itself in a situation where $A$ is the only functioning regular server, but its log doesn't contain the new committed log entry (made by $B$). The state in the witness prevents it from casting a vote for $A$, leaving the cluster in a situation without a leader. However, this doesn't violate the correctness of the extended Raft algorithm. And a new leader can always be elected once $B$ restarts.

1. $A$ is shut down

2. $B$ becomes leader

3. $B$ commits a new entry (hence writes states to witness)

4. $A$ starts up

5. $B$ is shut down before it replicating the new committed entry to $A$

In real-world scenarios, when a server shuts down, its subsystems do not cease functioning in the same time. For instance, NIC connecting to the witness ($NIC_w$) may stop after the other NIC that connects to regular servers ($NIC_s$). Moreover, the server's OS ceases operation after all NICs go down. Considering the same cluster as before, if $A$ was operating as leader in the extended Raft algorithm, unfunctional $NIC_s$ would prompt it to adjust its replication set and advance the subterm. A new log entry may then be committed in the new subterm, with its term and subterm recorded to the witness through the still-functioning $NIC_w$. Now after $A$ fully shuts down, $B$ won't be elected as the new leader, since the witness, containing more up-to-date states, rejects casting a vote for it.

This particular issue arises due to the limited role of the witness. However, it's a necessary trade-off to minimize the overall footprint and maintain reasonable performance when faults occur. It's important to note that this issue arises in very rare cases. Adding an extra delay in *HandleAppendEntriesToWitnessRequest* can greatly reduce the chance of the occurence. The impact to performance is trivial, as most commits do not involve writing to the witness.

## 4.4  Catching up upon replication set adjustment

When a follower server's log significantly lags behind the leader's, it may take a considerable amount of time for the follower server to catch up. Until then, new log entries in the leader will not be acknowledged by that follower server. If the live servers merely satisfy the minimal quorum number, the leader won't be able to commit until all servers catch up. This situation arises in Raft when a new server is added to a cluster, and the solution is to add the new server as a learner that receives the log but is not a part of the current membership configuration.

This situation can also occur in the extended Raft algorithm when subterm changes. When a server outside replication set recovers from fault state, the leader may rejoin it to the replication set. If there is only a quorum of live servers (including the witness) in the current replication set, the leader won't be able to commit new log entries until all servers catch up. For instance, in the cluster above, if $A$ commits many entries during $B$'s downtime, $B$'s log will significantly lag behind $A$. After $B$ recovers, $A$ will advance its subterm and change its replication set to $\{A, B\}$. Now, any new entry in $A$ will not be committed until it receives 2 acknowledgements, i.e., both $A$ and $B$ must

acknowledge. As a result, the cluster won't be able to commit any entry because it will take $B$ a long time to catch up with $A$ and acknowledge.

The solution to this issue is similar to the learner approach: leader does not change its replication set to add a regular server until that server's log catches up with the leader.

# 5    Specification

```
 1 ┌─────────────────── MODULE extendedraft ───────────────────┐
 2    This is the formal specification for the extended Raft consensus algorithm.
 3    < TODO: Copyright Info >
 4
 5    Copyright 2014 Diego Ongaro.
 6    This work is licensed under the Creative Commons Attribution − 4.0
 7    International License https://creativecommons.org/licenses/by/4.0/

 9  EXTENDS Naturals, FiniteSets, Sequences, TLC

11    The set of server IDs
12  CONSTANTS Server

14    The set of requests that can go into the log
15  CONSTANTS Value

17    ID of the witness
18  CONSTANTS WitnessID

20    Server states.
21  CONSTANTS Follower, Candidate, Leader

23    A reserved value.
24  CONSTANTS Nil

26    Message types:
27  CONSTANTS RequestVoteRequest, RequestVoteResponse,
28            AppendEntriesRequest, AppendEntriesResponse,
29            RequestWitnessVoteRequest, AppendEntriesToWitnessRequest

31  ├──────────────────────────────────────────────────────────────
32    Global variables
```

18

34  A bag of records representing requests and responses sent from one server

35  to another. *TLAPS* doesn't support the Bags module, so this is a function

36  mapping Message to *Nat*.

37  VARIABLE *messages*

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

40  implementation.

41  Keeps track of successful elections, including the initial logs of the

42  leader and voters' logs. Set of functions containing various things about

43  successful elections (see *BecomeLeader*).

44  VARIABLE *elections*

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

47  implementation.

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

49  VARIABLE *allLogs*

51  ├─────────────────────────────────────────────────────────────────────

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

54  The server's term number.

55  VARIABLE *currentTerm*

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

57  VARIABLE *state*

58  The candidate the server voted for in its current term, or

59  Nil if it hasn't voted for any.

60  VARIABLE *votedFor*

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

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

64  *log* entry. Unfortunately, the Sequence module defines $Head(s)$ as the entry

65  with index 1, so be careful not to use that!

66  VARIABLE *log*

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

68  VARIABLE *commitIndex*

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

71  The following variables are used only on candidates:

72  The set of servers from which the candidate has received a *RequestVote*

73  response in its $currentTerm$.

74  VARIABLE $votesResponded$

75  The set of servers from which the candidate has received a vote in its

76  $currentTerm$.

77  VARIABLE $votesGranted$

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

79  implementation.

80  Function from each server that voted for this candidate in its $currentTerm$

81  to that voter's $log$.

82  VARIABLE $voterLog$

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


85  The following variables are used only on leaders:

86  The next entry to send to each follower.

87  VARIABLE $nextIndex$

88  The latest entry that each follower has acknowledged is the same as the

89  leader's. This is used to calculate $commitIndex$ on the leader.

90  VARIABLE $matchIndex$

91  Leader's replication set

92  VARIABLE $replicationSet$

93  Leader's $subterm$ number

94  VARIABLE $currentSubterm$

95  The latest $subterm$ that this leader has written to the witness

96  VARIABLE $witnessSubterm$

97  $leaderVars \triangleq \langle nextIndex, matchIndex, elections,$

98  $\qquad\qquad replicationSet, currentSubterm, witnessSubterm \rangle$


100  The following variables are persisted only on witness:

101  The latest replication set sent from leader

102  VARIABLE $witnessReplicationSet$

103  The latest term of replicated $log$ entry

104  VARIABLE $witnessLastLogTerm$

105  The latest $subterm$ of replicated $log$ entry

106  VARIABLE $witnessLastLogSubterm$

107  $witnessVars \triangleq \langle witnessReplicationSet,$

108  $\qquad\qquad witnessLastLogTerm, witnessLastLogSubterm \rangle$

111    End of per server variables.

112 ├──────────────────────────────────────────────────────────────────┤

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

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

117 ├──────────────────────────────────────────────────────────────────┤

118    Helpers

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

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

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

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

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

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

128    $LastSubTerm(xlog) \triangleq \text{IF } Len(xlog) = 0 \text{ THEN } 0 \text{ ELSE } xlog[Len(xlog)].subterm$

130    Helper for *Send* and *Reply*. Given a message $m$ and bag of messages, return a

131    new bag of messages with one more $m$ in it.

