# Raft-based Replicated Block Store
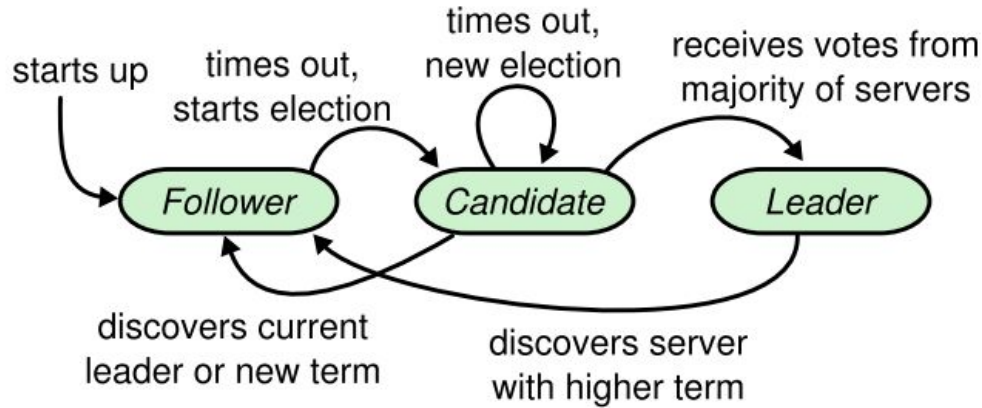
Shichun Yu, Yuting Yan, Zhikang Hao

# Outline

Implementation

- General architecture of Raft consensus algorithm
- Our special decisions and optimizations

Test

- Consistency Test
- Performance Test

# State machine of Raft Server



Persistent States: Term, Vote, Log[]

Volatile States: commitIndex, lastApplied, **leaderID**

Volatile State on Leader: nextIndex[], matchIndex[]

# Log Entry

```
struct entry {
    // 0: read, 1: write
    1: i32 command,
    2: i32 term,
    3: i64 address,
    4: string content,
}
```

In this project, for the 4K block store requests,

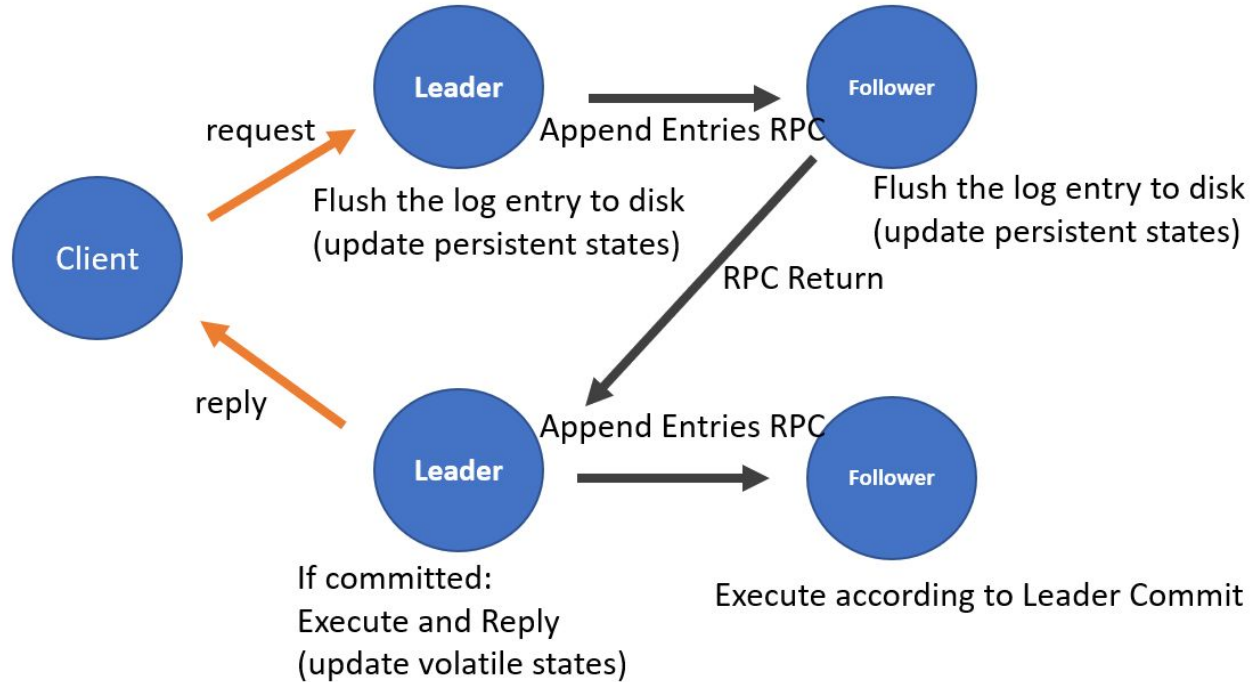the size of content will always be 4K.

