

Consensus



COS 518: *Advanced Computer Systems*
Lecture 4

Michael Freedman

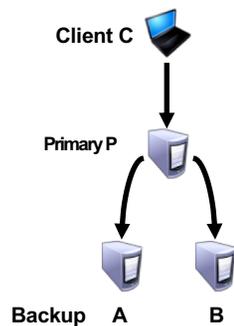
RAFT slides heavily based on those from Diego Ongaro and John Ousterhout

Recall: Linearizability (Strong Consistency)

- Provide behavior of a single copy of object:
 - Read should return the most recent write
 - Subsequent reads should return same value, until next write
- Telephone intuition:
 1. Alice updates Facebook post
 2. Alice calls Bob on phone: "Check my Facebook post!"
 3. Bob reads Alice's wall, sees her post

2

Two phase commit protocol

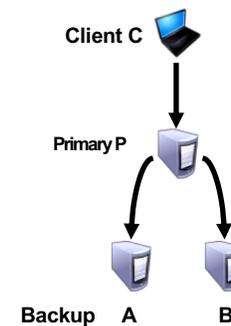


1. **C** → **P**: "request <op>"
2. **P** → **A, B**: "prepare <op>"
3. **A, B** → **P**: "prepared" or "error"
4. **P** → **C**: "result exec<op>" or "failed"
5. **P** → **A, B**: "commit <op>"

What if primary fails?
Backup fails?

3

Two phase commit protocol

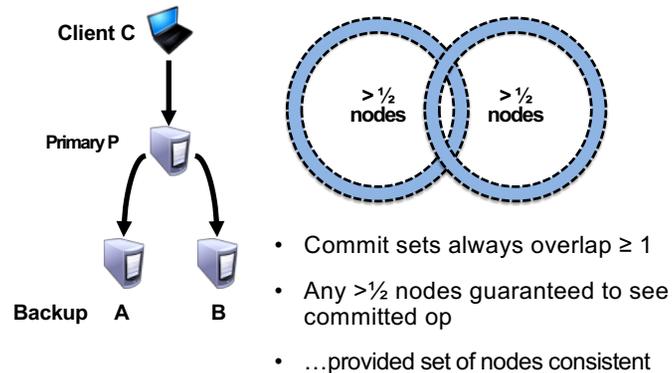


1. **C** → **P**: "request <op>"
2. **P** → **A, B**: "prepare <op>"
3. **A, B** → **P**: "prepared" or "error"
4. **P** → **C**: "result exec<op>" or "failed"
5. **P** → **A, B**: "commit <op>"

"Okay" (i.e., op is stable) if
written to > ½ nodes

4

Two phase commit protocol



5

Consensus

Definition:

1. A general agreement about something
2. An idea or opinion that is shared by all the people in a group

Origin: Latin, from *consentire*

6

Consensus used in systems

Group of servers attempting:

- Make sure all servers in group receive the same updates in the same order as each other
- Maintain own lists (views) on who is a current member of the group, and update lists when somebody leaves/fails
- Elect a leader in group, and inform everybody
- Ensure mutually exclusive (one process at a time only) access to a critical resource like a file

7

Paxos: the original consensus protocol

- Safety
 - Only a single value is chosen
 - Only a proposed value can be chosen
 - Only chosen values are learned by processes
- Liveness ***
 - Some proposed value eventually chosen if fewer than half of processes fail
 - If value is chosen, a process eventually learns it

8

Basic fault-tolerant Replicated State Machine (RSM) approach

1. Consensus protocol to elect leader
2. 2PC to replicate operations from leader
3. All replicas execute ops once committed

9

Why bother with a leader?

Not necessary, but ...

