

FAULT TOLERANCE AND CONSENSUS WITH RAFT

George Porter
March 2 and 7, 2023

ATTRIBUTION

- These slides incorporate material from:
 - Diego Ongaro and John Ousterhout, Stanford University
 - Distributed Systems, 2nd Edition, Sukumar Ghosh

REQUIRED READING



In Search of an Understandable Consensus Algorithm (Extended Version) by Diego Ongaro and John Ousterhout (<https://raft.github.io/raft.pdf>).

- Section 1, 2, 4, 5, 8, and 11 are required reading
- Sections 3, 6, 7, 9, 10, and 12 are optional and not necessary for your project
- You will **not** be implementing log compaction or membership changes!

To study for this topic, please refer to the paper. Consensus protocols are very subtle and studying these slides and/or rewatching the lecture will **NOT** be sufficient for obtaining a deep understanding of the RAFT protocol.

ROADMAP ON APPROACHES TO FAULT TOLERANCE

	# servers that can fail before data is lost	Accepts updates and serves clients during failures	# servers that need to be operational to accept updates and serve clients	Transaction coordinator (or leader) can fail
Single replica	0	No	1	N/A
N replicas w/ 2-PC	N-1	No	N	No
N replicas w/ RAFT algorithm	$N/2$ (N is odd)	Yes (up to $N/2$ failures)	$\text{Ceil}(N/2)$ (N is odd)	Yes

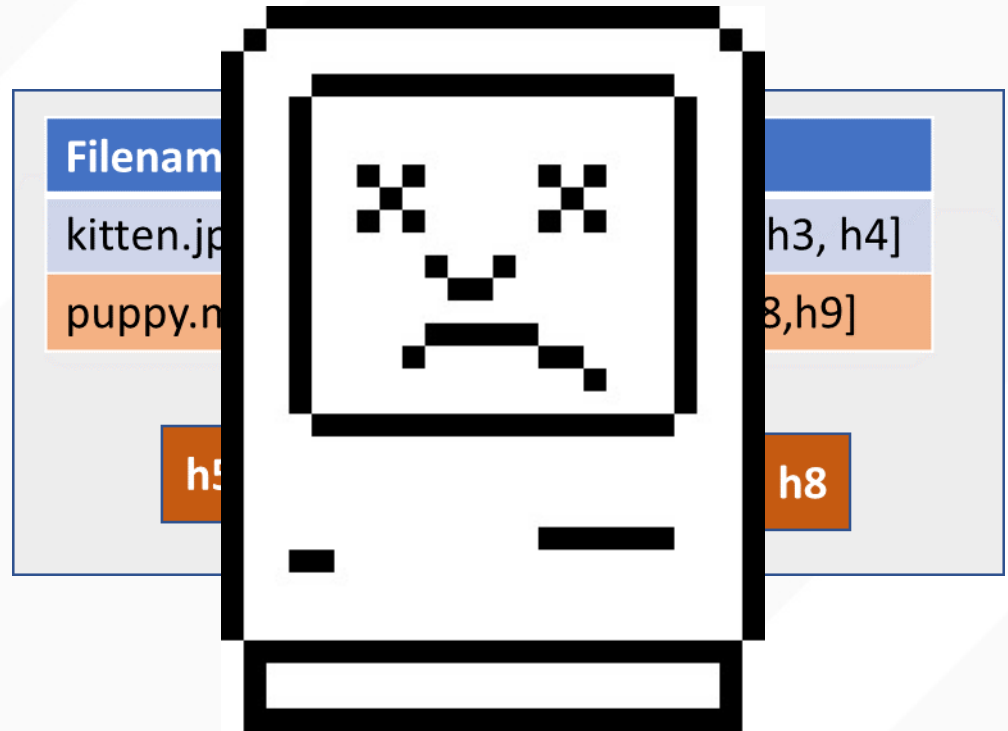
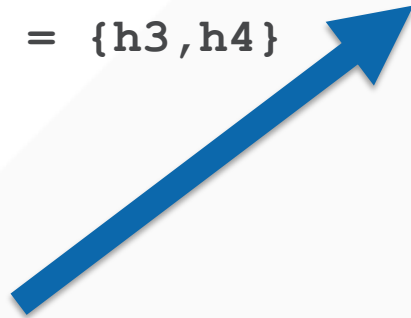
OUTLINE

1. Quorums
2. Formulate system logic as a state machine
3. Election to choose the leader
4. Leader replicates operations to multiple backup state machines
5. Mechanism to “clean up” system when leader fails



SURFSTORE METADATA SERVER PROBLEM

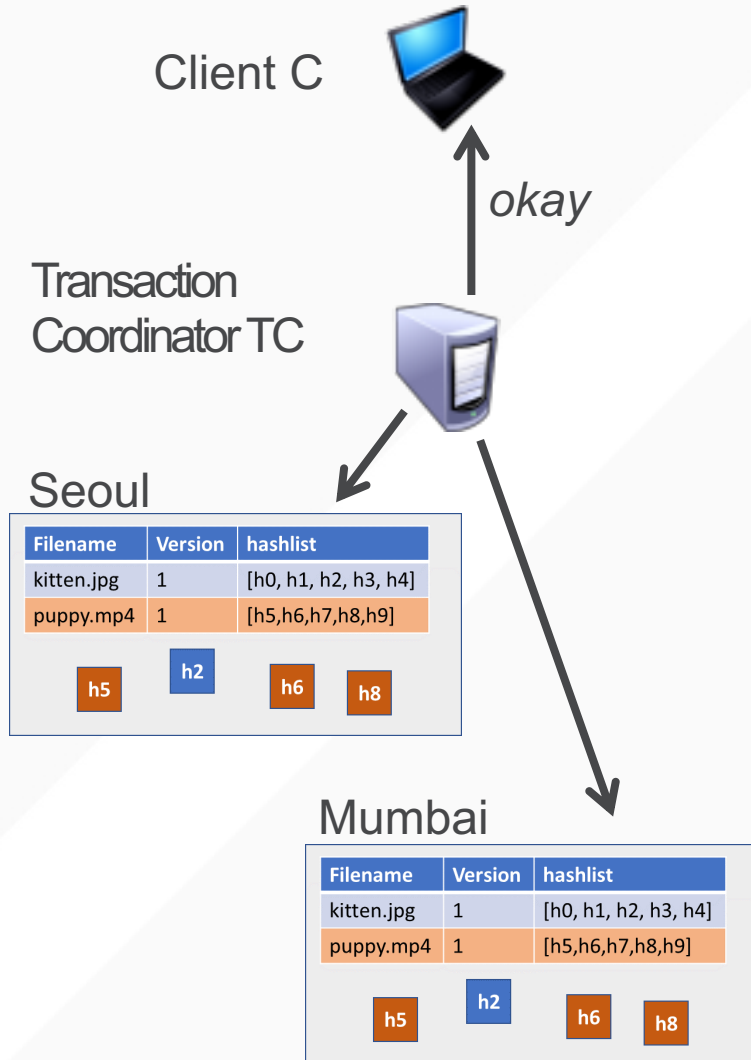
```
UpdateFile(  
  file="kitten.jpg",  
  ver=2,  
  hashlist = {h3,h4}  
);
```



All data is lost!

