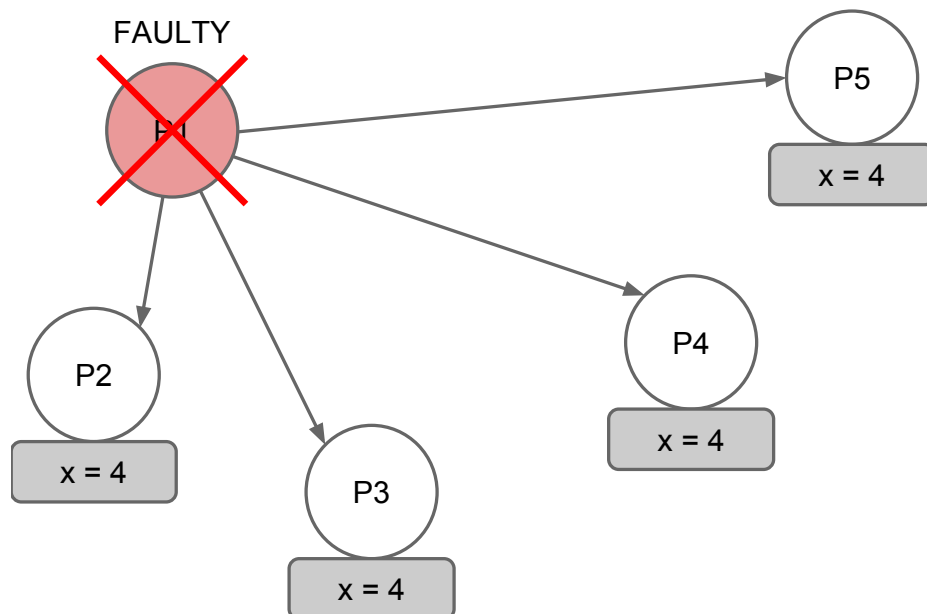


FINAL REPORT

RAFT CONSENSUS ALGORITHM



Joachim Kirkegaard Friis s093256
Anders Nielsen s103457

December 9, 2014

Preface

This report is for the final project in the course 02228 Fault-Tolerant Systems of fall 2014. The purpose of the project is to describe and evaluate the fault tolerance of a given system. In our project we have chosen to survey and implement a recently proposed solution to the consensus problem within distributed systems - Raft. Our approach to this is to implement a tool that applies Raft to a simulated distributed network and from this derive describe how the algorithm works and it's ways of providing fault tolerance. The problem about consensus will be presented after which Raft will be described in this context. The design and implementation will follow with a discussion on our experience through each step of the development. A conclusion will then summarize our result and experiences.

Contents

1. Introduction	1
1.1. Problem	1
1.2. Motivation	2
1.2.1. Two Generals Problem	2
1.2.2. Properties of Distributed Consensus	3
2. Raft	3
2.1. Components of Raft	4
3. Requirements	4
3.1. Functional requirements	5
3.2. Quality requirements	5
4. Problem Analysis	5
4.1. Existing Raft Implementations	6
4.2. Testing	6
4.3. Development Approach	6
4.4. Visualization Tool	7
5. Leader Election	7
5.1. Leader election in Raft	7
5.2. Implementation	10
5.2.1. Election Timer	10
5.3. RequestVote RPC	10
6. Log Replication	11
6.1. Log Replication in Raft	11
6.1.1. AppendEntries RPC/Heartbeat	12
6.1.2. Committing Log Entry	13
6.2. Handling Log Inconsistency	13
6.3. Implementation	13
7. Evaluation and Testing	14
7.1. White box tests	14
7.1.1. Tests	15
7.2. Black box test scenarios	15
7.2.1. Tests	15
8. Results	17
8.1. Running the simulator	18
8.2. Visualization	19
8.3. Source Code	19

9. Discussion	19
9.1. Requirements evaluation	20
9.2. Reflection on goal	20
9.3. Experiences	21
9.4. Future work	21
10. Conclusion	22
A. Results	24
B. Acceptance tests	26

1. Introduction

This is the final report for the project in the course “Fault tolerant systems” (02228). The purpose of this project is to learn document our experience in implementing Raft - a consensus algorithm. Our approach to this is use the specification of the Raft paper [5, p. 4] and iteratively construct a test for each behaviour the algorithm should have and implement its functionality that satisfies the given test.

In this introductory section we will present the fundamental problem when talking about consensus in a distributed system with a basic description of the Raft algorithm.

1.1. Problem

Reaching consensus in a distributed system means that all processes in the network eventually agree on some value or state of the entire system. This is often needed when processes might be faulty thus bringing reliability and availability at stake on single-point-of-failure. Upholding these properties in a given distributed system then relies on the architecture utilised and hardware implemented in the processes.

A simple solution to this could be to initiate a vote among all correct processes on to what the value is in which the value is the result of the majority vote and thereby mask the faulty processes. The figure 1 below illustrates an example of a distributed system consisting of a number of nodes. The value x is what the system must agree upon. Though here we have a faulty process $P1$ which is responsible to return the result. Taking the fact of connectivity of system aside, the system will become unavailable because of the now faulty process.

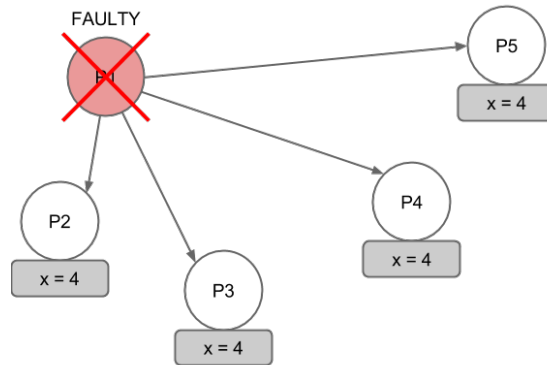


Figure 1: A distributed system consisting of a number of nodes with a faulty one.

The fundamental problem behind this is that you cannot rely on the individual processes to be reliable and thus solely store the value. The system needs more This boils down to a well-known problem - The Two Generals Problem.

1.2. Motivation

The elements of the problem presented serve as motivation behind consensus in a system of unreliable components. Because how do you rely on a value in a system where some components will eventually be faulty? The basic problem of reaching consensus in a distributed network can be illustrated by the Two Generals Problem analogy.

1.2.1. Two Generals Problem

Figure 2 is an illustration of two generals from the same army who want to attack an enemy army and they have to attack at the same time in order to win. They cannot communicate directly to each other since they are at different fronts of the battlefield. In the context of a distributed system the generals here can be seen as two processes trying to agree on a value. General 1 sends out a messenger to tell the second general that they should attack at dawn. The second general then receives this message, but the first general cannot be sure of this (the messenger might be captured or killed by the enemy on his way to the second general and vice versa). Again, in the context of distributed system, the unreliability of messenger can directly related back to the unreliability of message transmission in a normal distributed system [4].