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

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

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

135       ELSE

136          $msgs \,@@\, (m \!:> 1)$

138    Helper for *Discard* and *Reply*. Given a message $m$ and bag of messages, return

139    a new bag of messages with one less $m$ in it.

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

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

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

143       ELSE

144          $msgs$

146    Add a message to the bag of messages.

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

149    Remove a message from the bag of messages. Used when a server is done

150    processing a message.

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

153     Combination of *Send* and *Discard*

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

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

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

158   $Min(s) \triangleq$ CHOOSE $x \in s : \forall\, y \in s : x \leq y$

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

160   $Max(s) \triangleq$ CHOOSE $x \in s : \forall\, y \in s : x \geq y$

162 $\vdash$ ————————————————————————————————————————

163     Define initial values for all variables

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

166                     $\land\ allLogs\ \ \ \ = \{\}$

167                     $\land\ voterLog = [i \in Server \mapsto [j \in \{\} \mapsto \langle\rangle]]$

168   $InitServerVars \triangleq\ \land\ currentTerm = [i \in Server \mapsto 1]$

169                     $\land\ state\ \ \ \ \ \ \ \ \ = [i \in Server \mapsto Follower]$

170                     $\land\ votedFor\ \ \ \ \ = [i \in Server \mapsto Nil]$

171   $InitCandidateVars \triangleq\ \land\ votesResponded = [i \in Server \mapsto \{\}]$

172                       $\land\ votesGranted\ \ \ = [i \in Server \mapsto \{\}]$

173     The values $nextIndex[i][i]$ and $matchIndex[i][i]$ are never read, since the

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

175     in the functions.

176   $InitLeaderVars \triangleq\ \land\ nextIndex\ \ \ = [i \in Server \mapsto [j \in Server \mapsto 1]]$

177                     $\land\ matchIndex = [i \in Server \mapsto [j \in Server \mapsto 0]]$

178                     $\land\ replicationSet\ \ \ = Server \setminus \{WitnessID\}$

179                     $\land\ currentSubterm = 0$

180   $InitLogVars \triangleq\ \land\ log\ \ \ \ \ \ \ \ \ \ \ \ \ = [i \in Server \mapsto \langle\rangle]$

181                 $\land\ commitIndex\ \ = [i \in Server \mapsto 0]$

182   $InitWitnessVars \triangleq\ \land\ witnessReplicationSet = \{\}$

183                     $\land\ witnessLastLogTerm = 0$

184                     $\land\ witnessLastLogSubterm = 0$

185   $Init \triangleq\ \land\ messages = [m \in \{\} \mapsto 0]$

186          $\land\ InitHistoryVars$

187          $\land\ InitServerVars$

188          $\land\ InitCandidateVars$

189          $\land\ InitLeaderVars$

190          $\land\ InitLogVars$

191           $\wedge$ *InitWitnessVars*

193 $\vdash$────────────────────────────────────────────────

194    Define state transitions

196    Server $i$ restarts from stable storage.

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

198   $Restart(i) \triangleq$

199       $\wedge$  $i$                $\neq WitnessID$

200       $\wedge$  $state'$        $= [state \text{ EXCEPT } ![i] = Follower]$

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

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

203       $\wedge$  $voterLog'$         $= [voterLog \text{ EXCEPT } ![i] = [j \in \{\} \mapsto \langle\rangle]]$

204       $\wedge$  $nextIndex'$        $= [nextIndex \text{ EXCEPT } ![i] = [j \in Server \mapsto 1]]$

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

206       $\wedge$  $commitIndex'$     $= [commitIndex \text{ EXCEPT } ![i] = 0]$

207       $\wedge$  $currentSubterm'$ $= [currentSubterm \text{ EXCEPT } ![i] = 0]$

208       $\wedge$  $replicationSet'$    $= [replicationSet \text{ EXCEPT } ![i]$    $= Server]$

209       $\wedge$  $witnessSubterm'$ $= [witnessSubterm \text{ EXCEPT } ![i] = 0]$

210       $\wedge$ UNCHANGED $\langle messages, currentTerm, votedFor, log, elections, witnessVars \rangle$

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

213   $Timeout(i) \triangleq$  $\wedge i \neq WitnessID$

214                     $\wedge state[i] \in \{Follower, Candidate\}$

215                     $\wedge state' = [state \text{ EXCEPT } ![i] = Candidate]$

216                     $\wedge currentTerm' = [currentTerm \text{ EXCEPT } ![i] = currentTerm[i] + 1]$

217                     Most implementations would probably just set the local vote

218                     atomically, but messaging localhost for it is weaker.

219                     $\wedge votedFor' = [votedFor \text{ EXCEPT } ![i] = Nil]$

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

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

222                     $\wedge voterLog'$         $= [voterLog \text{ EXCEPT } ![i] = [j \in \{\} \mapsto \langle\rangle]]$

223                     $\wedge$ UNCHANGED $\langle messages, leaderVars, logVars, witnessVars \rangle$

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

226   $RequestVote(i, j) \triangleq$

227       $\wedge state[i] = Candidate$

228       $\wedge j \neq WitnessID$

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

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

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

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

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

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

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

236     $\land\ \textsc{unchanged}\ \langle serverVars,\ candidateVars,\ leaderVars,\ logVars,\ witnessVars\rangle$

<br>

238     Candidate $i$ sends witness a *RequestWitnessVote* request.

239 $RequestWitnessVote(i)\ \triangleq$

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

241     $\land\ WitnessID \notin votesResponded[i]$

242     $\land\ (votesGranted[i] \cup \{WitnessID\}) \in Quorum$

243     $\land\ Send([mtype \qquad\qquad\qquad \mapsto RequestWitnessVoteRequest,$

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

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

246            $mlastLogSubterm \quad\ \mapsto LastSubterm(log[i]),$

247            $mvotesGranted \qquad\ \mapsto votesGranted[i],$

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

249            $mdest \qquad\qquad\qquad\ \mapsto WitnessID])$

250     $\land\ \textsc{unchanged}\ \langle serverVars,\ candidateVars,\ leaderVars,\ logVars,$