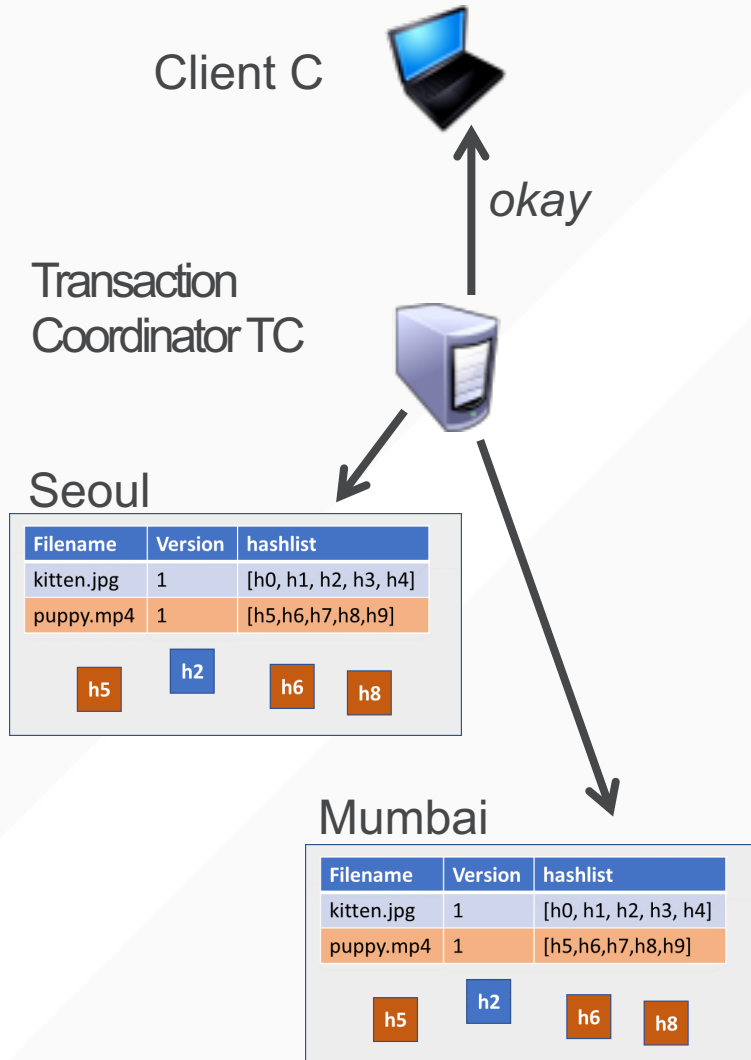
Surfstore
Client

IDEA 1: ADAPT TWO-PHASE COMMIT TO SAVE DATA



1. $C \rightarrow TC$: *“UpdateFile()”*
2. $TC \rightarrow \text{Seoul (S), Mumbai (M)}$: *“prepare!”*
3. $S, M \rightarrow P$: *“yes”* or *“wrong_version”*
4. $TC \rightarrow S, M$: *“commit!”* or *“abort!”*
 - TC sends **commit** if **both** say **yes**
 - TC sends **abort** if **either** say **no**
5. $TC \rightarrow C$: *“okay”* or *“failed”*
 - **S, M** commit on receipt of commit message

IDEA 2: ASSUME TC DOESN'T FAIL (FOR NOW)



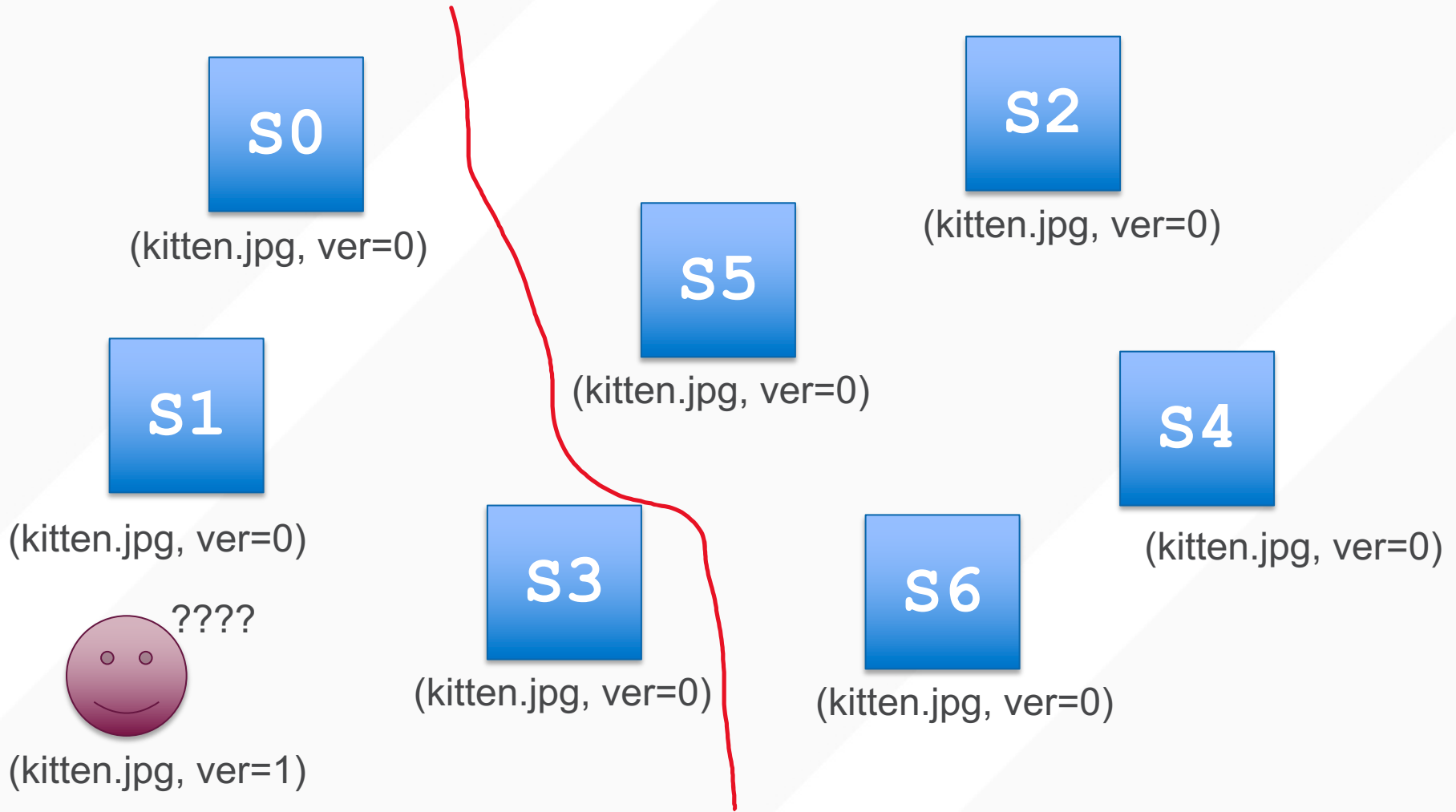
1. $C \rightarrow TC$: *"UpdateFile()"*
2. $TC \rightarrow \text{Seoul (S), Mumbai (M)}$: *"prepare!"*
3. $S, M \rightarrow P$: *"yes" [why always yes?]*
4. $TC \rightarrow S, M$: *"commit!"*
 - TC sends **commit**
5. $TC \rightarrow C$: *"okay"*
 - S, M commit on receipt of commit message
 - *Why do we still need the commit?*

NETWORK PARTITIONS

- Some failure (either network or host) keeps replicas from communicating with one another
- Two-phase commit (even if we assume all replicas agree) only works if all nodes can be contacted
- How to proceed with read/write transactions in case where not all replicas can be contacted?