# Leader Election

- Candidates send RequestVote RPC

To prevent the livelock, follow the guide of MIT raft project, quote

"only restart your election timer if a) you get an AppendEntries RPC from the current leader (i.e., if the term in the AppendEntries arguments is outdated, you should not reset your timer); b) you are starting an election; or c) you grant a vote to another peer."

Ref: https://thesquareplanet.com/blog/students-guide-to-raft/
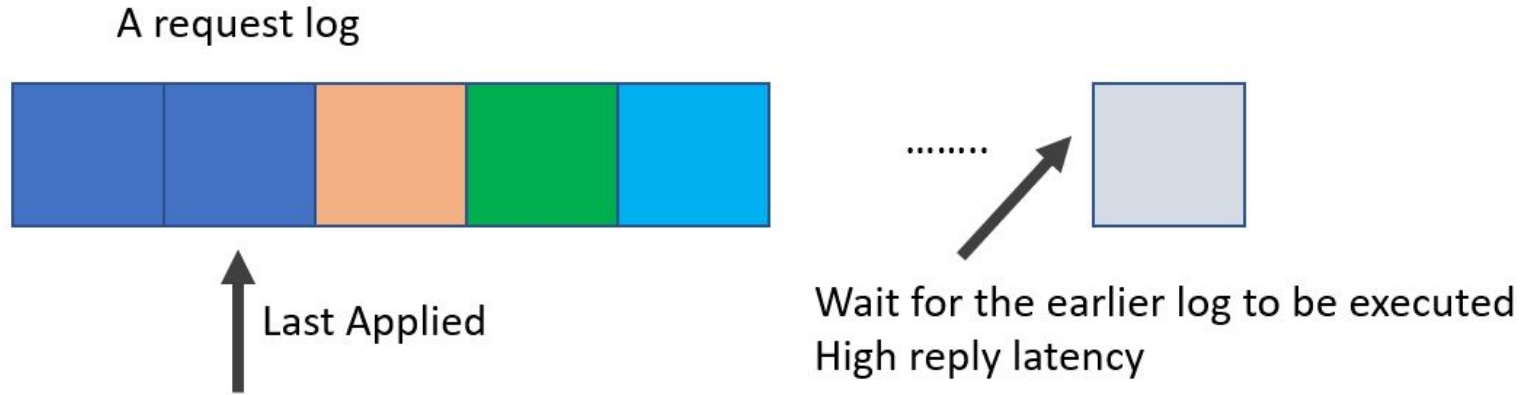
# Normal Operation

# Exactly-once Semantic (not fully implemented)

- Read logs are also replicated through the servers.
- Every request has a sequence number (monotonically non-decreasing for each client).
- Servers keep the completion record (a map from <client id, seq num> to execution result), but the CRs are not persistent.
- The Servers keep track of the largest seq number for each client, and if a request with lower seq num comes, reply immediately without executing twice.     (**only implemented in the leader**)
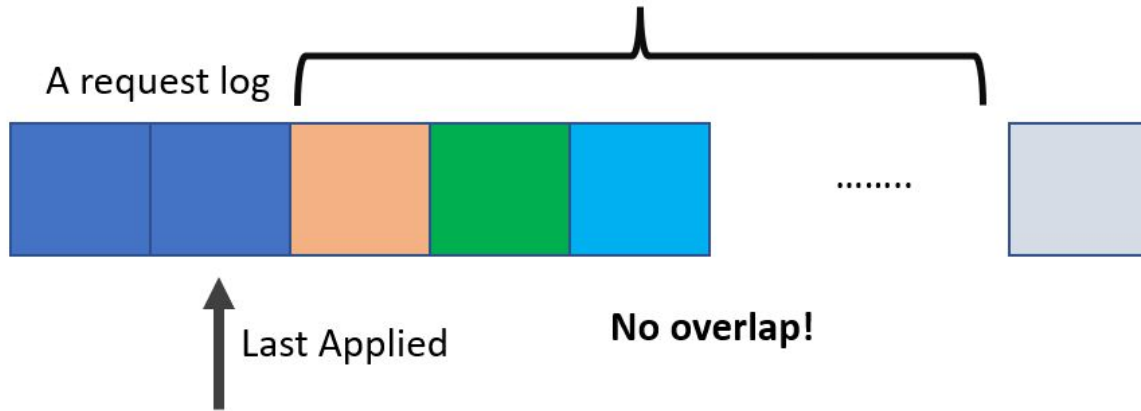
# Performance optimization

How to improve the latency when having large amount of requests?