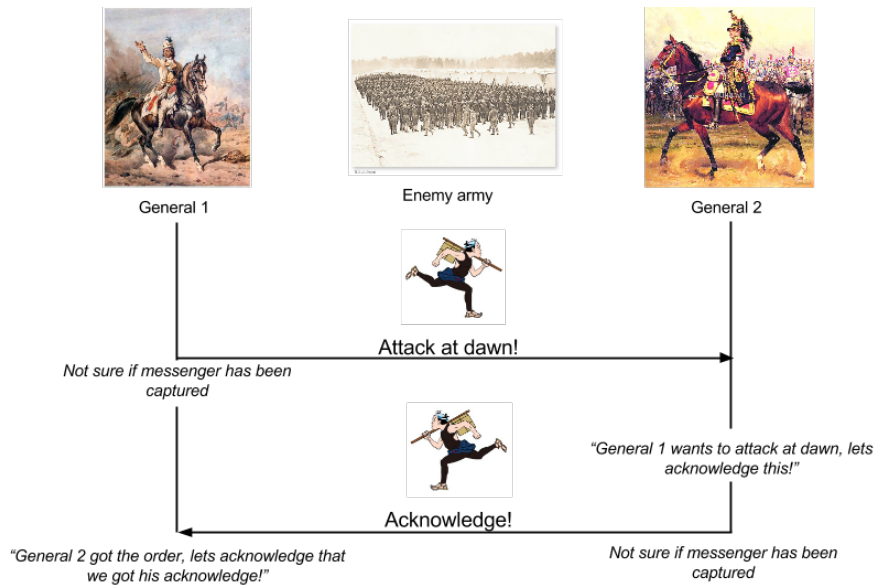


Figure 2: Two generals tries to agree on when to attack the enemy by sending a messenger, but they are not sure the messenger survives his trip between their camps.

Now the second general sends the messenger back in order to acknowledge this, but the first general also has to acknowledge this, resulting in a never ending run for the poor messenger - thus the generals can never agree on when to attack the enemy.

1.2.2. Properties of Distributed Consensus

When want to define consensus, the goal is to satisfy a set of requirements i.e. properties that the distributed system must uphold. These are used to describe the systems fault tolerant features related to faulty processes. A faulty process can either fail by crashing or experience a Byzantine failure. Such a failure in the context of distributes system occur when e.g. a process for some reason transmits incorrect or malicious data throughout the network. Due to the arbitrary results of these kinds of failures, properties of a system are often distinguished by either tolerating them or not. A main difference in terms of properties of a system whether it tolerates Byzantine failures or not, is the validity and integrity. The integrity property for a system that does not tolerate Byzantine failures is as following [2].

Byzantine failure can be related two the two generals problem in terms of the messenger being captured by the enemy and turned to spy on the generals i.e. given false information.

- non-Byzantine failure tolerant:
 - **Validity:** If all processes propose the same value v , then all correct processes decide v .
 - **Integrity:** Every correct process decides at most one value, and if it decides some value v , then v must have been proposed by some process.
- Byzantine failure tolerant:
 - **Validity:** If all correct processes propose the same value v , then all correct processes decide v .
 - **Integrity:** If a correct process decides v , then v must have been proposed by some correct process.

The clearest commonality here, is that you should always be able to say that all correct processes must be able to decide on the same value or state, as mentioned earlier. Though the main difference is that with Byzantine processes, you must be able to say if they are all correct before stating they can derive the same value. This fact differentiates many solutions to this problem, depending on which property set they can satisfy.

Also, as discussed in [5] a consensus algorithm for a non-Byzantine system has the following properties: safety, availability, timing interdependency, and majority vote on procedure calls. The safety property can be directly related back to our validity and integrity properties, as the system's safety and availability is measured be the correctness of return result upon a request.

2. Raft

As of writing this report, the most recent proposal to solving the consensus problem is Raft [5]. Diego Ongaro and John Ousterhout argue that most consensus algorithms, such as Paxos [3] suffer from poor understandability and are hard to understand and

implement. Diego and John introduce Raft as a simple and understandable solution to the consensus problem.

The authors of Raft argue that Raft is superior in both the aspects of educational uses but also in performance gained from a less complex implementation. Their main motivation behind Raft is their point of view on the frequently used consensus algorithm Paxos [3]. They describe Paxos as being simple by its foundation but limited which has resulted in more variations to solve more complex consensus scenarios. Another arguments they propose is that Paxos might well work in theory, but it is hard to integrate in a real system because of its too simple foundation i.e. “simple decree” decomposition (reaching agreement on a single decision) [5].

2.1. Components of Raft

Raft is described as having three main components where the first two describes the behaviour of algorithm and the last describes how the algorithm ensures consensus safety properties.

- **Leader election:** A strong leader is elected and responsible for keeping the rest of the system in consensus for a term that ends if it fails.
- **Log replication:** The leader must accept log entries from clients and replicate them across the cluster and forcing the other servers to agree with its own log [5].
- **Safety:** The specified design provides safety such that [5]:
 - **Election Safety:** There is always at most one leader in the system.
 - **Leader Append-Only:** A leader can only append new entries to its log i.e. it cannot overwrite existing entries.
 - **Log Matching:** The logs of all servers are matching up to their latest common index and term.
 - **Leader Completeness:** No newly elected leaders will overwrite other server’s log if its own log entry is not up to date.
 - **State Machine Safety:** When a server updates its log it must be up to date with the leader’s log.

The above description of Raft serves as a rough overview of the algorithm. A more detailed description will be given during the implementation in section 5 and 6 where we will describe the components in more details together with the design and implementation choices we have made during the implementation.

3. Requirements

The tool we are implementing should serve as an illustration tool to show how the Raft algorithm works for given scenarios of faulty processes in a distributed system. The requirements are split into two different aspects i.e. functional and quality attributes:

3.1. Functional requirements

The end product must satisfy the following functional requirements given the MoSCoW method:

1. The tool must illustrate a scenario of a given set of processes in a distributed system, in which a leader is elected.
2. The tool must be an implementation of the Raft algorithm as specified in [5].
3. The user should be able to input parameters to vary the scenario, such as the number of processes.
4. The user must then be able to disconnect any process.
 - a) If a leader is disconnected a new leader must be elected automatically.
5. The user must be able to request that a log entry is replicated to the system at run time.
6. The tool could be further extended with a visualisation of the given distributed system.