DATA REPLICATION

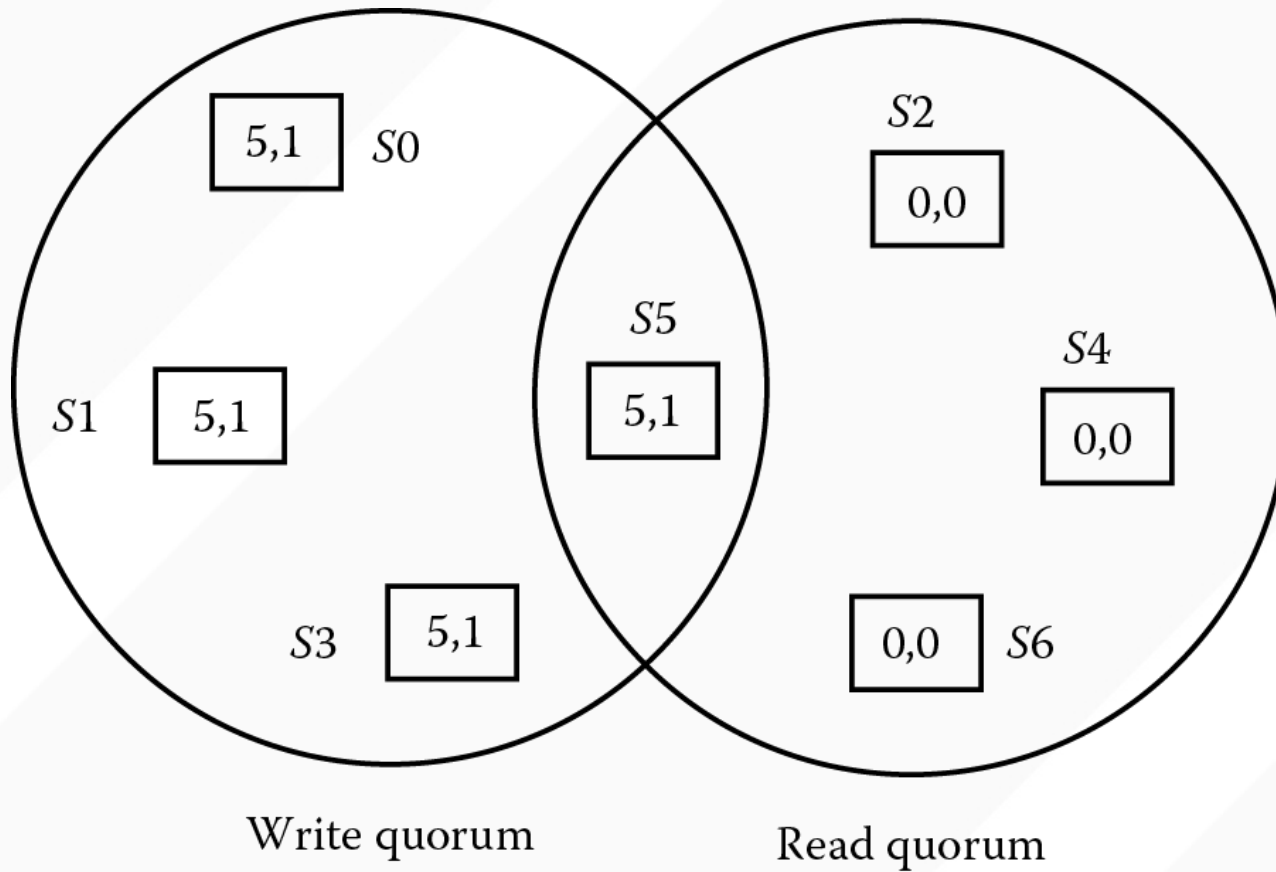
- Idea: Replicate data across multiple servers



QUORUM-BASED PROTOCOLS

- Idea: Tell client that a file's version is updated after a subset of SurfStoreServers get the update
- Form a “read quorum” of size N_R
 - Contact N_R servers and read all their versions
 - Select highest version as the “correct” version
- Form a “write quorum” of size N_W
 - Contact N_W servers
 - Increment the highest version from that set
 - Write out that new version to the servers in the write quorum

READ AND WRITE QUORUMS

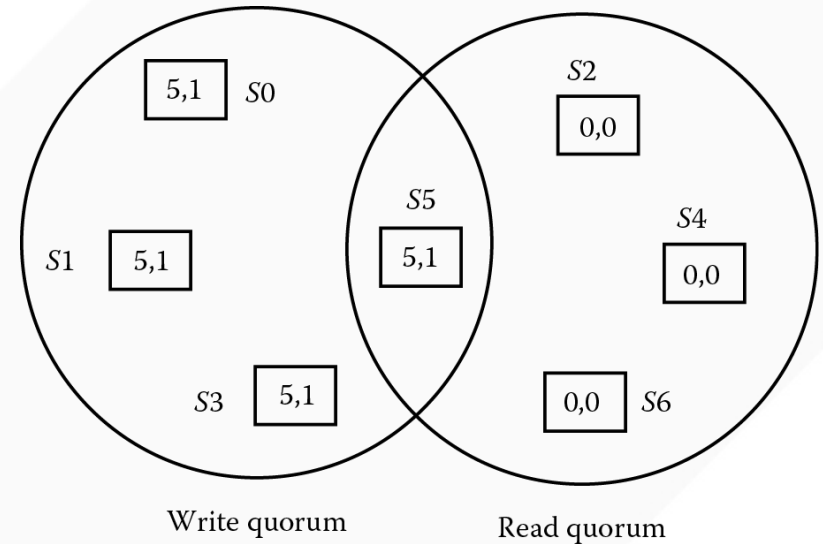


CONSTANTS AND CONSTRAINTS

- N : Total #Replicas
- N_R : #Replicas in Read Quorum
- N_W : #Replicas in Write Quorum
- Constraints:
 1. $N_R + N_W > N$
 2. $N_W > N/2$

READING AND WRITING TO QUORUMS


- To read:
 - Get “read locks” on N_r nodes
- To write:
 - Get “write locks” on a N_w nodes



QUORUM CONSENSUS

- Write operations can be propagated in background to replicas not in quorum
 - Assumes eventual repair of any network partition
- Operations are slowed by the necessity of first gathering a quorum
 - Though previously, all writes had to go to all replicas
 - With quorum system, must only contact subset of replicas

QUORUMS IN MICROSOFT ACTIVE DIRECTORY

 **Microsoft**

Windows IT Pro Center

Explore ▾ Docs ▾ Downloads ▾ Scripts Support

Docs / Windows Server / Failover Clustering / Deploy / Manage quorum and witnesses

Feedback Edit Share Da

Filter by title

Failover Clustering

What's New in Failover Clustering

> Understand

> Plan

▾ Deploy

Create a failover cluster

Deploy a two-node file server

> Prestage a cluster in AD DS

Manage quorum and witnesses


Deploy a Cloud Witness

Deploy a file share witness

Cluster operating system rolling upgrades

> Manage

Configure and manage quorum

01/17/2019 • 20 minutes to read • Contributors 

Applies to: Windows Server 2019, Windows Server 2016, Windows Server 2012 R2, Windows Server 2012

This topic provides background and steps to configure and manage the quorum in a Windows Server failover cluster.

Understanding quorum

The quorum for a cluster is determined by the number of voting elements that must be part of active cluster membership for that cluster to start properly or continue running. For a more detailed explanation, see the [understanding cluster and pool quorum doc](#).

