# Distributed Systems && Go

Connor Zanin

# /usr/bin/whoami

- Connor Zanin
- Live in Boulder
- Work for Newstore, Inc.
- Presented GoMR, 2020


- https://connorzanin.com
- https://github.com/cnnrznn/raft
- https://github.com/cnnrznn/fake-etcd

# Goal

1. Introduce DS concepts and Raft
2. Show how DS + Go == success
3. Demo the system
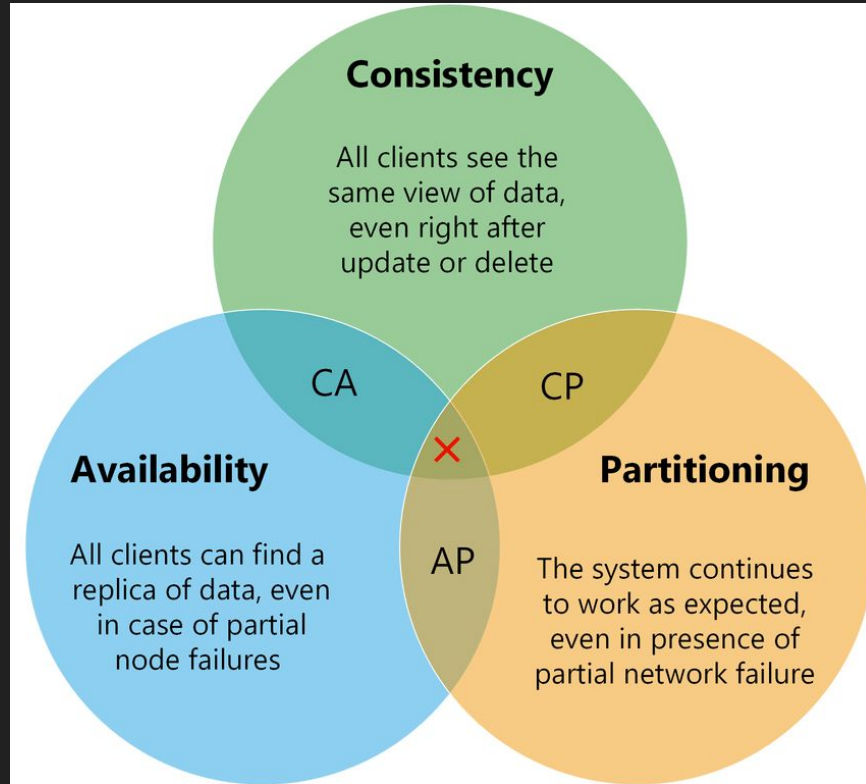
# Distributed Systems

# What is a distributed system?

- Set of machines connected over a network
- Working to some common goal

# Distributed systems problems

- Partial failure
- Fault tolerance + recovery
- Synchronization
    - Clocks
    - State machines
    - Order of events
- Vulnerable to CAP (consistency, availability, partitioning)
- **Consensus**

# CAP Theorem

# Raft High-Level

# What is Raft?

- Consensus protocol
- Replicated Log
- Crash Fault Tolerant (CFT) (fail-stop)
    - All nodes follow the protocol or crash
    - To tolerate $f$ failures, need $2f+1$ nodes
- Leader-Follower protocol (asymmetric)
- "Committed" log entries survive