3.2. Quality requirements

The quality requirements of the illustration tool are inherited from the properties that are provided by the Raft algorithm.

1. The Raft implementation should have the following properties: Safety, availability, timing independence, and that commands complete as soon as a majority has responded.
 - a) Safety: The properties described in section 2 about Raft.
 - b) Availability: The system should be available as long it is possible to elect a leader through a majority vote.
 - c) Timing independence: Safety should not depend on the timing of the system, i.e. the speed of a given process should provide it with an advantage in terms of the consensus of the system. For this, it is specified that:

$$broadcastTime \ll electionTimeout \ll MTBF [5].$$

4. Problem Analysis

In this section we will analyze and describe our approach on implementing Raft with relation to the purpose of this project and with relation to achieve a more fault-tolerant and robust implementation.

4.1. Existing Raft Implementations

The application of consensus algorithms such as Raft is many, but the most generic applications are databases and service discovery systems. As Raft mostly fits applications in the low level end of the stack, it will also be best suited to be implemented in a more low level and light language such as C, C++ or Go. Erlang would also be a great fit with its fault-tolerant features.

Raft have already been implemented in several versions and in many languages. Diego Ongaro himself has implemented Raft in C++ (LogCabin¹), CoreOS² uses etcd³, which is a Raft implementation in Go used for service discovery.

If a Raft implementation was needed for real usage, the solution would have been to extend some of the rich open source implementations already available in solid languages as Go and C. But as mentioned in the introduction, the scope of this project is to learn about Raft with relation to fault-tolerance by implementing it. To be able to cover as much of Raft as possible in the time scope, the language chosen for implementation is JavaScript, because it would allow more productivity.

There already exists solutions in JavaScript, but since one of the focus of this project is to learn about Raft and the consensus problem, we have chosen to implement a new version with our own approach. Already existings solutions in JavaScript such as **raftscope** and **skiff** have worked as inspirations during implementation.

4.2. Testing

In order to be able to verify that the implementation satisfies the requirements and that to ensure a higher level of robustness, the software will be thoroughly tested. To document requirement satisfaction, the implementation will be black-box tested on the level of acceptance and integration.

In order to achieve robustness the implementation will also be white-box tested on the unit test level in order to test the behaviour of the different objects in situations and cases that are hard to foresee.

4.3. Development Approach

The Raft paper includes an informal specification [5, page 4] with behaviour of the servers described by the state they are in or how to respond to requests. In order to have a more common language between the specification and the tests, the development approach Behaviour Driven Development (BDD) was chosen⁴ [6]. The Behaviour Driven approach is iterative and have been used in the report with the following steps:

1. Derive behaviour from the specification and write a failing test of the given behaviour.

¹<https://github.com/logcabin/logcabin>

²A new Linux operation system fork made for server deployments.

³<https://github.com/coreos/etcd>

⁴Originally described by Dan North in the blog post <http://dannorth.net/introducing-bdd/>

2. Write the simplest code that implements the behaviour and make the test pass.
3. Refactor the code
4. Repeat

This is close to Test Driven Development, but with focus on writing tests that describes behaviour instead of testing methods as in unit tests. One important aspect to make this approach work well is to derive the behaviour that seem most easy to implement first.

To facilitate writing tests with a BDD-approach, we used the node libraries `chai` together with `mocha`, which allows us to write tests written as behaviour such as the example seen in code example below.

Listing 1: Example of a JavaScript test with the libraries mocha and chai

```
describe("Server", function() {
  it("responds with false if log is outdated", function() {
    var outDatedLog = new Log();
    var server = new Server(outDatedLog);
    assert(server.getResponse(), false);
  });
});
```

4.4. Visualization Tool

An extended requirement specified in section 3 was to build a visualization of the simulation tool. To be able to use the Raft implementation such that the visualization tool does not affect other extensions it is important to build it independent from each other as seen on figure 3. As visualization platform we have chosen to use the command line to visualize the simulation since it is the simplest one.

5. Leader Election

In this section we will describe the component of leader election in Raft followed by a discussion of the design and implementation choices made during implementation.

5.1. Leader election in Raft

As mentioned in section 2 one of the two main components of Raft is its leader election. After a successful leader election the given leader is responsible of replicating commands it receives from a client to the rest of the followers, which is explained in the log replication section 6. All servers contain an implementation of Raft, having a log of commands requested by clients. A server can be in three states: “leader”, “candidate” or “follower”. The servers convert to different states in different scenarios defined by the Raft specification and illustrated in figure 4. As default the servers in the system will be in the “follower” state in which they act as an active part receiving heartbeat messages

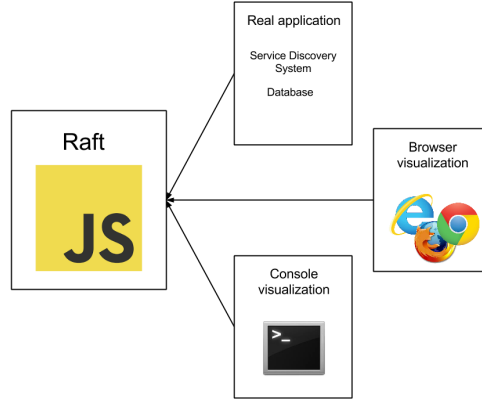


Figure 3: Illustrating how Raft should be implemented such that it is extendable for uses such as command line tools, browser visualizations or a real usage such as databases or service discovery systems.

from the current leader. At most one leader exists at any given election term. The election term is incremented each time a new leader is elected. In distributed systems the term is the systems logical clock in which any entry in the server’s log has a term attached to it in order to know when that value has been replicated to.

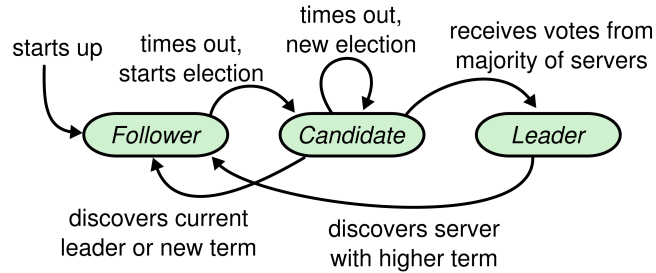


Figure 4: A graph illustrating the three different states of Raft servers and the situations in which they convert to other states. Figure taken from the Raft paper [5].