Quorum configuration options

The quorum model in Windows Server is flexible. If you need to modify the quorum

In this

Under quoru

Quoru config option

Gener recom for qu config

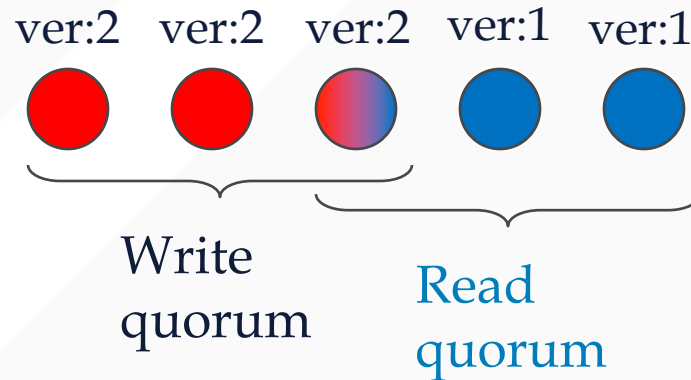
Config cluster

Recover startin quoru

Quoru consid disaster config

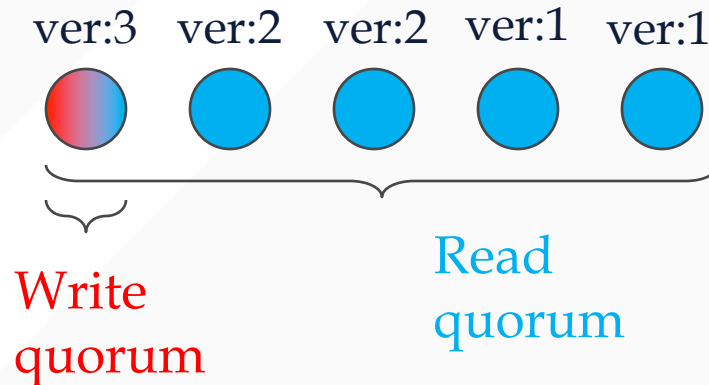
More

QUORUM EXAMPLE



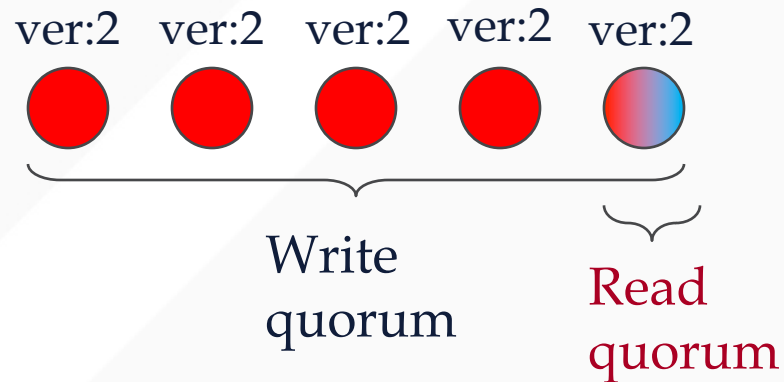
- 5 replicas, read quorum: 3, write quorum: 3
- $R+W > 5$ votes ensures overlap between any read/write quorum
- How does this perform for reads?
- How does this perform for writes?

QUORUM EXAMPLE



- 5 replicas, read quorum: 5, write quorum: 1
 - $R+W > 5$ votes ensures overlap between any read/write quorum
- How does this perform for reads?
- How does this perform for writes?

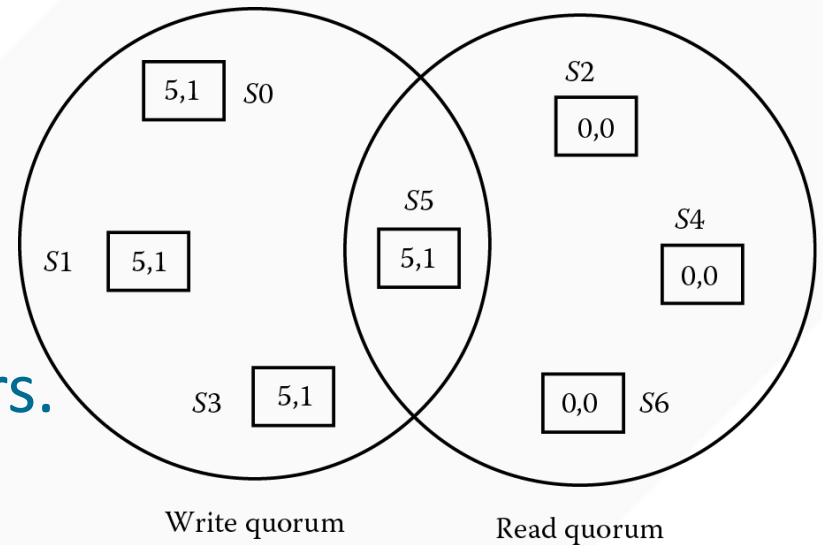
QUORUM EXAMPLE



- 5 replicas, read quorum: 1, write quorum: 5
 - $R+W > 5$ votes ensures overlap between any read/write quorum
 - Also called ROWA (read one, write all)
- How does this perform for reads?
- How does this perform for writes?

OBSERVATIONS

- Observation 1
- The system is resilient to the crash of $f \leq (N/2) - 1$ servers.



- Observation 2
- Since a read lock does not block readers, multiple readers can concurrently read: For example, one reader can read from the quorum {S0, S1, S2, S4}, while a second reader can read from the quorum {S1, S3, S5, S6}.

OBSERVATIONS

- Observation 3
 - Two different write operations cannot proceed at the same time, so all write operations are serialized. Furthermore, the intersection of the read quorum and the write quorum is nonempty, so reads do not overlap with writes. As a result, every read operation returns the latest version that was written

OUTLINE

1. Quorums
2. Formulate system logic as a state machine
3. Election to choose the leader
4. Leader replicates operations to multiple backup state machines
5. Mechanism to “clean up” system when leader fails



FORMULATE SYSTEM AS A STATE MACHINE

GENERAL CATALOG 2018–19 FEBRUARY 6, 2019 INTERIM UPDATE

UC San Diego

[Home](#) [Courses/Curricula/Faculty](#)

4. **General Science:** Two courses chosen from PHYS 2A, PHYS 2B, PHYS 4A, PHYS 4B, CHEM 6A or CHEM 6AH, CHEM 6B or CHEM 6BH, BILD 1, BILD 2, BILD 3, and BICD 100 (8 units)
5. **Probability and Statistics:** MATH 183 or ECE 109 or ECON 120A or CSE 103 (4 units)

2. Upper-Division Requirements

Students must complete 72 upper-division units: 44 units of core courses and 28 units of cluster and elective courses.

1. Core Courses:

- Data structures and programming: CSE 101
- Algorithms/theory: CSE 101 and CSE 105
- Software engineering: CSE 110
- Hardware/architecture: CSE 140 and CSE 141, along with CSE 140L and CSE141L
- Systems/networks: CSE 120 or CSE 123 or CSE 124
- PL/databases: CSE 130 or CSE 132A
- Security/cryptography: CSE 107 or CSE 127
- Learning/vision/graphics: CSE 150 or CSE 151 or CSE 152 or CSE 153 or CSE 158 or CSE 167

Students are expected to complete the majority of these courses by the end of their junior year.

REPLICATED DETERMINISTIC FINITE STATE MACHINES

PLAYER 1

1. GO NORTH
2. GO NORTH
3. GO EAST
4. GET SWORD
5. OPEN DOOR
6. GO SOUTH
7. FIGHT DRAGON
8. GET LAMP

PLAYER 2

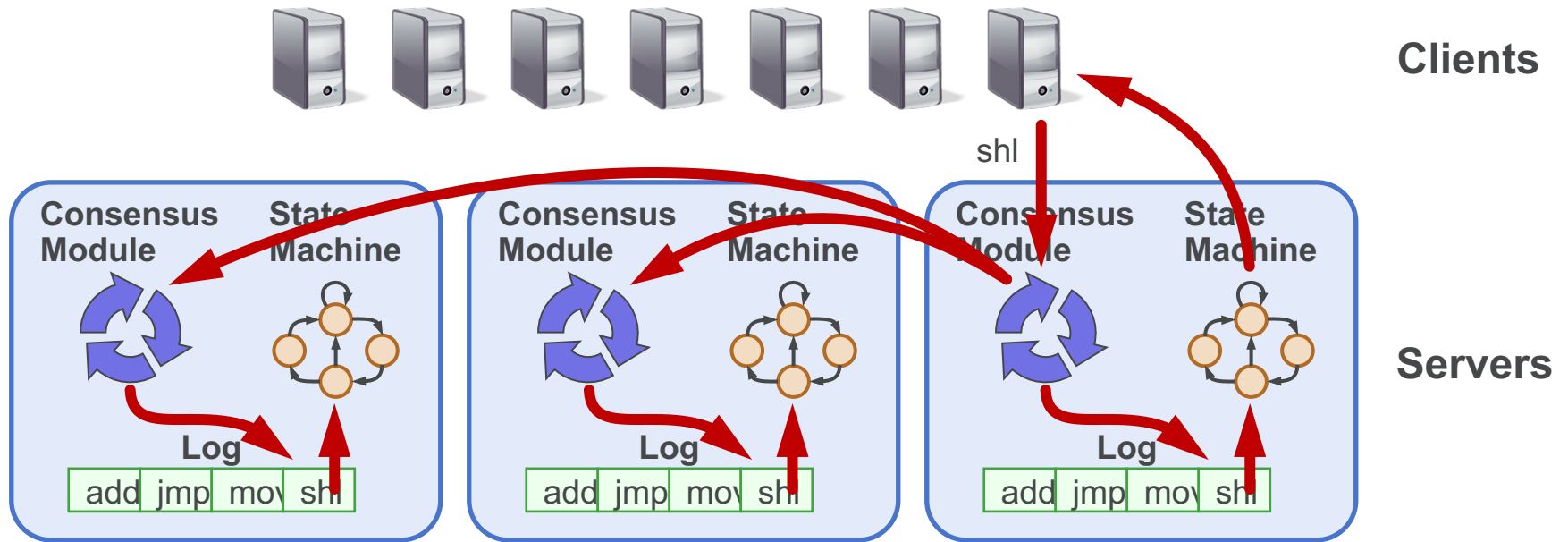
1. GO NORTH
2. GO NORTH
3. GO EAST
4. GET SWORD
5. OPEN DOOR
6. GO SOUTH
7. FIGHT DRAGON
8. GET LAMP

