

Raft Protocol Summary

Followers

- Respond to RPCs from candidates and leaders.
- Convert to candidate if election timeout elapses without either:
 - Receiving valid AppendEntries RPC, or
 - Granting vote to candidate

Candidates

- Increment currentTerm, vote for self
- Reset election timeout
- Send RequestVote RPCs to all other servers, wait for either:
 - Votes received from majority of servers: become leader
 - AppendEntries RPC received from new leader: step down
- Election timeout elapses without election resolution: increment term, start new election
- Discover higher term: step down

Leaders

- Initialize nextIndex for each to last log index + 1
- Send initial empty AppendEntries RPCs (heartbeat) to each follower; repeat during idle periods to prevent election timeouts
- Accept commands from clients, append new entries to local log
- Whenever last log index \geq nextIndex for a follower, send AppendEntries RPC with log entries starting at nextIndex, update nextIndex if successful
- If AppendEntries fails because of log inconsistency, decrement nextIndex and retry
- Mark log entries committed if stored on a majority of servers and at least one entry from current term is stored on a majority of servers
- Step down if currentTerm changes

Persistent State

Each server persists the following to stable storage synchronously before responding to RPCs:

<u>currentTerm</u>	latest term server has seen (initialized to 0 on first boot)
<u>votedFor</u>	<u>candidateId</u> that received vote in current term (or null if none)
<u>log[]</u>	log entries

Log Entry

<u>term</u>	term when entry was received by leader
<u>index</u>	position of entry in the log
<u>command</u>	command for state machine

RequestVote RPC

Invoked by candidates to gather votes.

Arguments:

<u>candidateId</u>	candidate requesting vote
<u>term</u>	candidate's term
<u>lastLogIndex</u>	index of candidate's last log entry
<u>lastLogTerm</u>	term of candidate's last log entry

Results:

<u>term</u>	<u>currentTerm</u> , for candidate to update itself
<u>voteGranted</u>	true means candidate received vote

Implementation:

- If term > currentTerm, currentTerm \leftarrow term (step down if leader or candidate)
- If term == currentTerm, votedFor is null or candidateId, and candidate's log is at least as complete as local log, grant vote and reset election timeout

AppendEntries RPC

Invoked by leader to replicate log entries and discover inconsistencies; also used as heartbeat.

Arguments:

<u>term</u>	leader's term
<u>leaderId</u>	so follower can redirect clients
<u>prevLogIndex</u>	index of log entry immediately preceding new ones
<u>prevLogTerm</u>	term of <u>prevLogIndex</u> entry
<u>entries[]</u>	log entries to store (empty for heartbeat)
<u>commitIndex</u>	last entry known to be committed

Results:

<u>term</u>	<u>currentTerm</u> , for leader to update itself
<u>success</u>	true if follower contained entry matching <u>prevLogIndex</u> and <u>prevLogTerm</u>

Implementation:

- Return if term < currentTerm
- If term > currentTerm, currentTerm \leftarrow term
- If candidate or leader, step down
- Reset election timeout
- Return failure if log doesn't contain an entry at prevLogIndex whose term matches prevLogTerm
- If existing entries conflict with new entries, delete all existing entries starting with first conflicting entry
- Append any new entries not already in the log
- Advance state machine with newly committed entries