- Decomposition: normal operation vs. leader changes
- Simplifies normal operation (no conflicts)
- More efficient than leader-less approaches
- Obvious place to handle non-determinism

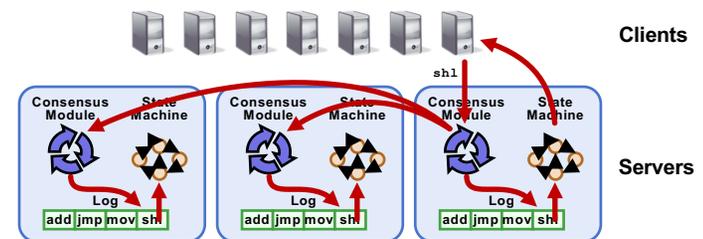
10

Raft: A Consensus Algorithm for Replicated Logs

Diego Ongaro and John Ousterhout
Stanford University

11

Goal: Replicated Log



- Replicated log => replicated state machine
 - All servers execute same commands in same order
- Consensus module ensures proper log replication

12

Raft Overview

1. Leader election
2. Normal operation (basic log replication)
3. Safety and consistency after leader changes
4. Neutralizing old leaders
5. Client interactions
6. Reconfiguration

13

Server States

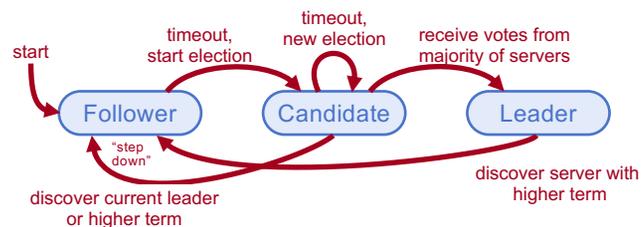
- At any given time, each server is either:
 - Leader: handles all client interactions, log replication
 - Follower: completely passive
 - Candidate: used to elect a new leader
- Normal operation: 1 leader, N-1 followers



14

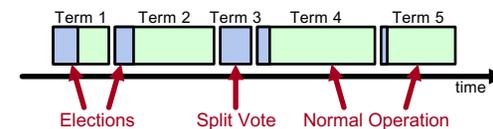
Liveness Validation

- Servers start as followers
- Leaders send **heartbeats** (empty AppendEntries RPCs) to maintain authority
- If **electionTimeout** elapses with no RPCs (100-500ms), follower assumes leader has crashed and starts new election



15

Terms (aka epochs)



- Time divided into terms
 - Election (either failed or resulted in 1 leader)
 - Normal operation under a single leader
- Each server maintains **current term** value
- **Key role of terms: identify obsolete information**

16

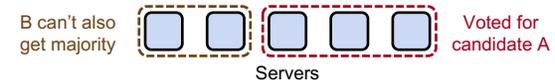
Elections

- **Start election:**
 - Increment current term, change to candidate state, vote for self
- **Send RequestVote to all other servers, retry until either:**
 1. Receive votes from **majority of servers:**
 - Become leader
 - Send AppendEntries heartbeats to all other servers
 2. Receive RPC from valid leader:
 - Return to follower state
 3. No-one wins election (election timeout elapses):
 - Increment term, start new election

17

Elections

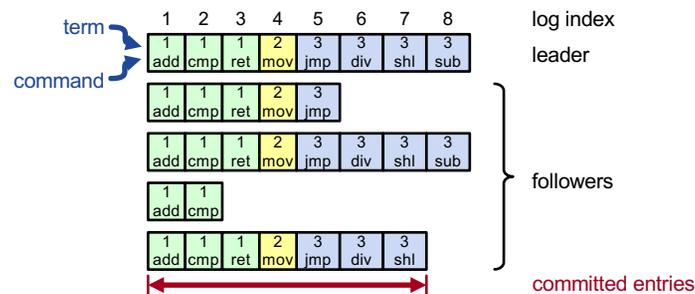
- **Safety: allow at most one winner per term**
 - Each server votes only once per term (persists on disk)
 - Two different candidates can't get majorities in same term



- **Liveness: some candidate must eventually win**
 - Each choose election timeouts randomly in $[T, 2T]$
 - One usually initiates and wins election before others start
 - Works well if $T \gg$ network RTT