Both players in same state after same invocation of commands

STATE MACHINE REPLICATION

- Each machine starts in the same initial state
- Executes the same requests (deterministic)
- Requires consensus to execute in same order
 - (GET SWORD; FIGHT DRAGON) has a very different outcome from (FIGHT DRAGON; GET SWORD)
- Produces the same output

Goal: Replicated Log



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

OUTLINE

1. Quorums
2. Formulate system logic as a state machine
3. Election to choose the leader
4. Leader replicates operations to multiple backup state machines
5. Mechanism to “clean up” system when leader fails



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 such as raw quorum replication**
- **Obvious place to handle non-determinism (leader chooses the sequence)**

Image courtesy of Reuters

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

Follower

Candidate

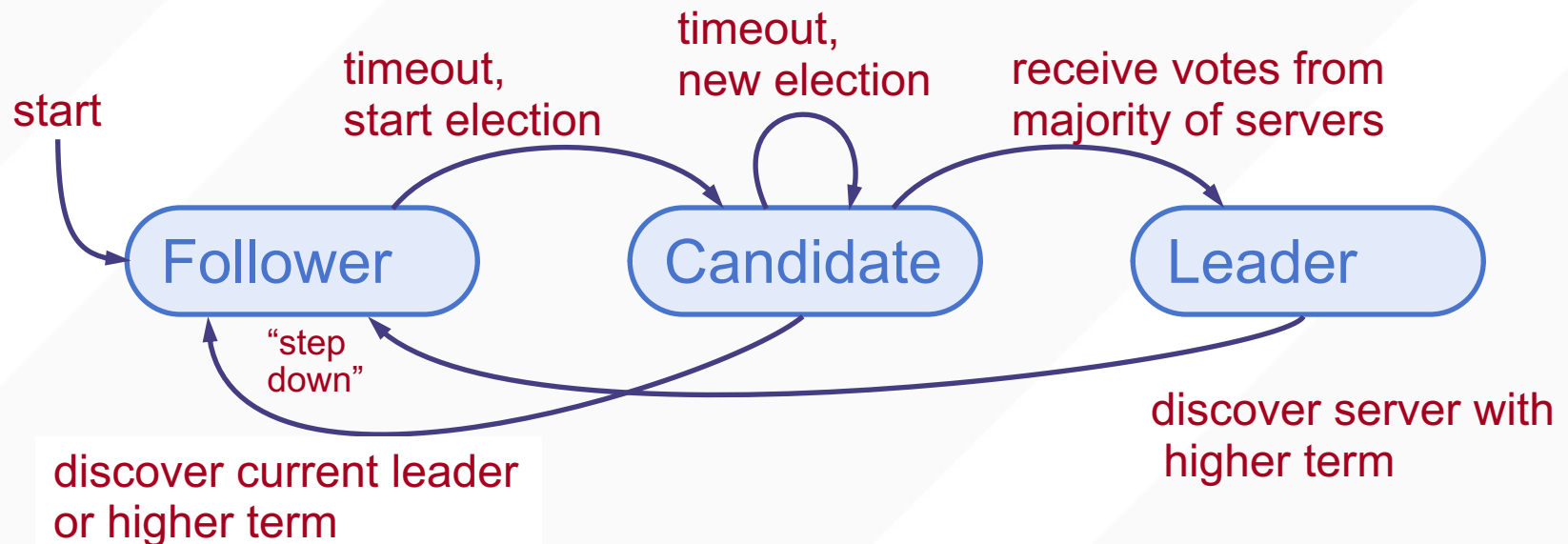
Leader

OPERATIONS

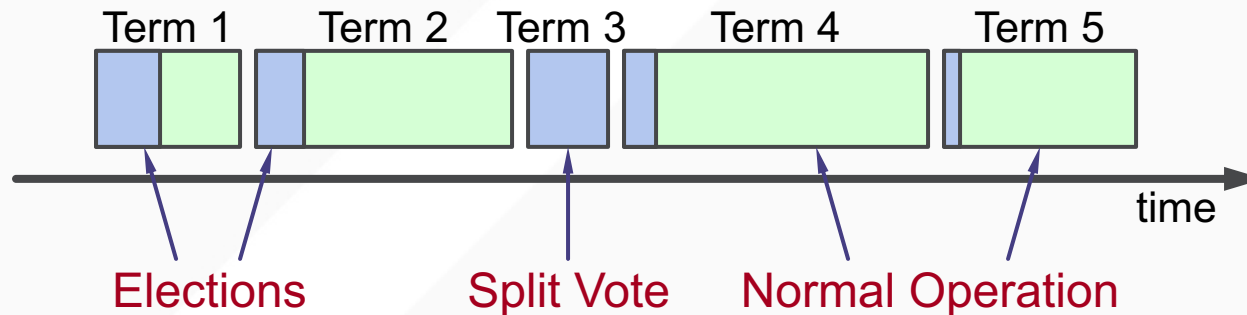
- AppendEntries()
 - The leader uses this to “push” new operations to the replicated state machines
 - Also used by the leader to tell the other nodes it is the leader
- RequestVote()
 - Used when the system starts up to select a leader
 - Used when the leader fails to elect a new leader
 - Used when the leader is unreachable due to a network partition to elect a new leader

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



TERMS (AKA EPOCHS)



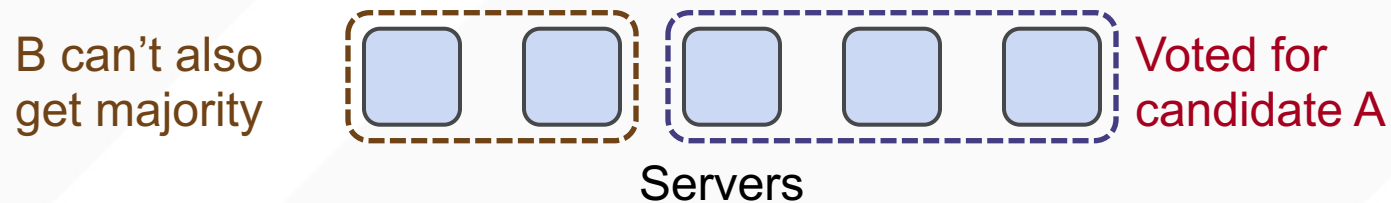
- Time divided into terms
 - Election (either failed or resulted in 1 leader)
 - Normal operation under a single leader
- Each server maintains current term value
- Key role of terms: identify obsolete information

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

ELECTION PROPERTIES

- **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 $T \gg \text{network RTT}$