# Uses

- In-memory key-value stores (etcd)
- Distributed configs (Zookeeper, Consul)
- Pub/sub, message queues
- Distributed file systems (GFS)
- Any in-memory fault-tolerant cache
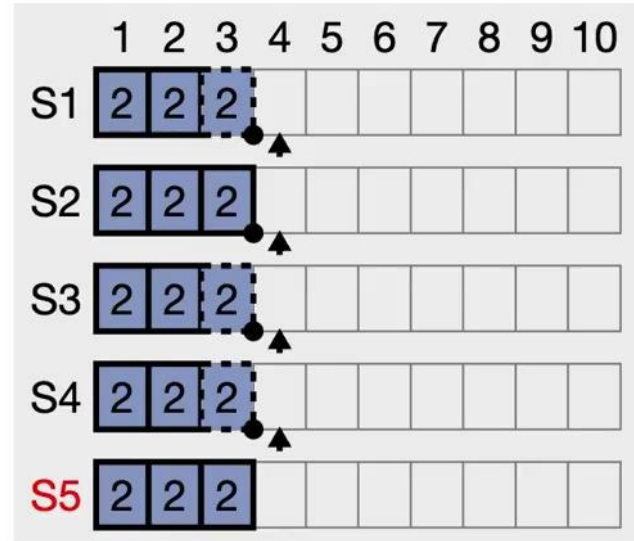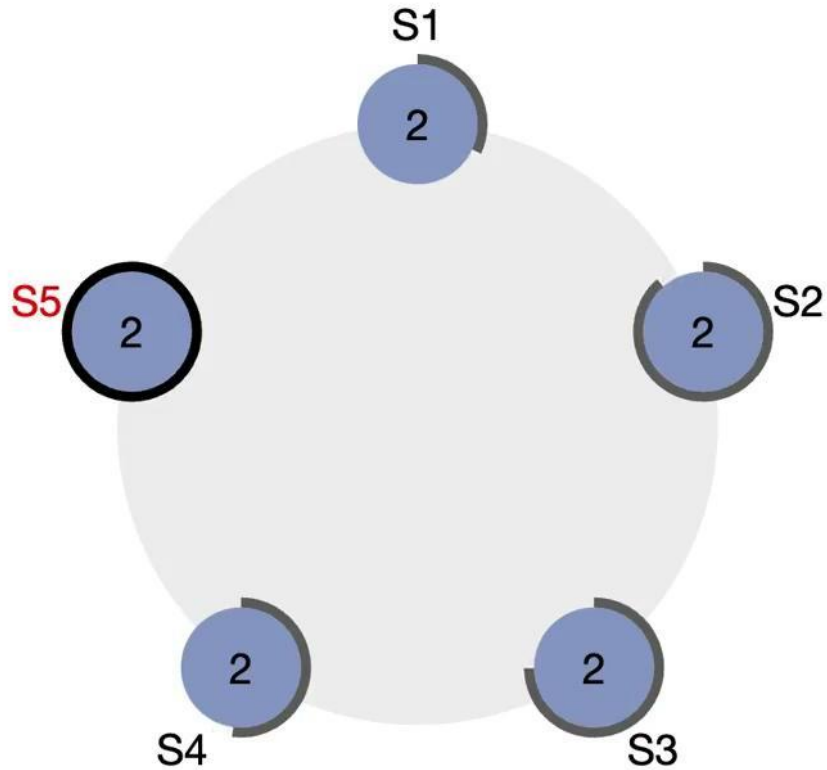
# Raft Protocol States

- Normal operation
    - System accepts log entries from clients
    - Log entries are replicated
- Leader election
    - Leader is detected to have failed
    - Remaining nodes vote on a new leader

# Raft Node States

At any time a node is either a

- Leader
    - Accepts user input and replicates to followers
- Follower
    - Listens for leader heartbeats
    - Replicates log entries
    - Detect leader failure
- Candidate
    - Claim leadership for term i+1
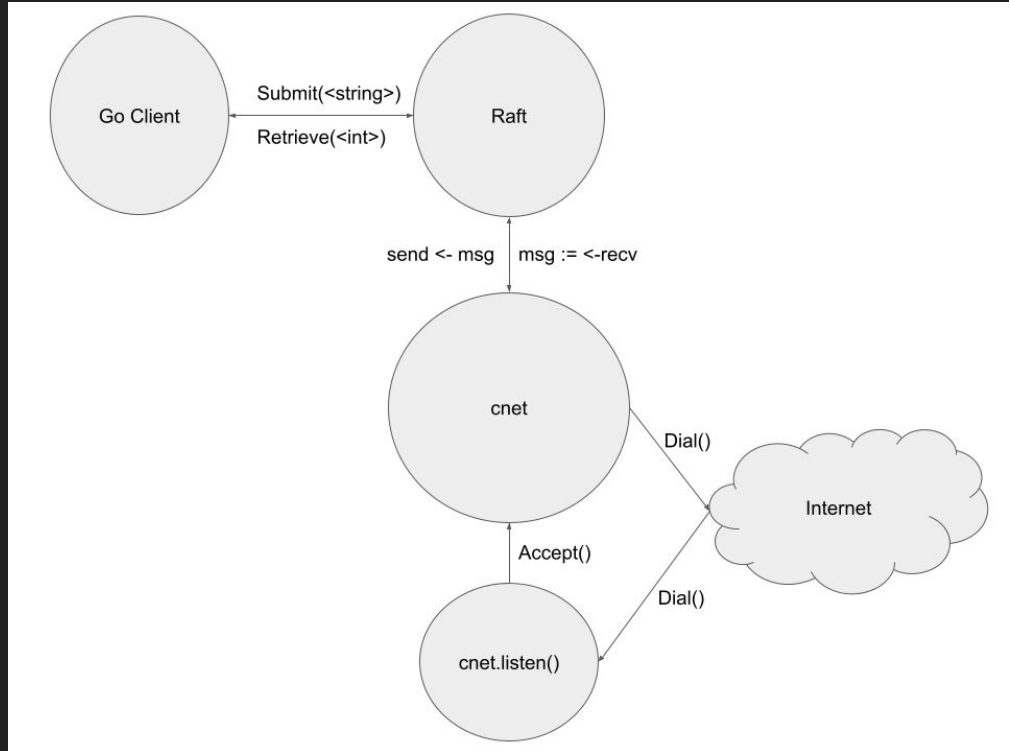    - Requests votes from other nodes