Figure 5 illustrates the scenario in which a leader is elected:

In the beginning of the illustration, on initial system startup, every server is a follower until $P1$ times out, because it has not received any valid heart beat messages from a leader. When $P1$ times out it becomes a “candidate”. If this candidate receives a majority of votes it transitions to the leader state after which it has the responsibilities of a leader as described in section 2. The leader claims its leadership by sending out heartbeats, where if a server receives a heartbeat request from another server it is instantly converted to follower and recognizes the source as leader.

Figure 6 shows how Raft will tolerate a faulty leader. Here we see that the leader $P1$ fails by crashing. $P4$ then times out and initiates a new election.

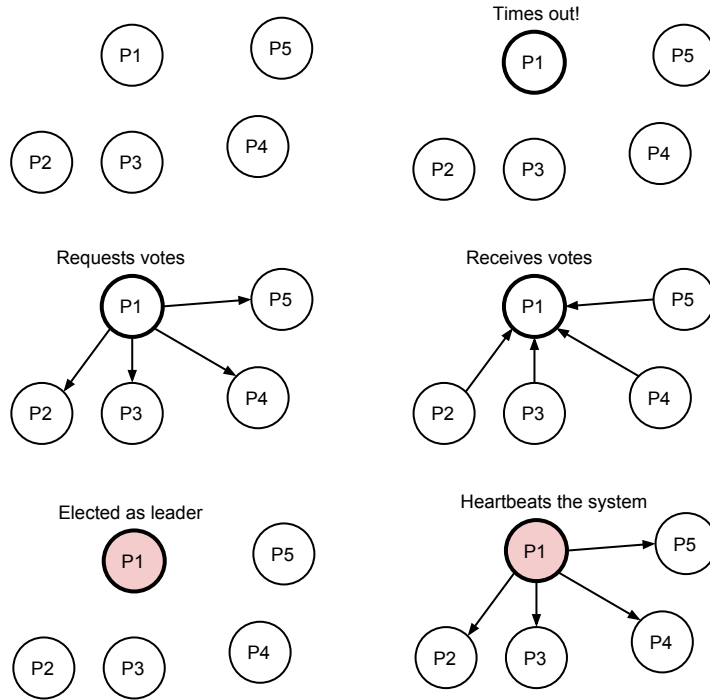


Figure 5: A simple scenario where 5 servers are to obtain consensus by first electing a leader.

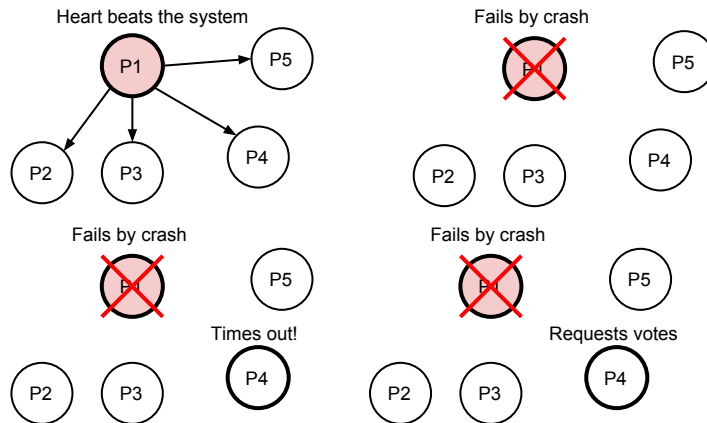


Figure 6: A simple scenario where a leader $P1$ crashes after which $P4$ times out and start a new election.

It should also be noted that upon transitioning to the candidate state and requesting votes, another server might have been requesting votes before this. In this situation the candidate would receive a heartbeat message from a leader with a term at least as large

as its own, after which it transitions back to follower. The situation in which a majority vote is not achieved by any candidate, the candidates will issue yet another election once they time out again. In order to avoid tied election such as these Raft uses a randomized time out within a given interval.

It should also be noted that since the availability, provided by Raft, relies on the majority votes and only non-Byzantine failures, the algorithm can thus only uphold these properties in a network where $3F \leq N$, where F is the amount of failures and N is the amount of processes in the system. [1]

Looking back at the Two Generals Problem, Raft is only able to solve the problem if a majority vote would be possible which cannot be the case with only two “servers” (or generals). In order to solve this problem when utilising Raft, one must add another general to the scenario such that there are more than two servers in the system. We should also modify the problem such that the messenger always tells the truth (i.e. cannot suffer from Byzantine failures). Now that we have three generals and an “uncapturable” messenger, they will reach consensus by finding a leader after which he decides when to launch the attack.

5.2. Implementation

Our Raft implementation consists of two major parts:

1. An election timer that is reset with a random time in a specified interval
2. RequestVote RPC implementation with methods for handling invocation and receiving RequestVote RPC request and response

5.2.1. Election Timer

The election timer is implemented by keeping state of the amount of milliseconds left and decrement this value every millisecond using the JavaScript function `setInterval`. The timer is then reset in five situations:

- On initialization
- On restart
- When receiving an AppendEntries RPC
- When it has started an election
- When it has voted for another candidate

5.3. RequestVote RPC

If the election timer reaches zero it will start an election by sending a Remote Procedure Call (RPC) request called *RequestVote RPC*, where each server responds whether or not

it will grant a vote to the candidate. As the votes are independent to each server they can be sent out in parallel as specified in the Raft paper. Since our implementation is a simulation with objects instead of processes or servers and JavaScript does not have threads, it is not possible to do true parallel computing. But since it is a simulation, we simulate parallel requests with the JavaScript method `setTimeout`, which delays the execution of a given function. To make the simulation more like real requests, they are delayed given by a random amount of milliseconds in a given interval.

This implementation is fine for the visualization, but since we also want to verify the implementation through tests, it is not convenient to make the election random since you cannot assert a random election result. It would be possible to assert a given property such as “assertion that there is at most one leader after x seconds”. Although this could be a viable test it is not convenient either, since it will make the tests slow and thereby slow down the development flow.

In order to easily test the execution of the RequestVoteRPC, we have made a protocol for communicating between servers in Raft. The protocol is an object that delegates a method call to a target and is injectable in the Server object. The delayed implementation explained above is seen in the protocol object `DirectAsync`. The protocol used in testing will call methods directly on the target object and is called `Direct`.

6. Log Replication

In this section we will describe how Raft handles log replication, discuss some of the design decisions we have taken within Raft and describe the challenges with the implementation.

6.1. Log Replication in Raft

The general purpose of the log replication feature in a consensus algorithm and in Raft is to keep a log history of commands applied to the state machine in order to compare integrity between servers.

Log entries are appended to the log when a client requests a given command to be applied. The flow of receiving the request from the client, committing the command to the log and responding the client is illustrated in figure 7. In overall the steps of committing a command in Raft can be written as:

1. Client sends request to the leader with command
2. Leader appends command to own log
3. Leader replicates log to followers in parallel
4. When the command is replicated and committed to a majority of the servers in the cluster, the leader will respond positively to the client and apply the command to the state machine.