18

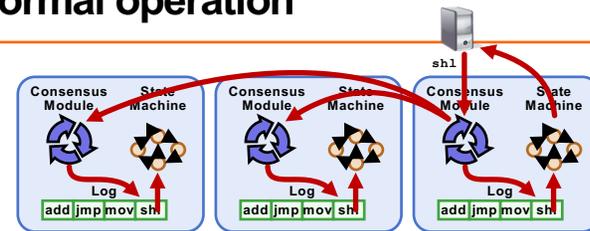
Log Structure



- Log entry = \langle index, term, command \rangle
- Log stored on stable storage (disk); survives crashes
- Entry **committed** if known to be stored on majority of servers
 - Durable / stable, will eventually be executed by state machines

19

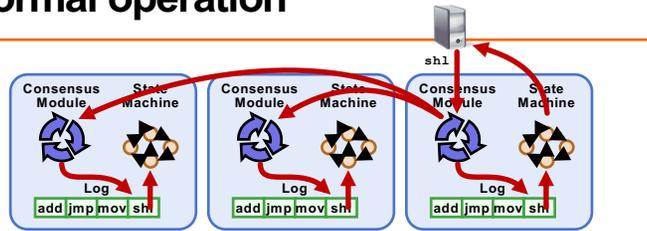
Normal operation



- Client sends command to leader
- Leader appends command to its log
- Leader sends AppendEntries RPCs to followers
- **Once new entry committed:**
 - Leader passes command to its state machine, sends result to client
 - Leader piggybacks commitment to followers in later AppendEntries
 - Followers pass committed commands to their state machines

20

Normal operation



- Crashed / slow followers?
 - Leader retries RPCs until they succeed
- Performance is optimal in common case:
 - One successful RPC to any majority of servers

21

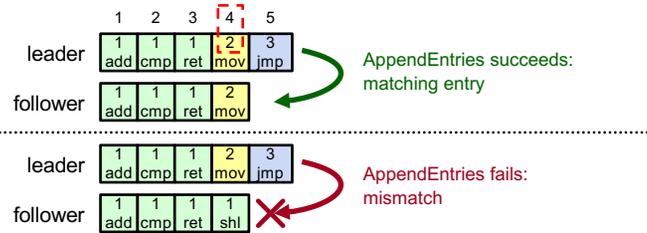
Log Operation: Highly Coherent

	1	2	3	4	5	6
server1	1 add	1 cmp	1 ret	2 mov	3 jmp	3 div
server2	1 add	1 cmp	1 ret	2 mov	3 jmp	4 sub

- If log entries on different server have same index and term:
 - Store the same command
 - Logs are identical in all preceding entries
- If given entry is committed, all preceding also committed

22

Log Operation: Consistency Check



- AppendEntries has <index,term> of entry preceding new ones
- Follower must contain matching entry; otherwise it rejects
- Implements an [induction step](#), ensures coherency

23

Leader Changes

- New leader's log is truth, no special steps, start normal operation
 - Will eventually make follower's logs identical to leader's
 - Old leader may have left entries partially replicated
- Multiple crashes can leave many extraneous log entries

	log index	1	2	3	4	5	6	7
S1	term	1	1	5	6	6	6	
S2		1	1	5	6	7	7	7
S3		1	1	5	5			
S4		1	1	2	4			
S5		1	1	2	2	3	3	3

24

Challenge: Log Inconsistencies

Leader for term 8: [1, 1, 1, 4, 4, 5, 5, 6, 6, 6]

Possible followers:

- (a) [1, 1, 1, 4, 4, 5, 5, 6, 6] ← Missing Entries
- (b) [1, 1, 1, 4]
- (c) [1, 1, 1, 4, 4, 5, 5, 6, 6, 6, 6] ← Extraneous Entries
- (d) [1, 1, 1, 4, 4, 5, 5, 6, 6, 6, 7, 7] ← Extraneous Entries
- (e) [1, 1, 1, 4, 4, 4, 4, 4]
- (f) [1, 1, 1, 2, 2, 2, 3, 3, 3, 3, 3]

Leader changes can result in log inconsistencies

25

Repairing Follower Logs

Leader for term 7: [1, 1, 1, 4, 4, 5, 5, 6, 6, 6]