# Normal Operation

# Leader Election

# Design and Implementation

# Implementation

- Design for a distributed system
- Interesting + key code segments

# Design

- Circle == goroutine

# Client API

Why does this block?

```go
func (r *Raft) Submit(msg string) Result {
    r.input <- Entry{
        Msg: msg,
        Id:  uuid.New(),
    }

    return <-r.output
}

func (r *Raft) Retrieve(start int) []Entry {
    start = max(start, 1)
    return r.log[start : r.commitIndex+1]
}
```

# Run() goroutine

- Same select per role
- Leader & Follower iteration speed

```
91   for {
92       r.scanForAwaiting()
93
94       var timeout time.Duration
95       var callback func(chan cnet.PeerMsg)
96       switch r.role {
97       case Leader:
98           timeout = 100 * time.Millisecond
99           callback = r.sendAppendMsg
100      case Follower, Candidate:
101          timeout = time.Duration(rand.Intn(500)+500) * time.Millisecond
102          callback = r.becomeCandidate
103      }
104
105      select {
106      // receive client command
107      case entry := <-r.input:
108          r.handleInput(entry, send)
109      // handle append responses
110      case am := <-appendChan:
111          r.handleAppendMsg(am, send)
112      // handle election
113      case lm := <-leaderChan:
114          r.handleLeaderMsg(lm, send)
115      // send regular updates faster than heartbeat timeout
116      case <-time.After(timeout):
117          callback(send)
118      }
119  }
```

# Message routing

```go
129  func route(recv chan cnet.PeerMsg, appendChan chan AppendMsg, leaderChan chan LeaderMsg) {
130      for {
131          // Read messages from recv
132          pm := <-recv
133
134          // Parse payload
135          // Forward message to correct channel
136          switch pm.Type {
137          case LEADER:
138              var lm LeaderMsg
139              err := json.Unmarshal(pm.Msg, &lm)
140              if err != nil {
141                  fmt.Println(err)
142                  continue
143              }
144              leaderChan <- lm
145          case APPEND:
146              var am AppendMsg
147              err := json.Unmarshal(pm.Msg, &am)
148              if err != nil {
149                  fmt.Println(err)
150                  continue
151              }
152              appendChan <- am
153          }
154      }
155  }
```

# CNET - Simple networking

- Payload agnostic
- Runs in its own goroutine
- Failure agnostic
- Room for optimization

```
10
11    type MessageType int
12
13    type PeerMsg struct {
14        Src, Dst string
15        Msg      []byte
16        Type     MessageType
17    }
18
```

```
33
34    func (n *Network) Run(send, recv chan PeerMsg) {
35        connChan := make(chan net.Conn, 100)
36        go n.listen(connChan)
37
38        for {
39            select {
40            case conn := <-connChan:
41                pm, err := recvMsg(conn)
42                conn.Close()
43                if err != nil {
44                    fmt.Println(err)
45                    continue
46                }
47                recv <- *pm
48            case pm := <-send:
49                n.sendMsg(pm)
50            }
51        }
52    }
53
```

# Demos

# Demo 1 - ./httpraft

- Code in network send() for delaying messages
- See the protocol updating in realtime

# Demo 2 - ./fake-etcd

1.   API definition
2.   Data store definition
3.   Interaction with the raft lib

# Want to contribute?

- Log modularity
    - Rip it out and give it its own interface
    - DB, files, in-memory
- Leader-forwarding
    - Forward a request to the current leader on behalf of the client
- Network optimization
    - TCP → UDP

# Summary

- Raft is a protocol for maintaining a replicated log
- Raft is resilient to crash-faults (fail-stop)
- General approach for DS software design using Go
- Demonstrated ease-of-use with `fake-etcd`