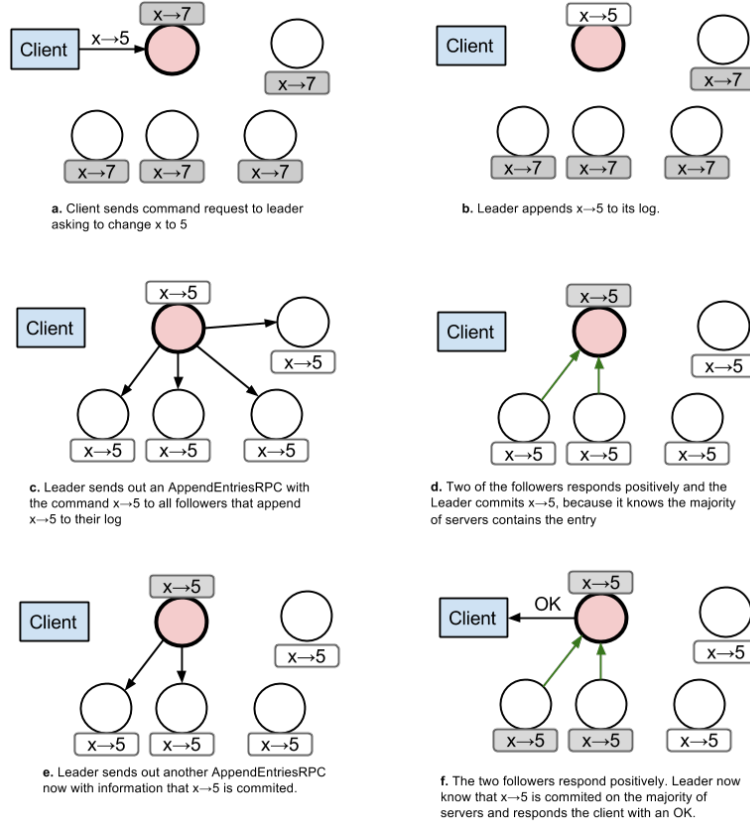


Figure 7: An illustration of how the leader successfully applies a command from the client. At first (a) the client requests the leader to apply the command $x \rightarrow 5$ to the state machine, then (b) it applies the command to its own log and (c) replicates it to the followers. (d) The followers respond positively if the command is applied to their log and the leader commits the entry from the client, when a majority of the servers have appended the log. (e) The followers will receive information that the entry is committed and (f) when the leader knows that a majority of the servers have applied the command, it will respond to the client positively and append it to the state machine

6.1.1. AppendEntries RPC/Heartbeat

The log replication mechanism explained in these steps are in Raft the same heartbeat mechanism used for claiming leadership as explained in section 5. The heartbeat/log replication request is a Remote Procedure Call (RPC) called *AppendEntries RPC*. The leader invokes the RPC with information for the follower to duplicate the leaders log and be in consensus. The follower then responds the leader with success if it is fully replicated or else failure.

6.1.2. Committing Log Entry

The first time a server receives an `AppendEntries` RPC it will just append it to the log. A log entry is committed by the leader when the leader knows that a majority of the servers have appended the given log entry to their log. The log entry is then committed on a follower, when the leader sends out information (through *`AppendEntries RPC`*) that the index of a given log entry is committed. The log entry is safe to apply to the state machine, when a majority of the servers have committed the entry.

6.2. Handling Log Inconsistency

The scenario explained this far is the successful one, but in distributed systems failures such as network partitions and single server crashes are not rare and if this failure happens in the middle of the steps in figure 7 it could result in inconsistent logs throughout the system. Raft has behaviour for handling different kinds of log inconsistencies, which compared to the true log of the leader can be:

1. A follower is missing entries
2. A follower has extra entries
3. A follower is both missing entries and has extra entries

All situations of inconsistencies are handled by forcing the followers to duplicate the leaders log. The first situation is handled through the `AppendEntries` RPC by finding the largest common committed index and replicate the missing entries from the leader. The second scenario is handled by finding the largest common committed index and delete all entries following. Thereafter the situation is as the first and handled so. The last scenario is handled as the second.

6.3. Implementation

The log and heartbeat are two implementation specific parts of the Raft algorithm that is not described in implementation details.

Normally the log should be persistent, but since the implementation in this project is only for visualization, persistency is not implemented. The log is implemented as an object `Log` that works as an abstraction on top of a simple array containing the entries. For a programmer this abstraction is important since Raft is documented and specified as being a 1-indexed system, where most databases and arrays in various programming languages are 0-indexed.

Besides replicating its log, the leader also need to keep state of the progress of the follower's logs. More precisely it need to know the next empty index of a given follower's log and the index of largest log entry that matches in the leader log. In our implementation this state is kept in the `LeaderState` object that is initialized (and reset) when a new leader is elected.

The heartbeat is implemented and used by the leader as an independent mechanism that sends an `AppendEntriesRPC` every t milliseconds using the JavaScript function `setInterval`. Alternatively the heartbeat could be dependent on the response of the follower by sending a heartbeat i seconds after the leader got a response, but this would introduce extra complexity as we would need to handle the case where the follower is down.

The `AppendEntriesRPC` is as the `RequestVoteRPC` explained in section 5.3 done in parallel. In order to both be able to test the behaviour involving `AppendEntriesRPC` and simulate real requests in the simulation we have added `AppendEntries` to the `Direct` and `DirectAsync` protocols.

An important implementation detail of the log replication is that the log entries need to contain the index and the action of appending a log entry need to be idempotent as specified in the paper [5, p. 10]. Our initial implementation was by using the array index of the log to control indexes, but this meant that if the delay of the RPCs was close to the heartbeat interval the leader would send two `AppendEntriesRPCs` to the same active follower without getting the result of the first `AppendEntriesRPC` which results in a race condition that makes the follower append the same entries twice. To avoid this the leader assigns an incremental index to the log entry when received by a client and secondly the method of appending a log entry is made idempotent. The code for the idempotent log append method is seen in code example 2.

Listing 2: The method used by followers to append entries from the leader to the log. This method is idempotent, which means that appending the same entry to the exact same log twice (or more) will give the same result.

```
/**
 * Appends an entry to the log if it is not already appended.
 * @param {object} logEntry - An array of or single log entry
 */
Log.prototype.appendEntry = function(logEntry) {
  if (logEntry.index > this.lastIndex()) {
    this.logEntries.push(logEntry);
  }
};
```

7. Evaluation and Testing

As mentioned in the analysis section4 we are to evaluate our solution with white- and black box tests. Their description will thus be provided and followed by the results.