251                   $witnessVars\rangle$

<br>

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

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

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

257 $AppendEntries(i, j)\ \triangleq$

258     $\land\ i \neq j$

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

260     $\land\ \textsc{let}\ prevLogIndex\ \triangleq\ nextIndex[i][j] - 1$

261           $prevLogTerm\ \triangleq\ \textsc{if}\ prevLogIndex > 0\ \textsc{then}$

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

263                   $\textsc{else}$

264                      $0$

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

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

$$267 \qquad\qquad entries \;\triangleq\; SubSeq(log[i],\; nextIndex[i][j],\; lastEntry)$$

$268 \qquad$ IN $\quad Send([mtype \qquad\qquad\qquad \mapsto AppendEntriesRequest,$

$269 \qquad\qquad\qquad\qquad mterm \qquad\qquad\quad \mapsto currentTerm[i],$

$270 \qquad\qquad\qquad\qquad mprevLogIndex \;\mapsto prevLogIndex,$

$271 \qquad\qquad\qquad\qquad mprevLogTerm \;\;\mapsto prevLogTerm,$

$272 \qquad\qquad\qquad\qquad mentries \qquad\qquad \mapsto entries,$

$273 \qquad\qquad\qquad\qquad$

$274 \qquad\qquad\qquad\qquad$

$275 \qquad\qquad\qquad\qquad mlog \qquad\qquad\qquad \mapsto log[i],$

$276 \qquad\qquad\qquad\qquad mcommitIndex \quad \mapsto Min(\{commitIndex[i],\; lastEntry\}),$

$277 \qquad\qquad\qquad\qquad msource \qquad\qquad \mapsto i,$

$278 \qquad\qquad\qquad\qquad mdest \qquad\qquad\quad \mapsto j])$

$279 \qquad \wedge$ UNCHANGED $\langle serverVars,\; candidateVars,\; leaderVars,\; logVars,$

$280 \qquad\qquad\qquad\qquad witnessVars\rangle$

$282 \qquad$

$283 \qquad$

$284 \quad AppendEntriesToWitness(i) \;\triangleq$

$285 \qquad \wedge state[i] = Leader$

$286 \qquad \wedge WitnessID \in replicationSet[i]$

$287 \qquad \wedge$ LET $Agree(index) \;\triangleq$

$288 \qquad\qquad\qquad \{i\} \cup \{k \in replicationSet[i] : matchIndex[i][k] \geq index\}$

$289 \qquad\qquad IsAgreed(k) \;\triangleq$

$290 \qquad\qquad\qquad \wedge log[i][k].term \quad\;\; = currentTerm[i]$

$291 \qquad\qquad\qquad \wedge log[i][k].subterm = currentSubterm[i]$

$292 \qquad\qquad\qquad \wedge (\{WitnessID\} \cup Agree(k)) \in Quorum$

$293 \qquad\qquad agreeIndexes \;\triangleq\; \{k \in 1 \,..\, Len(log[i]) : IsAgreed(k)\}$

$294 \qquad\qquad lastEntry \qquad \triangleq\; Max(agreedIndex)$

$295 \qquad$ IN $\quad \wedge agreeIndexes \neq \{\}$

$296 \qquad\qquad \wedge \vee$

$297 \qquad\qquad\qquad \wedge currentSubterm[i] > witnessSubterm$

$298 \qquad\qquad\qquad \wedge Send([mtype \qquad\qquad\;\; \mapsto AppendEntriesToWitnessRequest,$

$299 \qquad\qquad\qquad\qquad mterm \qquad\qquad\;\; \mapsto currentTerm[i],$

$300 \qquad\qquad\qquad\qquad mlogTerm \qquad\quad \mapsto log[i][lastEntry].term,$

$301 \qquad\qquad\qquad\qquad mlogSubterm \quad\;\; \mapsto log[i][lastEntry].subterm,$

$302 \qquad\qquad\qquad\qquad mreplicationSet \mapsto replicaitonSet[i],$

$303 \qquad\qquad\qquad\qquad mindex \qquad\qquad\;\; \mapsto lastEntry,$

```
304                              mlog and mentries are used as history variable for
305                              the proof. They do not exist in a real implementation.
306                     mlog            ↦ log[i],
307                     mentries        ↦ SubSeq(log[i], 1, lastEntry),
308                     msource         ↦ i,
309                     mdest           ↦ WitnessID])
310            ∧ UNCHANGED ⟨leaderVars⟩
311         ∨  shortcut replication
312            ∧ currentSubterm[i] = witnessSubterm[i]
313            ∧ Send([mtype              ↦ AppendEntriesResponse,
314                    mterm              ↦ currentTerm[i],
315                    msuccess           ↦ TRUE,
316                    mmatchIndex        ↦ lastEntry,
317                    msource            ↦ WitnessID,
318                    mdest              ↦ i])
319      ∧ UNCHANGED ⟨serverVars, candidateVars, logVars, witnessVars⟩

321    Candidate i transitions to leader.
322  BecomeLeader(i) ≜
323       ∧ state[i] = Candidate
324       ∧ votesGranted[i] ∈ Quorum
325       ∧ state'       = [state EXCEPT ![i] = Leader]
326       ∧ nextIndex'   = [nextIndex EXCEPT ![i] =
327                             [j ∈ Server ↦ Len(log[i]) + 1]]
328       ∧ matchIndex' = [matchIndex EXCEPT ![i] =
329                             [j ∈ Server ↦ 0]]
330       ∧ elections'   = elections ∪
331                             {[eterm     ↦ currentTerm[i],
332                               eleader   ↦ i,
333                               elog      ↦ log[i],
334                               evotes    ↦ votesGranted[i],
335                               evoterLog ↦ voterLog[i]]}
336      ∧ UNCHANGED ⟨messages, currentTerm, votedFor, candidateVars, logVars, witnessVars⟩

338    Leader i receives a client request to add v to the log.
339  ClientRequest(i, v) ≜
340       ∧ state[i] = Leader
```

341      $\wedge$ LET $\;entry\;\;\;\;\stackrel{\Delta}{=}\;[term\;\;\;\;\;\;\mapsto currentTerm[i],$

342                                    $subterm\;\;\;\;\mapsto currentSubterm[i],$

343                                    $value\;\;\;\;\;\;\;\mapsto v]$

344                $newLog\;\;\stackrel{\Delta}{=}\;Append(log[i],\;entry)$

345        IN   $log' = [log$ EXCEPT $![i] = newLog]$

346      $\wedge$ UNCHANGED $\langle messages,\;serverVars,\;candidateVars,$

347                           $leaderVars,\;commitIndex,\;witnessVars\rangle$

349    Leader $i$ advances its $commitIndex$.

350    This is done as a separate step from handling $AppendEntries$ responses,

351    in part to minimize atomic regions, and in part so that leaders of

352    single-server clusters are able to mark entries committed.