Followers:

- (a) [1, 1, 1, 4]
- (b) [1, 1, 1, 2, 2, 2, 3, 3, 3, 3]

- **New leader must make follower logs consistent with its own**
 - Delete extraneous entries
 - Fill in missing entries
- **Leader keeps nextIndex for each follower:**
 - Index of next log entry to send to that follower
 - Initialized to (1 + leader's last index)
- If AppendEntries consistency check fails, decrement nextIndex, try again

Repairing Follower Logs

Leader for term 7: [1, 1, 1, 4, 4, 5, 5, 6, 6, 6]

Before repair (f): [1, 1, 1, 2, 2, 2, 3, 3, 3, 3]

After repair (f): [1, 1, 1, 4]

Safety Requirement

Once log entry applied to a state machine, no other state machine must apply a different value for that log entry

- **Raft safety property:** If leader has decided log entry is committed, entry will be present in logs of all future leaders
- Why does this guarantee higher-level goal?
 1. Leaders never overwrite entries in their logs
 2. Only entries in leader's log can be committed
 3. Entries must be committed before applying to state machine

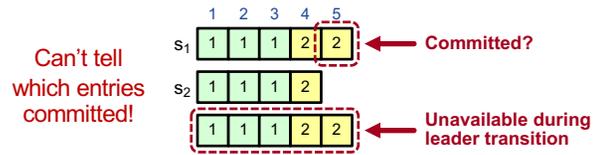
Committed → Present in future leaders' logs

Restrictions on commitment → Committed

Committed → Restrictions on leader election

28

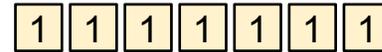
Picking the Best Leader



- Elect candidate most likely to contain all committed entries
 - In RequestVote, candidates incl. index + term of last log entry
 - Voter V denies vote if its log is “more complete”:
 - pick log whose last entry has the **higher term**
 - if last log term is the same, then pick **longer log**
 - Leader will have “most complete” log among electing majority

29

Which one is more complete?



30

Which one is more complete?



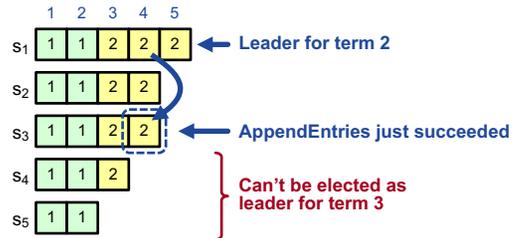
31

Which one is more complete?



32

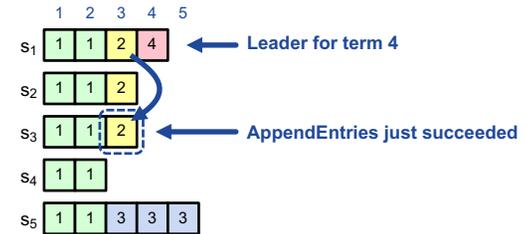
Committing Entry from Current Term



- **Case #1:** Leader decides entry in current term is committed
- **Safe:** leader for term 3 must contain entry 4

33

Committing Entry from Earlier Term



- **Case #2:** Leader trying to finish committing entry from earlier
- Entry 3 **not safely committed:**
 - S5 can be elected as leader for term 5
 - If elected, it will overwrite entry 3 on S1, S2, and S3

34

Linearizable Reads?

- **Not yet...**
 - 5 nodes: A (leader), B, C, D, E
 - A is partitioned from B, C, D, E
 - B is elected as new leader, commits a bunch of ops
 - But A still thinks he's the leader = can answer reads
 - If a client contacts A, the client will get **stale values!**
- **Fix:** Ensure you can contact majority before serving reads
 - ... by committing an extra log entry for each read
 - This guarantees you are still the rightful leader

Monday lecture

1. Consensus papers
2. From single register consistency to multi-register transactions

36

Additional Slides

37

Neutralizing Old Leaders

Leader temporarily disconnected

- other servers elect new leader
- old leader reconnected
- old leader attempts to commit log entries