7.1. White box tests

The development of our solution was to iterate based on the behaviour of Raft. This included writing a test for each behaviour described in the condensed summary of Raft [5].

For each of these tests its required functionality was then implemented in order to satisfy it.

7.1.1. Tests

Although the functional tests, constructed during implementation, can be directly read in the code, a small example will be presented for completeness's sake.

The following code snippet shows the test for the AppendEntries RPC attribute that a receiver should reply false if the given RPC call contains a term \neq currentTerm.

Listing 3: The test for the AppendEntries RPC Receiver implementation

```
describe("Receiver AppendEntries Implementation", function() {  
  // Rule 1  
  it("replies false if term  $\neq$  currentTerm", function() {  
    var server1 = new Server(2, 'leader');  
    var server2 = new Server(3, 'candidate');  
    server2.currentTerm = 2;  
    updatePeers([server1, server2]);  
    var result = server1.invokeAppendEntries(server2.id);  
    assert.equal(result.success, false);  
  });  
});
```

7.2. Black box test scenarios

Verifying that our solution actually is an implementation of Raft requires more in-depth tests in the form of validating that we uphold the properties that are argued to provide consensus. This is done by constructing a set of scenarios that reflect the behaviour of the components of Raft: Leader Election, Log Replication, and Safety. The scenarios are as following with success criteria and a simple use case:

7.2.1. Tests

- **Leader Election:**

1. A leader is elected upon initiate start up of the system.
 - a) The user starts the program.
 - b) Server S times out.
 - c) S converts to candidate.
 - d) S receives a majority vote and converts to leader.
2. A new leader must be elected if the current leader fails.
 - a) The user starts the program.
 - b) A server L is elected as leader.

- c) The user crashes L at run-time.
- d) Server S times out and starts a new election.

• **Log Replication:**

1. Log entry is replicated through the system.
 - a) The user starts the program.
 - b) A server L is elected as leader.
 - c) The user requests that a value x is replicated.
 - d) L replicates x such that each log of each server contains x in their log at the correct term and index.
2. Awoken follower get its log updated.
 - a) The user starts the program.
 - b) A server L is elected as leader.
 - c) The user crashes a server S .
 - d) The user requests that a value y is replicated.
 - e) The user wakes up S .
 - f) L updates S 's log such that it contains y in its log at the correct index and term.

An example of a leader election tests looks as following:

Listing 4: The test for the leader election

```
describe("LeaderElection: the system", function() {
  it("elects a new leader on initialization", function(done) {
    var simulator = raft.buildCluster({
      clusterSize: 5,
      heartbeatDelay: 1,
      electionTimerInterval: [150, 300],
      protocol: new Direct(null, {})
    });
    setTimeout(function() {
      assertOneLeader(simulator.servers);
      done();
    }, 1000);
  });
});
```

Here a cluster of 5 servers is set up in which the servers talk directly to each other without delay hence giving more reliable behaviour. After a short delay of 1 second the it is then assert that a leader should have been elected. The `assertOneLeader` function simply checks for all servers in the cluster that only one of them is in the leader state.

In appendix B a list of acceptance tests for the rest of the components of Raft are proposed, which are not implemented and run in our solution. Please note that there is no acceptance test for state machine safety. This is because there's is no actual state machine in our implementation, thus testing this would be impossible. Further details about this in section 9.

And referring back to requirements in section 3. The solution should also satisfy the requirements about availability and timing independence. The success criteria and scenarios for these requirements are as listed in appendix B.

The last scenarios about safety, availability, and especially timing independence are hard to formally verify since the program simulates a concurrent system, in which the utilisation of a formal verification tool would be necessary, which is out of scope of this project. A more detailed discussion of this will be done in section 9.

8. Results

In this section the results of the tests described in section 7 will be discussed. These results should provide knowledge about the correctness of the implementation in relation to the requirements specified in section 3.

As mentioned, the white box tests were derived from Raft's behaviour in its condensed summary in the paper. The output of these tests is listed in appendix 8. The black box tests, however, were only done for a given set of scenarios in order to show that the fault tolerant features of the implementation reflects those described in the paper.

8.1. Running the simulator

On top of the tests, a simulator that visualizes the implementation was done in order to get illustrative feature to show how Raft solves the consensus problem in presence of faulty servers. The screenshot in figure 8 belows shows a run of the simulator.

```
----- Running Raft -----
Election timeout: between 3000 and 4500 ms
Heartbeat: every 500 ms
RPC Delay: between 100 and 120 ms
-----
Server 1 (follower) term: 1 logEntries: 0 commitIndex: 0 945
[]
Server 2 (follower) term: 1 logEntries: 0 commitIndex: 0 1032
[]
Server 3 (follower) term: 1 logEntries: 0 commitIndex: 0 1833
[]
Server 4 (follower) term: 1 logEntries: 0 commitIndex: 0 1031
[]
Server 5 (follower) term: 1 logEntries: 0 commitIndex: 0 1465
[]
```

Figure 8: Example run of the simulator in the state where a leader has not yet been elected.

The simulator lists the servers by their id, state, the amount of entries in its log, its latest commit index, and its time out which decrements. Server 4 is colored red to illustrates that it has crashed. The screenshot in figure 9 then shows the situation where a server times out and starts an election by transitioning to the candidate state and requesting votes from the other servers.

```
----- Running Raft -----
Election timeout: between 10000 and 15000 ms
Heartbeat: every 4000 ms
RPC Delay: between 1000 and 1500 ms
-----
Server 1 (candidate) term: 1 logEntries: 0 commitIndex: 0 10068
[]
Server 2 (follower) term: 1 logEntries: 0 commitIndex: 0 10943
[]
Server 3 (follower) term: 1 logEntries: 0 commitIndex: 0 11667
[]
Server 4 (follower) term: 1 logEntries: 0 commitIndex: 0 14444 votedFor: 1
[]
Server 5 (follower) term: 1 logEntries: 0 commitIndex: 0 12629 votedFor: 1
[]
```

Figure 9: Example run of the simulator in the state where a server has started an election.