353 $AdvanceCommitIndex(i)\;\stackrel{\Delta}{=}$

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

355     $\wedge$ LET   The set of servers that agree up through index.

356           $Agree(index)\;\stackrel{\Delta}{=}\;\{i\} \cup \{k \in Server :$

357                                        $matchIndex[i][k] \geq index\}$

358           The maximum indexes for which a quorum agrees

359           $agreeIndexes\;\stackrel{\Delta}{=}\;\{index \in 1\,..\,Len(log[i]) :$

360                              $Agree(index) \in Quorum\}$

361           New value for $commitIndex'[i]$

362           $newCommitIndex\;\stackrel{\Delta}{=}$

363             IF  $\wedge\;agreeIndexes \neq \{\}$

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

365            THEN

366               $Max(agreeIndexes)$

367            ELSE

368               $commitIndex[i]$

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

370     $\wedge$ UNCHANGED $\langle messages,\;serverVars,\;candidateVars,\;leaderVars,\;log,\;witnessVars\rangle$

372    Adjust replication set on leader $i$. This action updates replication set

373    by swapping items inside and outside. While implementations may change

374    replication set in various ways, the spec uses this simple swapping to

375    minimize atomic regions without loss of generality.

376 $AdjustReplicationSet(i)\;\stackrel{\Delta}{=}$

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

378    $\land\ replicationSet[i] \neq Server$

380    $\land\ \text{LET}\ in\ \ \stackrel{\Delta}{=}\ \text{CHOOSE}\ x \in replicationSet[i] : x \neq i$

381    $out\ \stackrel{\Delta}{=}\ \text{CHOOSE}\ x \in Server \setminus replicationSet[i] : \text{TRUE}$

382    $\text{IN}\ \ \ \ replicationSet'\ \ \ \ =$

383    $[replicationSet\ \text{EXCEPT}\ ![i] = (@ \setminus \{in\}) \cup \{out\}]$

384    $\land\ currentSubTerm' = [currentSubTerm\ \text{EXCEPT}\ ![i] = @ + 1]$

385    $\land\ \text{UNCHANGED}\ \langle messages,\ serverVars,\ candidateVars,\ nextIndex,$

386    $matchIndex,\ witnessSubterm,\ log,\ commitIndex,$

387    $witnessVars\rangle$

389 ├─────────────────────────────────────────────────────────────────┤

395    $HandleRequestVoteRequest(i,\ j,\ m)\ \stackrel{\Delta}{=}$

396    $\text{LET}\ logOk\ \stackrel{\Delta}{=}\ \lor\ m.mlastLogTerm > LastTerm(log[i])$

397    $\lor\ \land\ m.mlastLogTerm = LastTerm(log[i])$

398    $\land\ m.mlastLogIndex \geq Len(log[i])$

399    $grant\ \stackrel{\Delta}{=}\ \land\ m.mterm = currentTerm[i]$

400    $\land\ logOk$

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

402    $\text{IN}\ \ \ \land\ m.mterm \leq currentTerm[i]$

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

404    $\lor\ \neg grant \land \text{UNCHANGED}\ votedFor$

405    $\land\ Reply([mtype\ \ \ \ \ \ \ \ \ \ \ \mapsto RequestVoteResponse,$

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

407    $mvoteGranted \mapsto grant,$

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

411    $msource\ \ \ \ \ \ \mapsto i,$

412    $mdest\ \ \ \ \ \ \ \ \ \mapsto j],$

413    $m)$

414    $\land\ \text{UNCHANGED}\ \langle state,\ currentTerm,\ candidateVars,\ leaderVars,\ logVars,\ witnessVars\rangle$

416

417

418  $HandleRequestWitnessVoteRequest(j, m) \triangleq$

419      LET $logOk \triangleq \lor m.mlastLogTerm > witnessLastLogTerm$

420                          $\lor \land m.mlastLogTerm = witnessLastLogTerm$

421                              $\land m.mlastLogSubterm > witnessLastLogSubterm$

422                          $\lor \land m.mlastLogTerm = witnessLastLogTerm$

423                              $\land m.mlastLogSubterm = witnessLastLogSubterm$

424                              $\land m.mvotesGranted \subseteq witnessReplicationSet$

425          $grant \triangleq \land m.mterm = currentTerm[WitnessID]$

426                   $\land logOk$

427                   $\land votedFor[WitnessID] \in \{Nil, j\}$

428      IN   $\land m.mterm \leq currentTerm[WitnessID]$

429            $\land \lor grant \quad \land votedFor' = [votedFor \text{ EXCEPT } ![WitnessID] = j]$

430                $\lor \neg grant \land \text{UNCHANGED } votedFor$

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

432                         $mterm \qquad\qquad \mapsto currentTerm[WitnessID],$

433                         $mvoteGranted \mapsto grant,$

434                            

435                            

436                         $mlog \qquad\qquad \mapsto log[WitnessID],$

437                         $msource \qquad \mapsto WitnessID,$

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

439                         $m)$

440            $\land \text{UNCHANGED } \langle state, currentTerm, candidateVars, leaderVars,$

441                         $logVars, witnessVars\rangle$

443

444

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

446      

447      

448      $\land m.mterm = currentTerm[i]$

449      $\land votesResponded' = [votesResponded \text{ EXCEPT } ![i] =$

450                           $votesResponded[i] \cup \{j\}]$

451      $\land \lor \land m.mvoteGranted$

452          $\land votesGranted' = [votesGranted \text{ EXCEPT } ![i] =$

29

453          $votesGranted[i] \cup \{j\}]$

454       $\wedge\ voterLog' = [voterLog\ \text{EXCEPT}\ ![i] =$

455          $voterLog[i] @@ (j :> m.mlog)]$

456     $\vee\ \wedge\ \neg m.mvoteGranted$

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

458     $\wedge\ Discard(m)$

459     $\wedge\ \text{UNCHANGED}\ \langle serverVars,\ votedFor,\ leaderVars,\ logVars,\ witnessVars \rangle$

461 Server $i$ receives an $AppendEntries$ request from server $j$ with

462 $m.mterm \leq currentTerm[i]$. This just handles $m.entries$ of length 0 or 1, but

463 implementations could safely accept more by treating them the same as

464 multiple independent requests of 1 entry.