A request log

# Simple latency optimization

Sometimes there is no need to wait!

Execute and reply immediately if there is no overlap to the requests come earlier.



A request log

Last Applied

No overlap!

Insert the executed request log index to a set.

Update "lastApplied" accroding to this set.

# Test RPCs support

Implement extra RPCs to make debug and test easier.

We implemented Two RPCs

to make comparison of states between leader and followers.

- CompareStates (Log, term, vote)
- CompareBlockStore  (compare the persistent Block Store)

# Correctness Tests

Case 1: follower crash

Proof: resist one crash

& block any modification with 2

crashes (Leader completeness)

Recovered follower will be

Consistent after several seconds

Method: manually crash servers

# Correctness Tests

Case 2: General Leader crash case

Proof: new leader will be elected

Logs in other servers keeps consistent

Recovered server will become follower (we update terms)

Logs are consistent after and state machine is up-to-date

# Correctness Tests

Video

# Correctness Tests

Case 3: Election behavior

Proof: the server with highest

vote priority will be the leader
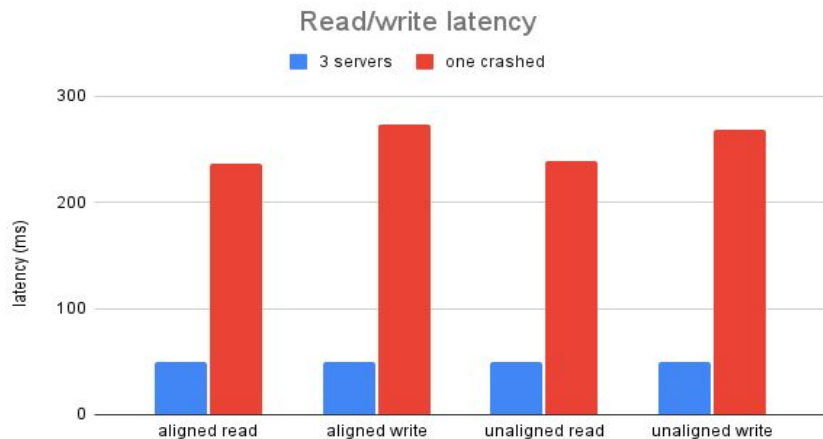
# Performance Tests

Environment:

4 cloudlab c220g1 linux servers, 1Gb cable connection
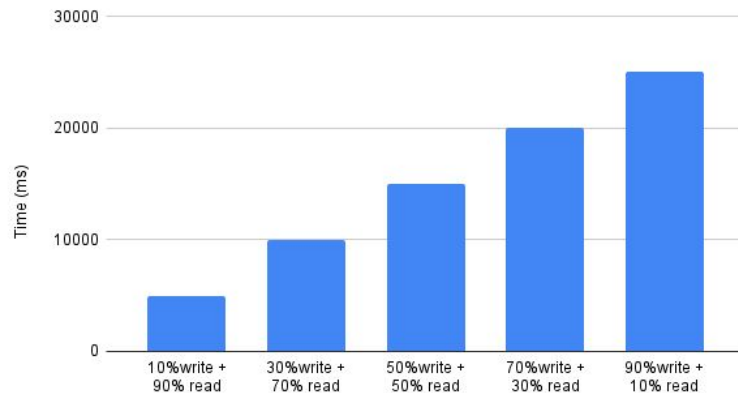
3 distinct raft nodes