And lastly the screenshot in figure 10 shows the situation in which a crashed server hasn't had its log updated since the leader does not get replies from the heartbeat messages it send to it. In the moment a log entry from the client is appended there is to underline on the log entries in the visualization. On figure 10 all log entries have an underline since they are committed.

```

----- Running Raft -----
Election timeout: between 1500 and 3000 ms
Heartbeat: every 200 ms
RPC Delay: between 10 and 20 ms

Server 1 (follower) term: 1 logEntries: 1 commitIndex: 1 -52319
[v->x->5, t->1]

Server 2 (follower) term: 1 logEntries: 2 commitIndex: 2 2288
[v->x->5, t->1], [v->y->5, t->1]

Server 3 (leader) term: 1 logEntries: 2 commitIndex: 2 2051
[v->x->5, t->1], [v->y->5, t->1]

Server 4 (follower) term: 1 logEntries: 2 commitIndex: 2 2307
[v->x->5, t->1], [v->y->5, t->1]

Server 5 (follower) term: 1 logEntries: 2 commitIndex: 2 2835
[v->x->5, t->1], [v->y->5, t->1]

```

Figure 10: Example run of the simulator where a leader has been elected and request from a the client has been replicated to all correct servers.

These screenshots of a run of the simulator serves as results of the acceptance tests. Furthermore a video of running the simulator can be seen in the following link:

https://www.dropbox.com/s/75kpa65yykld2mv/simulator_run_video.mov?dl=0.

The output results of the both the white box and black box test are given in appendix A.

8.2. Visualization

The raft implementation, raft simulator and command line visualization are seperated in three reuseable components. This means that a another GUI visualization such as in a browser would be easy as the raw raft implementation and simulator can be reused together with a browser visualization module.

The command line visualization is implemented using an HTTP-server that runs with the raft visualization process. To manipulate with the visualization the raft client sends an HTTP-request, which the command server receives and invokes the relevant method on the visualization process.

8.3. Source Code

The source code has been attached as appendix on CampusNet as a ZIP-file, but it can also be accessed through GitHub at the url:

<http://github.com/anderslime/ft-raft>

9. Discussion

In this section we will look back at the problem introduced in section 1 and evaluate the implementation of Raft. This will be done by comparing our test results with the requirements stated in section 3. Doing this, we can reflect on the overall goal of the project i.e. to gain deeper knowledge of Raft. Further extensions and improvements to the end product will then be proposed and discussed.

9.1. Requirements evaluation

The requirements stated that we should have tool that illustrated a given distributed system with a number of servers, implementing Raft as a solution to finding consensus in the presence of faulty servers. We evaluate each functional requirement by referring to the results in section 8 as following:

1. The tool illustrates a distributed system where the user can issue commands both related to Raft but also for simulation purposes.
2. The tool provides fault tolerance by utilising Raft, in the sense that it will tolerate servers that crash in a non-Byzantine manner.
3. The user can give parameters for: time between heartbeats, how many servers the system should consist of, latency for RPC's, and interval for random election timers.
4. The user can crash a server through the client by issuing the correct command for it in the console.
 - a) The implementation of Raft ensures that a new leader is elected, once a candidate achieves a majority vote.
5. The user can request that an entry is replicated through the system by issuing the correct command for it in the console.
6. The tool visualized the system through a console GUI.

It is hard to verify whether we have met the quality requirement or not. These it due to the lag of static behaviour of the system caused by the randomized election time outs and general concurrent aspects of the system. This is discussed in further detail in the next section.

9.2. Reflection on goal

The goal of this project was to learn about what is the state of the art in terms of solving the consensus problem. In order to do this, our approach was to gain knowledge about Raft through implementing it in a simulated environment. The requirements for the implementation stated that it should provide a tool for visualising the behaviour of Raft.

We divided our tests into white- and black box tests in order to document our development process better by testing the behaviour in the white box tests and come with example tests when verifying the solution itself in the black box tests. But the black box tests do not verify the correctness of the implementation, since they do not cover to whole space of possible scenarios. Another way to verify our implementation would have been to use a formal verification tool such as SPIN, which is used for testing models of concurrent systems. But this would require us to develop a whole new solution in

another language i.e. Promela which conflicts with our main goal, which was to utilize the flexibility and small overhead provided by Javascript.

In terms of the correctness of the implementation of Raft, referring back to the properties for non-Byzantine failure tolerant systems listed in section 1, the choice of platform again provides a poor basis for a formal verification. But it can be argued that by adopting the behaviour specified in the Raft paper thus also the safety properties, that our implementation upholds the properties of validity and integrity in the cases we have chosen to test.

9.3. Experiences

The biggest challenge of implementing the Raft algorithm is to implement and test the situations in parallel (such as mentioned with log replication and RPC's in elections) and the leader election driven by a randomly initialized timer.

In the situations in parallel we have suggested the method of using a protocol supporting both direct and delayed communicating between servers in section 5.3. This made it easy to test the behaviour involving RPC's, but there are challenges in parallel computing that cannot be hidden with abstractions. The first one is race conditions which we experienced in the simulation when the RPC delay is too high in relation to the heartbeat interval. If the leader e.g. replicates a log entry $x \rightarrow 5$ to a follower and before the follower is able to respond to the RPC (e.g. because of latency) a new heartbeat interval is started and the leader blindly tries to replicate the same log entry with $x \rightarrow 5$. This will make the follower append the same entry twice. The situation is illustrated in figure 11. Although we discovered this implementation issue it did not ruin the consensus due to the safety rules of Raft and the problem was fixed by making the log appending method of followers idempotent as described in section 6.3, but it illustrates that the solution of abstracting away parallel operations does not solve the problem of testing fully.

The random election timer used in leader election in Raft as explained in section 5.3 made it hard to test the leader election. The solution for white box tests was to invoke methods to simulate the random timer. To test the leader election with the random timer we used a delay as explained in section 7.2.1. One problem with tests as these is that if the tests fail in a given situation it is impossible to replicate the random situation. Another problem is that if the system testing relies on faults in random situation it is not a system you should trust. The last problem is that slow tests slow down the development flow such that every test will be a tradeoff of much slower test suites.

9.4. Future work

The code for this project is open source and can easily be extended. One interesting extension would be to visualize the server in a graph with visualized RPC-requests to see the communication between the servers and maybe the data send back and forward.

As mentioned it is hard to verify a Raft implementation because some parts of algorithm rely on parallelism, random values and timing. A future project could be about

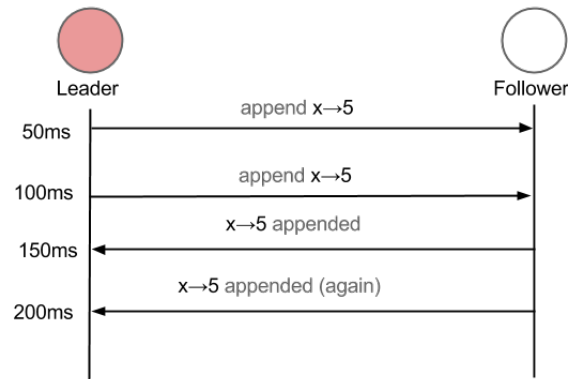


Figure 11: A situation with possible race condition when idempotent log appends is not implemented. This figure illustrates that some situations are only present when in parallel (or asynchronous).

how such components can be verified either formally or by using tests.