465 $HandleAppendEntriesRequest(i,\ j,\ m)\ \triangleq$

466     $\text{LET}\ logOk\ \triangleq\ \vee\ m.mprevLogIndex = 0$

467            $\vee\ \wedge\ m.mprevLogIndex > 0$

468              $\wedge\ m.mprevLogIndex \leq Len(log[i])$

469              $\wedge\ m.mprevLogTerm = log[i][m.mprevLogIndex].term$

470    $\text{IN}\ \ \ \wedge\ m.mterm \leq currentTerm[i]$

471       $\wedge\ \vee\ \wedge$   reject request

472          $\vee\ m.mterm < currentTerm[i]$

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

474           $\wedge\ state[i] = Follower$

475           $\wedge\ \neg logOk$

476        $\wedge\ Reply([mtype$            $\mapsto AppendEntriesResponse,$

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

478             $msuccess$         $\mapsto \text{FALSE},$

479             $mmatchIndex$    $\mapsto 0,$

480             $msource$          $\mapsto i,$

481             $mdest$            $\mapsto j],$

482             $m)$

483        $\wedge\ \text{UNCHANGED}\ \langle serverVars,\ logVars \rangle$

484      $\vee$   return to follower state

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

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

487        $\wedge\ state' = [state\ \text{EXCEPT}\ ![i] = Follower]$

488        $\wedge\ \text{UNCHANGED}\ \langle currentTerm,\ votedFor,\ logVars,\ messages \rangle$

489      $\vee$   accept request

490                $\land\ m.mterm = currentTerm[i]$

491                $\land\ state[i] = Follower$

492                $\land\ logOk$

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

494          $\text{IN}$     $\lor$    already done with request

495                         $\land\ \lor\ m.mentries = \langle\rangle$

496                           $\lor\ \land\ m.mentries \neq \langle\rangle$

497                              $\land\ Len(log[i]) \geq index$

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

499                         This could make our *commitIndex* decrease (for

500                         example if we process an old, duplicated request),

501                         but that doesn't really affect anything.

502                         $\land\ commitIndex' = [commitIndex\ \text{EXCEPT}\ ![i] =$

503                                         $m.mcommitIndex]$

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

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

506                                 $msuccess \qquad\quad\ \mapsto \text{TRUE},$

507                                 $mmatchIndex \qquad \mapsto m.mprevLogIndex +$

508                                             $Len(m.mentries),$

509                                 $msource \qquad\qquad \mapsto i,$

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

511                                 $m)$

512                         $\land\ \text{UNCHANGED}\ \langle serverVars, log\rangle$

513                 $\lor$    conflict: remove 1 entry

514                         $\land\ m.mentries \neq \langle\rangle$

515                         $\land\ Len(log[i]) \geq index$

516                         $\land\ log[i][index].term \quad \neq m.mentries[1].term$

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

518                                        $log[i][index2]]$

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

520                         $\land\ \text{UNCHANGED}\ \langle serverVars, commitIndex, messages\rangle$

521                 $\lor$    no conflict: append entry

522                         $\land\ m.mentries \neq \langle\rangle$

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

524                         $\land\ log' = [log\ \text{EXCEPT}\ ![i] =$

525                                   $Append(log[i], m.mentries[1])]$

526                         $\land\ \text{UNCHANGED}\ \langle serverVars, commitIndex, messages\rangle$

31

527            $\land$ UNCHANGED $\langle candidateVars,\ leaderVars,\ witnessVars \rangle$

529    Witness receives an $AppendEntriesToWitnessRequest$ from server $j$ with

530    $m.mterm \leq currentTerm[WitnessID]$.

531   $HandleAppendEntriesToWitnessRequest(j,\ m) \triangleq$

532      $\land\ m.mterm \leq currentTerm[WitnessID]$

533      $\land\ \lor$   reject request

534          $\land\ m.mterm < currentTerm[WitnessID]$

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

536                    $mterm \qquad\qquad\ \mapsto currentTerm[WitnessID],$

537                    $msuccess \qquad\ \mapsto$ FALSE,

538                    $mmatchIndex \quad\ \mapsto 0,$

539                    $msource \qquad\quad \mapsto WitnessID,$

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

541                    $m)$

542          $\land$ UNCHANGED $\langle serverVars,\ logVars,\ witnessVars \rangle$

543        $\lor$   accept request, always no conflict.

544          $\land\ m.mterm = currentTerm[WitnessID]$

545          $\land\ \lor\ m.mlastLogTerm > witnessLastLogTerm[WitnessID]$

546            $\lor\ \land\ m.mlastLogTerm = witnessLastLogTerm[WitnessID]$

547               $\land\ m.mlastLogSubterm > witnessLastLogSubterm[WitnessID]$

548          $\land\ witnessReplicationSet' =$

549            $[witnessReplicationSet$ EXCEPT $![WitnessID] = m.mreplicationSet]$

550          $\land\ witnessLastLogTerm' =$

551            $[witnessLastLogTerm$ EXCEPT $![WitnessID] = m.mlastLogTerm]$

552          $\land\ witnessLastLogSubterm' =$

553            $[witnessLastLogSubterm$ EXCEPT $![WitnessID] = m.mlastLogSubterm]$

554            $log$ will not be modified in real implementation. It is only

555            used here for the proof.

556          $\land\ log' = [log$ EXCEPT $![WitnessID] = m.mentries]$

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

558                    $mterm \qquad\qquad\ \mapsto currentTerm[WitnessID],$

559                    $msuccess \qquad\ \mapsto$ TRUE,

560                    $mmatchIndex \quad\ \mapsto m.mindex,$

561                    $msource \qquad\quad \mapsto WitnessID,$

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

563                    $m)$

564          $\land$ UNCHANGED $\langle serverVars,\ log \rangle$

565      $\land$ UNCHANGED $\langle candidateVars,\ leaderVars \rangle$

567     Server $i$ receives an $AppendEntries$ response from server $j$ with

568     $m.mterm = currentTerm[i]$.