• Terms used to detect stale leaders (and candidates)

- Every RPC contains term of sender
- Sender's term < receiver:
 - Receiver: Rejects RPC (via ACK which sender processes...)
- Receiver's term < sender:
 - Receiver reverts to follower, updates term, processes RPC

• Election updates terms of majority of servers

- Deposed server cannot commit new log entries

38

Client Protocol

- **Send commands to leader**
 - If leader unknown, contact any server, which redirects client to leader
- **Leader only responds after command logged, committed, and executed by leader**
- **If request times out (e.g., leader crashes):**
 - Client reissues command to new leader (after possible redirect)
- **Ensure **exactly-once semantics** even with leader failures**
 - E.g., Leader can execute command then crash before responding
 - Client should embed unique ID in each command
 - This client ID included in log entry
 - Before accepting request, leader checks log for entry with same id

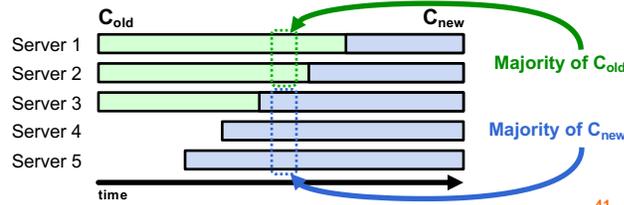
39

Reconfiguration

40

Configuration Changes

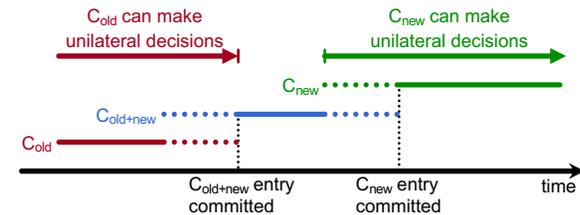
- **View configuration:** { leader, { members }, settings }
- **Consensus must support changes to configuration**
 - Replace failed machine
 - Change degree of replication
- **Cannot switch directly from one config to another: conflicting majorities could arise**



41

2-Phase Approach via Joint Consensus

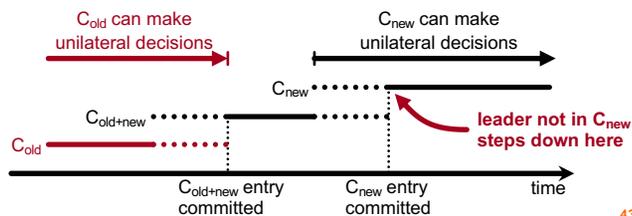
- **Joint consensus** in intermediate phase: need majority of **both** old and new configurations for elections, commitment
- Configuration change just a log entry; applied immediately on receipt (committed or not)
- Once joint consensus is committed, begin replicating log entry for final configuration



42

2-Phase Approach via Joint Consensus

- Any server from either configuration can serve as leader
- If leader not in C_{new} , must step down once C_{new} committed



43

Viewstamped Replication:

A new primary copy method to support highly-available distributed systems

Oki and Liskov, PODC 1988

44

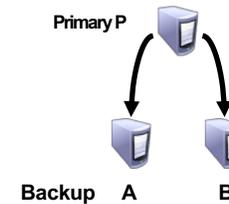
Raft vs. VR

- **Strong leader**
 - Log entries flow only from leader to other servers
 - Select leader from limited set so doesn't need to "catch up"
- **Leader election**
 - Randomized timers to initiate elections
- **Membership changes**
 - New joint consensus approach with overlapping majorities
 - Cluster can operate normally during configuration changes

45

View changes on failure

1. Backups monitor primary
2. If a backup thinks primary failed, initiate **View Change** (leader election)



46

View changes on failure

1. Backups monitor primary
2. If a backup thinks primary failed, initiate **View Change** (leader election)
3. Intuitive safety argument:
 - View change requires $f+1$ agreement
 - Op committed once written to $f+1$ nodes
 - At least one node both saw write and in new view
4. More advanced: Adding or removing nodes ("reconfiguration")

Requires $2f + 1$ nodes
to handle f failures



47