A third extension for further work would be to build or use a monitoring tool to measure the availability of the implementation. The simulation tools build in this project could then be exposed to realistic loads and operations in a longer period of time in which the monitoring tool could measure and estimate a Mean time between failures to calculate the availability and easier verify the Timing independence quality requirement specified in section 3.

10. Conclusion

The result of the project is a simulation tool for the command line that simulates a cluster of servers, where the user is able to crash, restart and request the leader to append entries to the log. The Raft simulation implementation ensures that the leader replicates the entry to all other servers and if the server crashes and restarts such that the servers will have consensus and the system is always be available (as long as no more than a third of the servers are active).

One of the focuses was to test and verify the implementation, which has been difficult because parts of Raft depends on parallelism, random values and timing. In order to make these parts more testable we have used specific protocols for testing which made it easier to test, but the abstraction did not succeed in taking the challenges in parallel computing and timing into account.

The implementation, which is in JavaScript, has provided us with a much deeper understanding of the consensus problem and Raft as a solution. Raft was implemented with an iterative approach by constructing a test describing behaviour from the Raft specification followed by an implementation of the specified behaviour that satisfies the test. Practicing this behavioural development process allowed us to structurally

implement Raft without having to give much thought to the overall design and made it easy to refactor the code that was backed by behavioural tests of the whole system.

References

- [1] Michael J. Fischer, Nancy A. Lynch, and Michael S. Paterson. Impossibility of distributed consensus with one faulty process. *J. ACM*, 32(2):374–382, April 1985.
- [2] Tim Kindberg George Coulouris, Jean Dollimore and Gordon Blair. *Disitrbuted Systems, Concepts and Design*.
- [3] Leslie Lamport. The part-time parliament. *ACM Trans. Comput. Syst.*, 16(2):133–169, May 1998.
- [4] Leslie Lamport, Robert Shostak, and Marshall Pease. The byzantine generals problem. *ACM Trans. Program. Lang. Syst.*, 4(3):382–401, July 1982.
- [5] Diego Ongaro and John Ousterhout. In search of an understandable consensus algorithm. In *Proceedings of the 2014 USENIX Conference on USENIX Annual Technical Conference*, USENIX ATC’14, pages 305–320, Berkeley, CA, USA, 2014. USENIX Association.
- [6] C. Solis and Xiaofeng Wang. A study of the characteristics of behaviour driven development. In *Software Engineering and Advanced Applications (SEAA), 2011 37th EUROMICRO Conference on*, pages 383–387, Aug 2011.

A. Results

Rules for Followers

- Respond to RPCs from candidates and leaders
 - responds to RequestVote RPCs from candidate (88ms)
 - responds to AppendEntries RPC from leader

Rules for Candidates

- On conversion to candidate, start election:
 - increments currentTerm
 - votes for self
 - resets election timer
 - sending RequestVote RPCs to all other servers
 - is granted a vote
 - is not granted a vote if vote term is less than target server's term
 - is converted to follower and updates term when outdated
 - is not granted a vote when target server that has already voted for another server
 - is not granted a vote when the candidate log is not up to date
- If votes received from majority of servers, become leader
 - with 5 peers in network
 - becomes leader when receiving 3 votes
 - does not become leader when receiving 2 votes
 - with 4 peers in network
 - becomes leader when receiving 3 votes
 - does not become leader when only receiving 2 votes
- If AppendEntries RPC received from new leader: convert to follower
 - becomes follower when receiving AppendEntries from leader
- Receiver AppendEntries Implementation
 - replies false if term < currentTerm
 - replies true if term >= currentTerm
 - reply false if log doesn't contain entry at prevLogIndex
 - reply false if log doesn't contain entry at prevLogIndex whose terms matches prevLogTerm
 - reply true if log does contain entry at prevLogIndex whose terms matches prevLogTerm
 - appends any new entries not already in the log
 - sets commitIndex=leaderCommit if leaderCommit > commitIndex and leaderCommit < index of last new entry
 - sets commitIndex=index of last new entry if leaderCommit > commitIndex and index of last new entry < leaderCommit
 - sets commitIndex=leaderCommit of last new entry if leaderCommit > commitIndex and index of last new entry == leaderCommit
 - If an existing entry conflicts with a new one
 - deletes existing entry and all that follows

Client Log Entry Request

- client sends log request to server

client sends log request to a follower

General rules for servers

If successful: update nextIndex and matchIndex for follower.

If AppendEntries fails because of log inconsistency: decrement nextIndex and
retry.

If there exists an N such that $N > \text{commitIndex}$ (advance commit index)

updates commitIndex if majority of servers matches that index

does not update commitIndex if majority of servers does not match that index

31 passing (120ms)

B. Acceptance tests

- **Safety:**

1. There is at most one leader at any given term.
 - a) The user starts the program in which two servers S_1 and S_2 have the same election time out.
 - b) A leader L , being the only one with an election time out smaller than both S_1 and S_2 's, gets elected.
 - c) The user crashes L .
 - d) Both S_1 and S_2 start their elections.
 - e) Either S_1 or S_2 is elected as the new leader.
2. Leader never overwrites its own log but only appends on to it.
 - a) The user starts the program.
 - b) A server L is elected as leader.
 - c) The user requests that a value x is replicated.
 - d) The user requests that a value y is replicated.
3. Logs of every server match each other up to their latest common index in their logs.
 - a) The user starts the program.
 - b) A server L is elected as leader.
 - c) The user requests that a value x is replicated.
 - d) The user crashes a server S .
 - e) The user requests that a value y is replicated.
 - f) All the servers have a commonality that they know of x but not y .
4. The log of a leader contains all previous log entries of leader of previous terms.
 - a) The user starts the program with a server L with the lowest election time out and S with the second lowest.
 - b) L is elected as leader.
 - c) The user crashes a server S .
 - d) The user requests that a value x is replicated.
 - e) The user wakes up S
 - f) The user crashes L
 - g) S has an up to date log.

- **Availability:**

1. As long as a majority vote is possible, a leader must always be elected. If not the system should become unavailable.
 - a) The user starts the program with five servers.
 - b) A server L is elected as leader.
 - c) The user crashes L and another server.
 - d) A new leader L' is elected as leader.
 - e) The user crashes L' .
 - f) The system becomes unavailable.