1 distinct client node

# Performance Tests

- Read/Write latency, appendEntry interval with 50ms
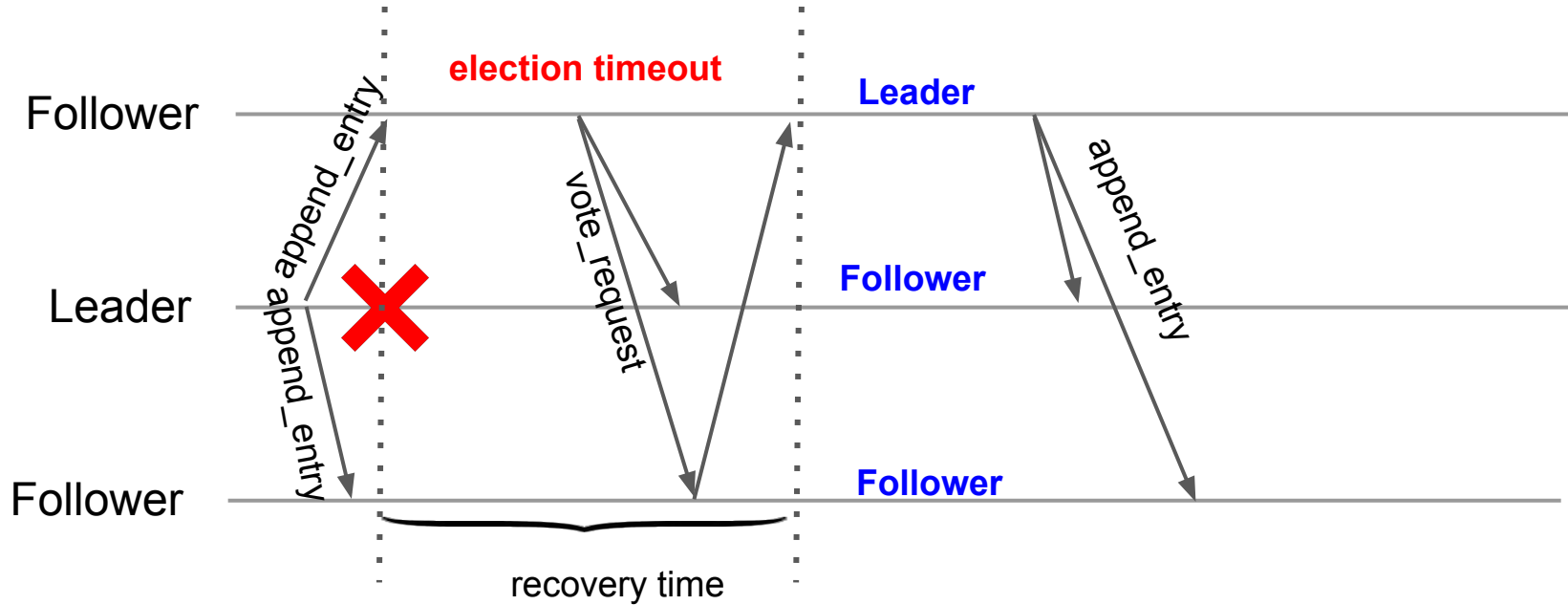


Read/write latency



Latency for 100 runs with different proportion of read and write

- Aligned operation or unaligned operation make no much difference in latency;
- Crashed server in a cluster would increase latency;
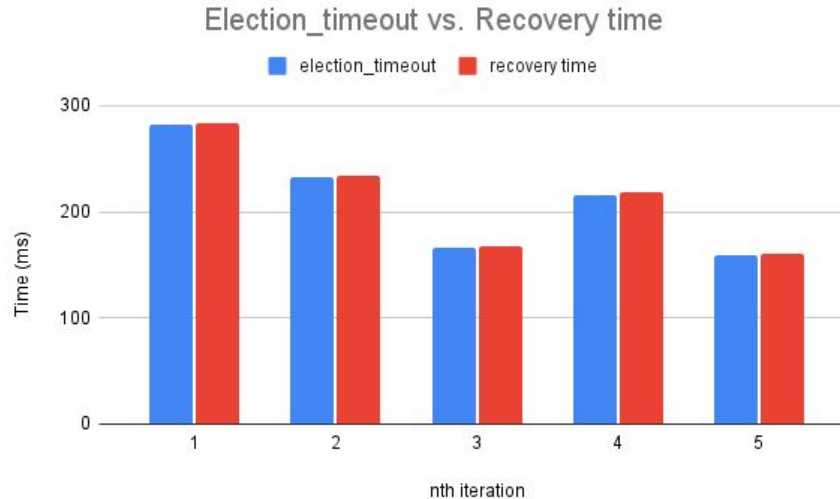- Write latency is larger than read.

# Performance Tests

- Crash Recovery process

# Performance Tests

- Recovery Time



Election_timeout vs. Recovery time

- election_timeout ■ recovery time

● Recovery time is highly related to election timeout.

# Conclusion

- Write latency is larger than read;
- A failed server will impact the performance of system.
- The recovery time is highly related to server's election timeout.

# Thank you!

Github Repo:
https://github.com/SmileIsThinking/Raft-based-Replicated-Block-Store