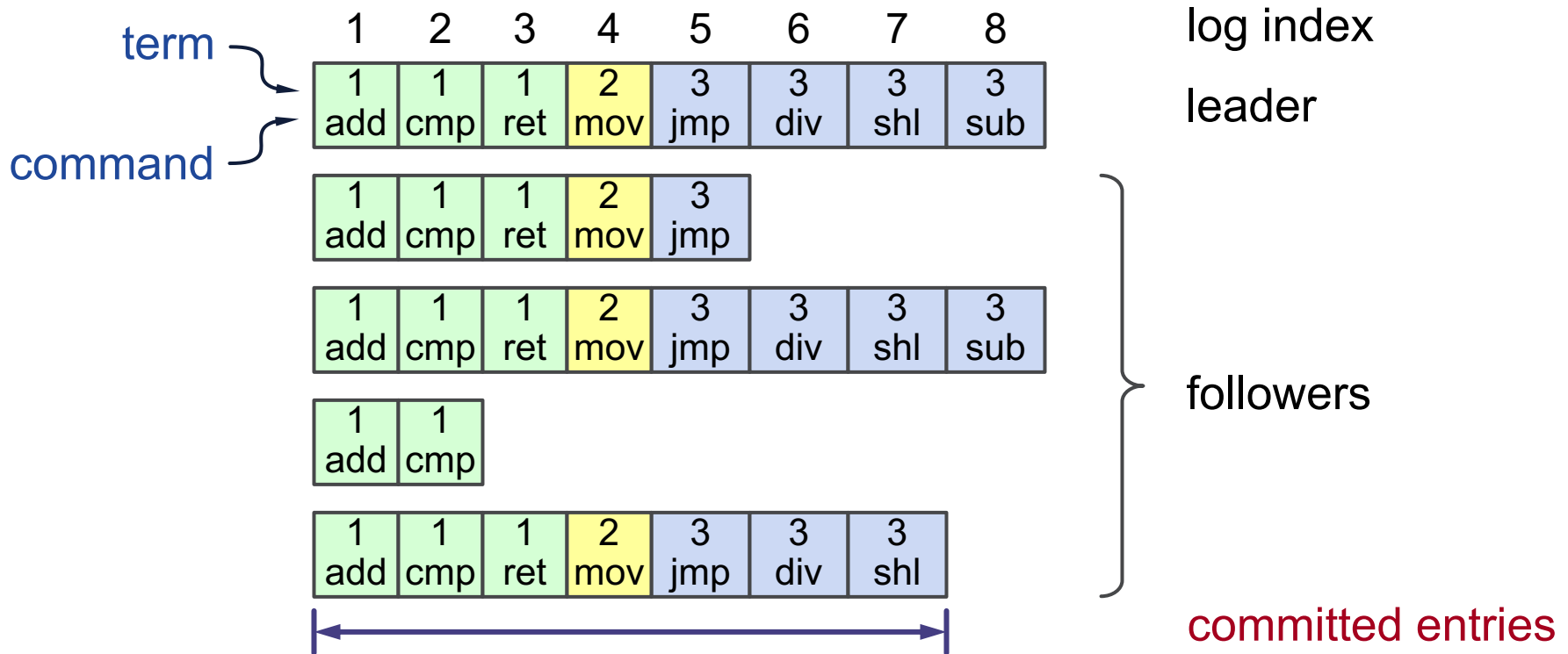
WEB SIMULATOR/DEMO

<https://raft.github.io/raftscope/index.html>

RAFT OVERVIEW

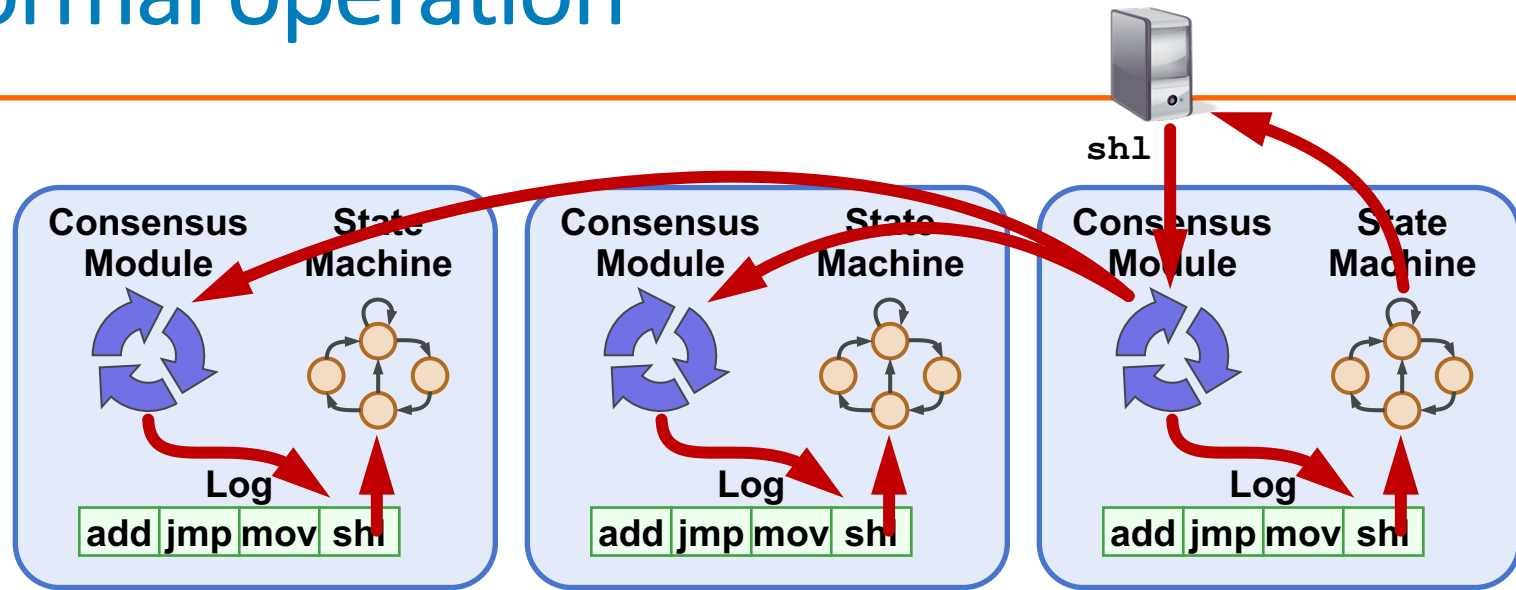
- ~~1. Leader election~~
2. Normal operation (basic log replication)
3. Safety and consistency after leader changes
4. Neutralizing old leaders
5. Client interactions
- ~~6. Reconfiguration~~