569 $HandleAppendEntriesResponse(i,\ j,\ m)\ \triangleq$

570      $\land\ m.mterm = currentTerm[i]$

571      $\land\ \lor\ \land\ m.msuccess$   successful

572          $\land\ nextIndex'\ \ = [nextIndex\ \ \text{EXCEPT}\ ![i][j]\ = m.mmatchIndex + 1]$

573          $\land\ matchIndex' = [matchIndex\ \text{EXCEPT}\ ![i][j] = m.mmatchIndex]$

574       $\lor\ \land\ \neg m.msuccess$   not successful

575          $\land\ nextIndex' = [nextIndex\ \text{EXCEPT}\ ![i][j] =$

576                     $Max(\{nextIndex[i][j] - 1,\ 1\})]$

577          $\land$ UNCHANGED $\langle matchIndex \rangle$

578      $\land\ Discard(m)$

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

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

582 $UpdateTerm(i,\ j,\ m)\ \triangleq$

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

584      $\land\ currentTerm'\ \ \ = [currentTerm\ \text{EXCEPT}\ ![i] = m.mterm]$

585      $\land\ state'\ \ \ \ \ \ \ \ \ \ \ \ = [state\ \ \ \ \ \ \ \ \text{EXCEPT}\ ![i]\ \ \ = Follower]$

586      $\land\ votedFor'\ \ \ \ \ \ \ = [votedFor\ \ \ \ \ \text{EXCEPT}\ ![i]\ \ = Nil]$

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

588      $\land$ UNCHANGED $\langle messages,\ candidateVars,\ leaderVars,\ logVars \rangle$

590     Responses with stale terms are ignored.

591 $DropStaleResponse(i,\ j,\ m)\ \triangleq$

592      $\land\ m.mterm < currentTerm[i]$

593      $\land\ Discard(m)$

594      $\land$ UNCHANGED $\langle serverVars,\ candidateVars,\ leaderVars,\ logVars \rangle$

596     Receive a message.

597 $Receive(m)\ \triangleq$

598      LET $i\ \triangleq\ m.mdest$

599         $j\ \triangleq\ m.msource$

600      IN      Any $RPC$ with a newer term causes the recipient to advance

601             its term first. Responses with stale terms are ignored.

```
602              ∨ UpdateTerm(i, j, m)
603              ∨ ∧ m.mtype = RequestVoteRequest
604                ∧ HandleRequestVoteRequest(i, j, m)
605              ∨ ∧ m.mtype = RequestWitnessVoteRequest
606                ∧ HandleRequestWitnessVoteRequest(j, m)
607              ∨ ∧ m.mtype = RequestVoteResponse
608                ∧ ∨ DropStaleResponse(i, j, m)
609                  ∨ HandleRequestVoteResponse(i, j, m)
610              ∨ ∧ m.mtype = AppendEntriesRequest
611                ∧ HandleAppendEntriesRequest(i, j, m)
612              ∨ ∧ m.mtype = AppendEntriesToWitnessReques
613                ∧ HandleAppendEntriesToWitnessRequest(j, m)
614              ∨ ∧ m.mtype = AppendEntriesResponse
615                ∧ ∨ DropStaleResponse(i, j, m)
616                  ∨ HandleAppendEntriesResponse(i, j, m)

618     End of message handlers.
619 ├──────────────────────────────────────────────────────────────────
620     Network state transitions

622     The network duplicates a message
623     DuplicateMessage(m) ≜
624         ∧ Send(m)
625         ∧ UNCHANGED ⟨serverVars, candidateVars, leaderVars, logVars⟩

627     The network drops a message
628     DropMessage(m) ≜
629         ∧ Discard(m)
630         ∧ UNCHANGED ⟨serverVars, candidateVars, leaderVars, logVars⟩

632 ├──────────────────────────────────────────────────────────────────
633     Defines how the variables may transition.
634     Next ≜  ∧ ∨ ∃ i ∈ Server : Restart(i)
635              ∨ ∃ i ∈ Server : Timeout(i)
636              ∨ ∃ i, j ∈ Server : RequestVote(i, j)
637              ∨ ∃ i ∈ Server : BecomeLeader(i)
638              ∨ ∃ i ∈ Server, v ∈ Value : ClientRequest(i, v)
639              ∨ ∃ i ∈ Server : AdvanceCommitIndex(i)
```

34

| | |
|---|---|
| 640 | $\lor \exists\, i,\, j \in Server : AppendEntries(i,\, j)$ |
| 641 | $\lor \exists\, i\ \ \in Server : AdjustReplicationSet(i)$ |
| 642 | $\lor \exists\, i\ \ \in Server : RequestWitnessVote(i)$ |
| 643 | $\lor \exists\, i\ \ \in Server : AppendEntriesToWitness(i)$ |
| 644 | $\lor \exists\, m \in \text{DOMAIN } messages : Receive(m)$ |
| 645 | $\lor \exists\, m \in \text{DOMAIN } messages : DuplicateMessage(m)$ |
| 646 | $\lor \exists\, m \in \text{DOMAIN } messages : DropMessage(m)$ |
| 647 | History variable that tracks every *log* ever: |
| 648 | $\land\ allLogs' = allLogs \cup \{log[i] : i \in Server\}$ |
| | |
| 650 | The specification must start with the initial state and transition according |
| 651 | to *Next*. |
| 652 | $Spec \;\triangleq\; Init \land \Box[Next]_{vars}$ |
| | |
| 654 | |

## 6  Formal proof

In this section, we aim to prove the key properties of the Raft algorithm (namely, Election Safety, Leader Append-Only, Log Matching, Leader Completeness, and State Machine Safety) within the context of the extended Raft algorithm. These properties will be proven to hold true consistently in the extended Raft algorithm. Our proof will build upon the existing proof of the Raft algorithm, and we will reference lemmas from the Raft proof in the format of "Lemma R.n," where 'n' represents the lemma's index in the Raft proof.

**Lemma 1.** *Following lemmas are true in the extended Raft algorithm.*

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

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

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

4. *Lemma R.4. An $\langle < \rangle index,\, term >$ pair identifies a log prefix.*

5. *Lemma R.5. When a follower appends an entry to its log, its log after the append is a prefix of the leader's log at the time the leader sent the AppendEntries request.*

35

6. *Lemma R.6. A server's current term is always at least as large as the terms in its log.*

7. *Lemma R.7. The terms of entries grow monotonically in each log.*

*Proof.*    1. Witness, similar to other regular servers, monotonically increases its term upon receiving a message with a larger *mterm* (action *UpdateTerm*). This immediately validates Lemma R.1 as per the specification.

2. In the extended Raft algorithm, leader is elected by a quorum. Neither regular servers nor the witness vote for different servers within the same term, thus maintaining the truth of Lemma R.2 according to its proof.

3. Leader in the extended Raft algorithm manages its log identically to the Raft algorithm, validating Lemma R.3 since

   (a) Actions *BecomeLeader*, *ClientRequest*, and *AppendEntries* behave same as those in the Raft algorithm. Although the log entry is associated with the current subterm, this is trivial to the proof.

   (b) *AppendEntriesToWitness* is an implementation of the *AppendEntries* action, with batching log prefix to a log entry (that is associated with the current subterm and has received subquorum acknowledgments from the current replication set).

4. Lemma R.4 and Lemma R.5 are valid in the extended Raft algorithm given that:

   (a) Regular servers replicate logs identically to the Raft algorithm. The existence of a subterm in a log entry is trivial to the proof.

   (b) Leader always sends its log's prefix to the witness.

   (c) Witness replaces its log with the log prefix from the leader, hence $log'[WitnessID]$ remains a prefix of $m.mlog$ for $m$ sent to the witness.

5. The validation of Lemma R.6 in the Raft algorithm relies on the servers' behavior in log replication. Therefore, Lemma R.6 holds true in the extended Raft algorithm as:

   (a) Regular servers replicate logs in an identical manner.

   (b) Witness accepts the request with the necessary condition $currentTerm[WitnessID] = m.mterm$.

(c) Witness always replaces its log with the log prefix sent from the leader.

6. The proof of Lemma R.7 depends on Lemma R.5 and Lemma R.6. Thus, Lemma R.7 is also valid in the extended Raft algorithm, as both Lemma R.5 and Lemma R.6 are proven true.

$\square$

Now we have Election Safety property (Lemma R.2), Leader Append-Only property (Lemma R.3), and Log Matching property (Lemma R.4) hold in the extended Raft algorithm.

Similarly, we have following lemmas specific for the extended Raft algorithm.

**Lemma 2.** *Each leader's currentSubterm monotonically increases during each term:*

$$\forall i \in Server :$$
$$currentTerm[i] = currentTerm'[i]$$
$$\implies currentSubterm[i] \leq currentSubterm'[i]$$

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

**Lemma 3.** *The subterms of entries grow monotonically during its termin each log*

$$\forall l \in allLogs :$$
$$\forall index \in 1..(Len(l) - 1) :$$
$$l[index].term = l[index + 1].term$$
$$\implies l[index].subterm \leq l[index + 1].subterm$$

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

2. Inductive step: logs change in one of the following ways:

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

   i. The new entry's subterm is *currentSubterm(leader)*.

   ii. *currentSubterm* monotonically increases during *currentTerm*.

   iii. Thus, the new entry's subterm is at least as larger as the subterms during *currentTerm* in this log, since leader's log monotonically grows during its term.

37

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

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

(c) Case: a follower adds one entry (AppendEntries request), or witness replaces its log by log prefix from leader

    i. $log'[follower]$ is a prefix of $m.mlog$ (by Lemma R.5)

    ii. $m.mlog \in allLogs$

    iii. By the inductive hypothesis, the subterms in $m.mlog$ monotonically grow during its term, so the subterms in $log'[follower]$ monotonically grow during its term.

$\square$

*Lemma R.8 is part of the Leader Completeness property in Raft algorithm. Given the existence of shortcut replication, and the fact that witness uses different rules to cast vote, we present proof to the equivalent* **Lemma** *4 following the same idea as that in the Raft algorithm.*

**Lemma 4.** *Immediately committed entries are committed*

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

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

  2. Define

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

  3. Let *election* be an element in *Contradicting* with a minimal *term* field. That is,

$$\forall\ e \in Contradicting : election.eterm \leq e.eterm.$$

If more than one election has the same term, choose the earliest one. (The specification does not allow this to happen, but it is safe for a leader to step down and become leader again in the same term)

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

5. Let $\langle index', term, subterm \rangle$ be any entry that exists in logs of a quorum during $term$, where $index' \leq index$. Such entry must exist if the subterm's replication set contains the witness, since:

   (a) Case $\langle index, term, subterm \rangle$ is immediately committed with no shortcut replication: $\langle index', term, subterm \rangle$ exists in logs of a quorum during $term$, following the specification.

   (b) Case $\langle index, term, subterm \rangle$ is immediately committed with shortcut replication: Following the specification, shortcut replication can only be executed after a log entry in the same subterm being replicated to the witness. Thus $\langle index', term, subterm \rangle$ exists for some $index' < index$.

6. Let $voter$ be either

   - any regular server that both votes in $election$ and contains $\langle index, term, subterm \rangle$ in its log during $term$

   - witness that both votes in $election$ and contains $\langle index', term, subterm \rangle$ in its log during $term$

   Such server must exist, since:

   (a) Case: $\langle index, term, subterm \rangle$ is immediately committed with no shortcut replication.
   $voter$ must exists, since:

       i. $\langle index, term, subterm \rangle$ exists in logs of a quorum during $term$.
       ii. A quorum of servers voted in $election$ to make it succeed.
       iii. Two quorums always overlap.
       iv. $voter$ can be either regular server or witness.

   (b) Case: $\langle index, term, subterm \rangle$ is immediately committed with shortcut replication. $voter$ must exists, since:

       i. A subquorum of regular servers contains $\langle index, term, subterm \rangle$ in its log during $term$. And the witness server contains $\langle index', term, subterm \rangle$ in its log during $term$. These make a quorum.
       ii. A quorum of servers voted in $election$ to make it succeed. The server is either regular server or witness.

      iii. Two quorums always overlap.

7. Let $voterLog \triangleq election.evoterLog[voter]$, the voter's log at the time it casts its vote.

8. For any entry $\langle k, term, subterm \rangle (k \leq index)$ that the voter contains during $term$, it is also contained by the voter when it casts its vote during $election.eterm$. That is, $\langle k, term, subterm \rangle \in voterLog$, where $k \leq index$:

  (a) $\langle k, term, subterm \rangle$ was in the $voter$'s log during $term$.

  (b) The $voter$ must have stored the entry in $term$ before voting in $election.eterm$, since:

     i. $election.eterm > term$.

     ii. The $voter$ rejects requests with terms smaller than its current term, and its current term monotonically increases.

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

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

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

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

       B. The leader of $term$ created and therefore contains the entry.

       C. The leader would not send any conflicting requests to $voter$ during $term$.

     iii. Case: No $AppendEntriesRequest$ with $mterm > term$ removes the entry from the voter's log, since:

       A. Case: $mterm > election.eterm$:
        This can't happen, since $currentTerm[voter] > election.eterm$ would have presented the voter from voting in $term$.