Log Structure



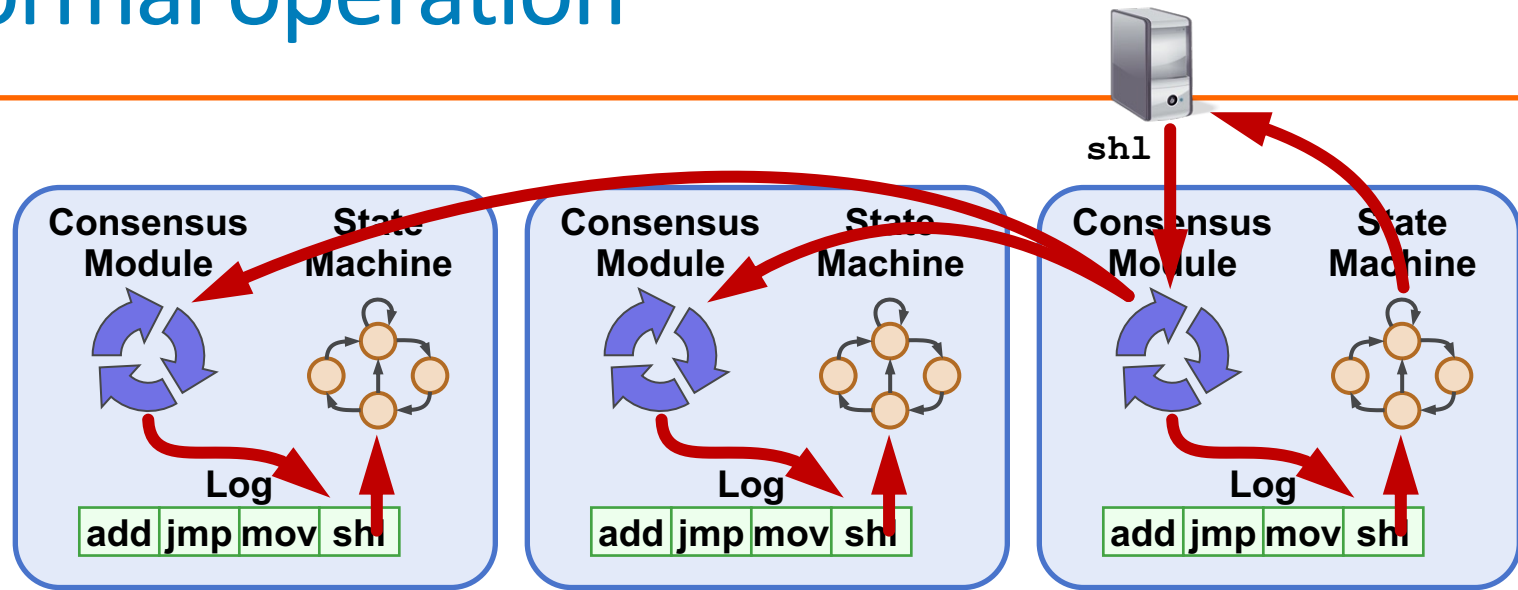
- Log entry = < index, term, command >
- Log stored on stable storage (disk); survives crashes
- Entry **committed** if known to be stored on majority of servers
 - Durable / stable, will eventually be executed by state machines

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

Normal operation



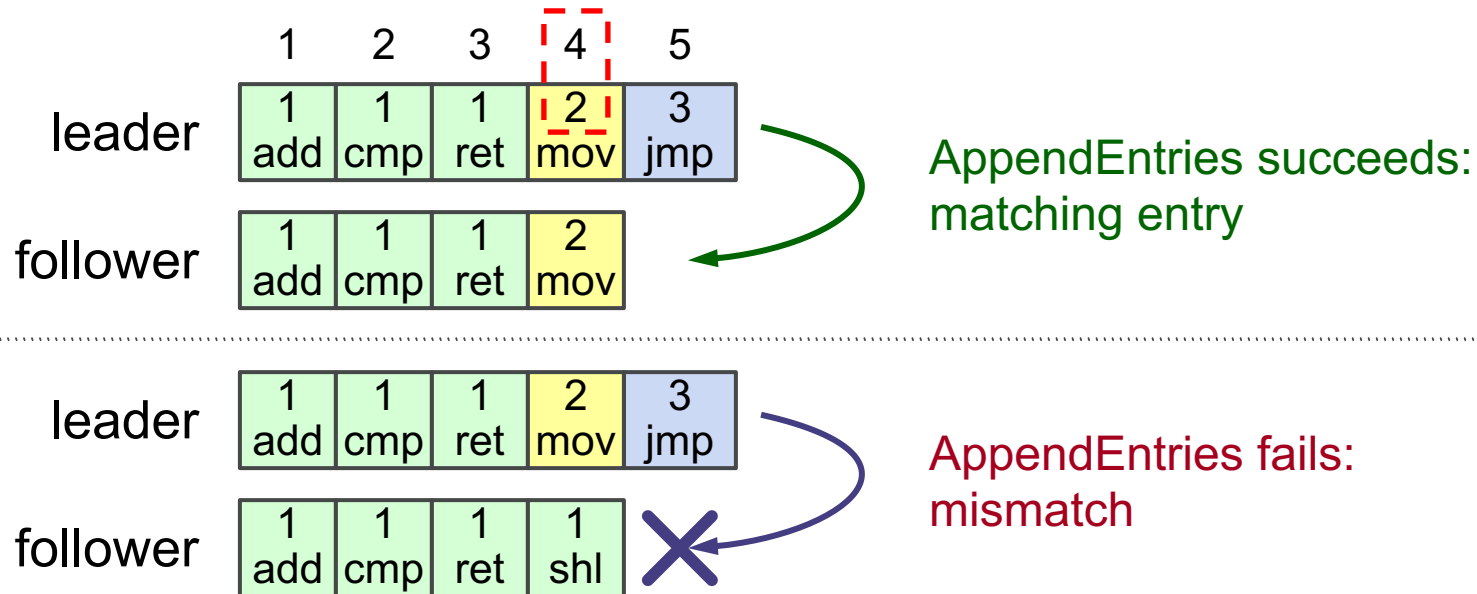
- Crashed / slow followers?
 - Leader retries RPCs until they succeed
- Performance is optimal in common case:
 - One successful RPC to any majority of servers

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

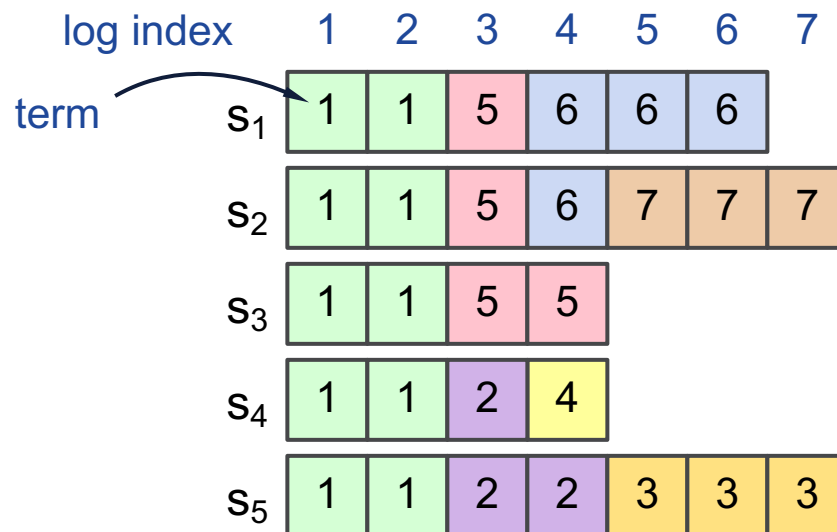
Log Operation: Consistency Check



- AppendEntries has $\langle \text{index}, \text{term} \rangle$ of entry preceding new ones
- Follower must contain matching entry; otherwise it rejects
- Implements an induction step, ensures coherency