       B. Case: $mterm = election.eterm$:
        Since there is at most one leader per term, this request would have to come from $election.eleader$ as a result of an earlier election in the same term ($election.eterm$).
        Because a leader's log grows monotonically during its term, the leader could not have had $\langle k, term, subterm \rangle$ in its log at

the start of its term, hence it could not have had $\langle index, term, subterm \rangle$ because $index \geq k$. Then there exists an earlier election with the same term in *Contradicting*; this is a contradiction.

   C. Case $mterm < election.eterm$:

The leader of $mterm$ must have contained $\langle index, term, subterm \rangle$ (otherwise its election would also be *Contradicting* but have a smaller term than *election*, which is a contradiction). Thus, the leader of $mterm$ could not send any conflicting entries to the voter for this index, nor could it send any conflicting entries for prior indexes: that it has this entry implies that it has the entire prefix before it (Raft lemma 4).

9. Thus, we have $\langle index, term, subterm \rangle \in voterLog$ for non-witness voter, and $\langle index', term, subterm \rangle \in voterLog$ for witness voter.

We show contradiction in following cases which will then be used to show contradiction in each case of voting rules.

10. Case: $LastTerm(election.elog) = LastTerm(voterLog) \wedge Len(election.elog) \geq index$

   (a) The leader of $LastTerm(voterLog)$ monotonically grew its log during its term (by Lemma R.3).

   (b) The same leader must have had $election.elog$ as its log at some point, since it created the last entry.

   (c) Thus, $voterLog[1..index]$ is a prefix of $election.elog$.

   (d) Then $\langle index, term, subterm \rangle \in election.elog$, since $\langle index, term, subterm \rangle \in voterLog[1..index]$.

   (e) But $election \in Contradicting$ implies that $\langle index, term, subterm \rangle \notin election.elog$.

11. Case: $LastTerm(election.elog) > term$

   (a) $election.eterm > LastTerm(election.elog)$ since servers increment their $currentTerm$ when starting an election, and Lemma R.6 states that a server's $currentTerm$ is at least as large as the terms in its log.

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

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

$$term \leq$$
$$LastTerm(election.elog) = prior.eterm <$$
$$election.eterm$$

(d) $\langle index, term, subterm \rangle \in prior.elog$, since $prior \notin Contradicting$ (*election* was assumed to have the lowest term of any election in *Contradicting*, and $prior.eterm < election.eterm$).

(e) $prior.elog$ is a prefix of $election.elog$ since:

    i. $prior.eleader$ creates entries with $prior.eterm$ by appending them to its log, which monotonicallygrows during $prior.eterm$ from $prior.elog$.

    ii. Thus, any entry with term $prior.eterm$ must follow $prior.elog$ in all logs (by Lemma R.4).

    iii. $LastTerm(election.elog) = prior.eterm$

(f) $\langle index, term, subterm \rangle \in election.elog$. Note that this is true no matter $\langle index, term, subterm \rangle$ existing in *voterLog* or not.

(g) This is a contradiction, since $election.elog$ was assumed to not contain the committed entry ($election \in Contradicting$).

12. The log comparison during elections states the following, since *voter* granted its vote during *election*:

$$\lor \land voter \neq WitnessID$$
$$\land \lor LastTerm(election.elog) > LastTerm(voterLog)$$
$$\lor \land LastTerm(election.elog) = LastTerm(voterLog)$$
$$\land Len(election.elog) \geq Len(voterLog)$$
$$\lor \land voter = WitnessID$$
$$\land \lor LastTerm(election.elog) > LastTerm(voterLog)$$
$$\lor \land LastTerm(election.elog) = LastTerm(voterLog)$$
$$\land LastSubterm(election.elog) > LastSubterm(voterLog)$$
$$\lor \land LastTerm(election.elog) = LastTerm(voterLog)$$
$$\land LastSubterm(election.elog) = LastSubterm(voterLog)$$
$$\land election.evoters \setminus WitnessID \subset ReplicationSet$$

$ReplicationSet$ is replication set of $LastSubterm(voterLog)$ in $LastTerm(voterLog)$ if $voter$ is witness.

We now investigate each voting rule in following cases, and use the result of above two cases to show contradiction.

13. Case: $voter$ is a regular server, and $LastTerm(election.elog) = LastTerm(voterLog) \wedge Len(election.elog) \geq Len(voterLog)$

    (a) $Len(election.elog) \geq Len(voterLog) \implies Len(election.elog) \geq index$

    (b) This leads to contradiction following case 10

14. Case: $LastTerm(election.elog) > LastTerm(voterLog)$, no matter $voter$ being a regular server or witness

    (a) $LastTerm(election.elog) > LastTerm(voterLog) \implies LastTerm(election.elog) > term$

    (b) This leads to contradiction following case 11

15. Case: $voter$ is witness, and $LastTerm(election.elog) = LastTerm(voterLog) \wedge LastSubterm(election.elog) > LastSubterm(voterLog)$

    (a) $LastSubterm(election.elog) > LastSubterm(voterLog) \implies Len(election.elog) > Len(voterLog) > term$, following Lemma 3

    (b) This leads to contradiction following case 10

16. Case: $voter$ is witness, and

    $$\wedge\ LastTerm(election.elog) = LastTerm(voterLog)$$
    $$\wedge\ LastSubterm(election.elog) = LastSubterm(voterLog)$$
    $$\wedge\ election.evoters \cup WitnessID \subset ReplicationSet$$

    (a) $LastTerm(voterLog) = term \wedge LastSubterm(voterLog) = subterm$

        i. There must be $voter'$ which is a regular server that contains $\langle index, term, subterm \rangle$ and votes in $election$, since

            A. $\langle index, term, subterm \rangle$ exists in at least subquorum regular servers in replication set of $term$ and $subterm$, no matter if it is shortcut replicated or not.

            B. A subquorum in replication set of $term$ and $subterm$ vote in $election$.

43

C. Subquorums of regular servers overlap in a replication set.

    ii. This leads to contradiction following case 13 and 14

(b) $LastTerm(voterLog) = term \wedge LastSubterm(voterLog) > subterm$

    i. $LastTerm(voterLog) = term \wedge LastSubterm(voterLog) > subterm \implies Len(voterLog) > index$, per Lemma 3

    ii. This leads to contradiction following case 10

(c) $LastTerm(voterLog) > term$. This leads to contradiction following case 11

Having Lemma 4 proved, Lemma R.9 and Theorem R.1 are true, following the same proof in Raft algorithm. Now we have Leader Completeness and State Machine Safety proved for the extended Raft algorithm. □