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



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

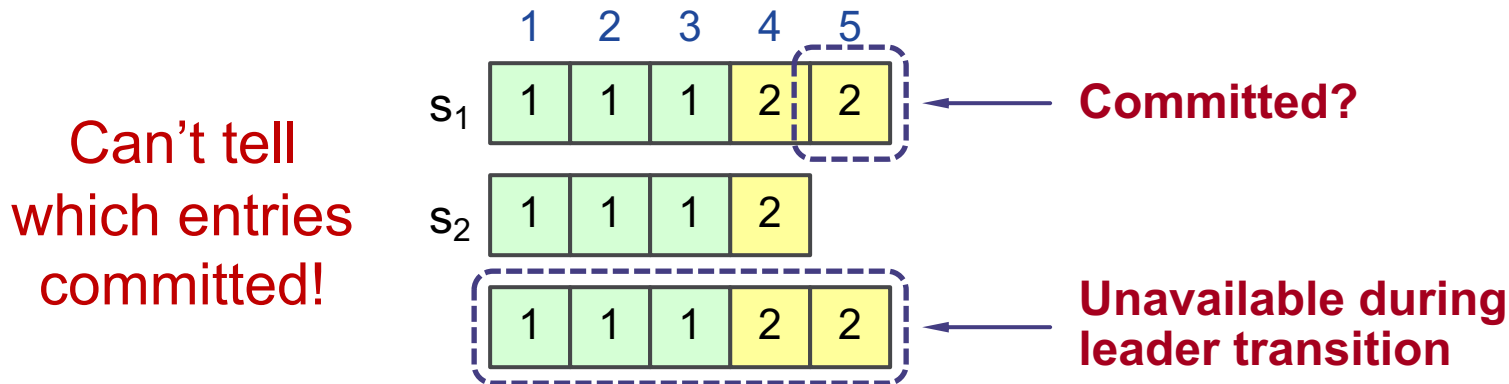
Restrictions on
commitment



Restrictions on
leader election

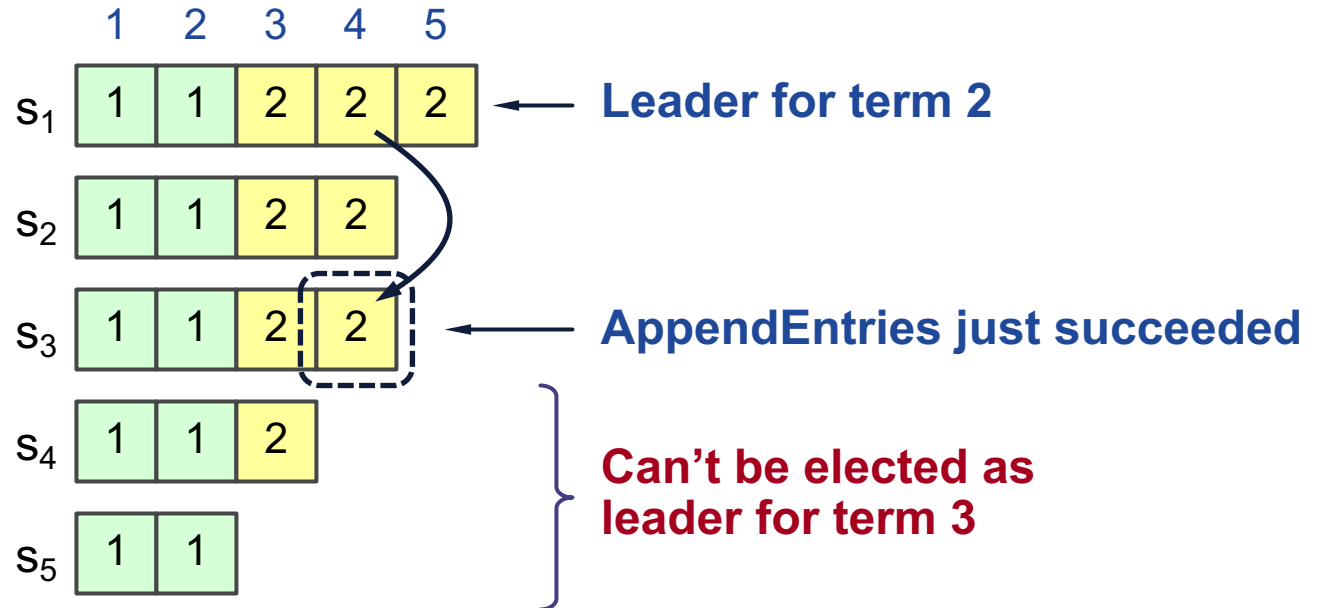


Picking the Best Leader



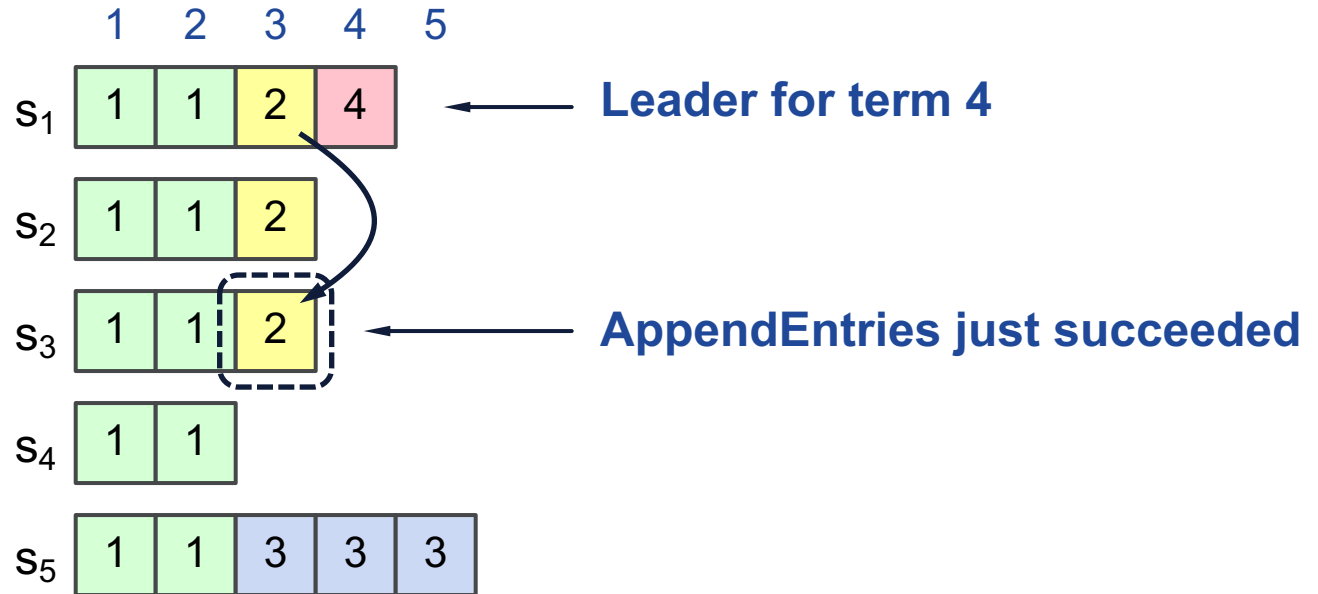
- Elect candidate most likely to contain all committed entries
 - In RequestVote, candidates incl. index + term of last log entry
 - Voter V denies vote if its log is “more complete”: (newer term) or (entry in higher index of same term)
 - Leader will have “most complete” log among electing majority

Committing Entry from Current Term



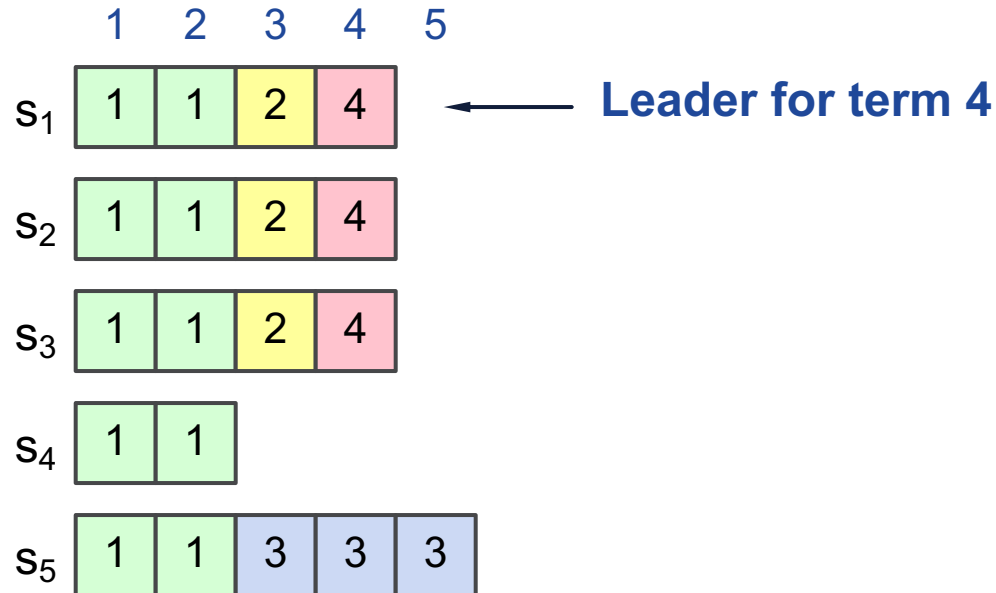
- Case #1: Leader decides entry in current term is committed
- **Safe:** leader for term 3 must contain entry 4

Committing Entry from Earlier Term



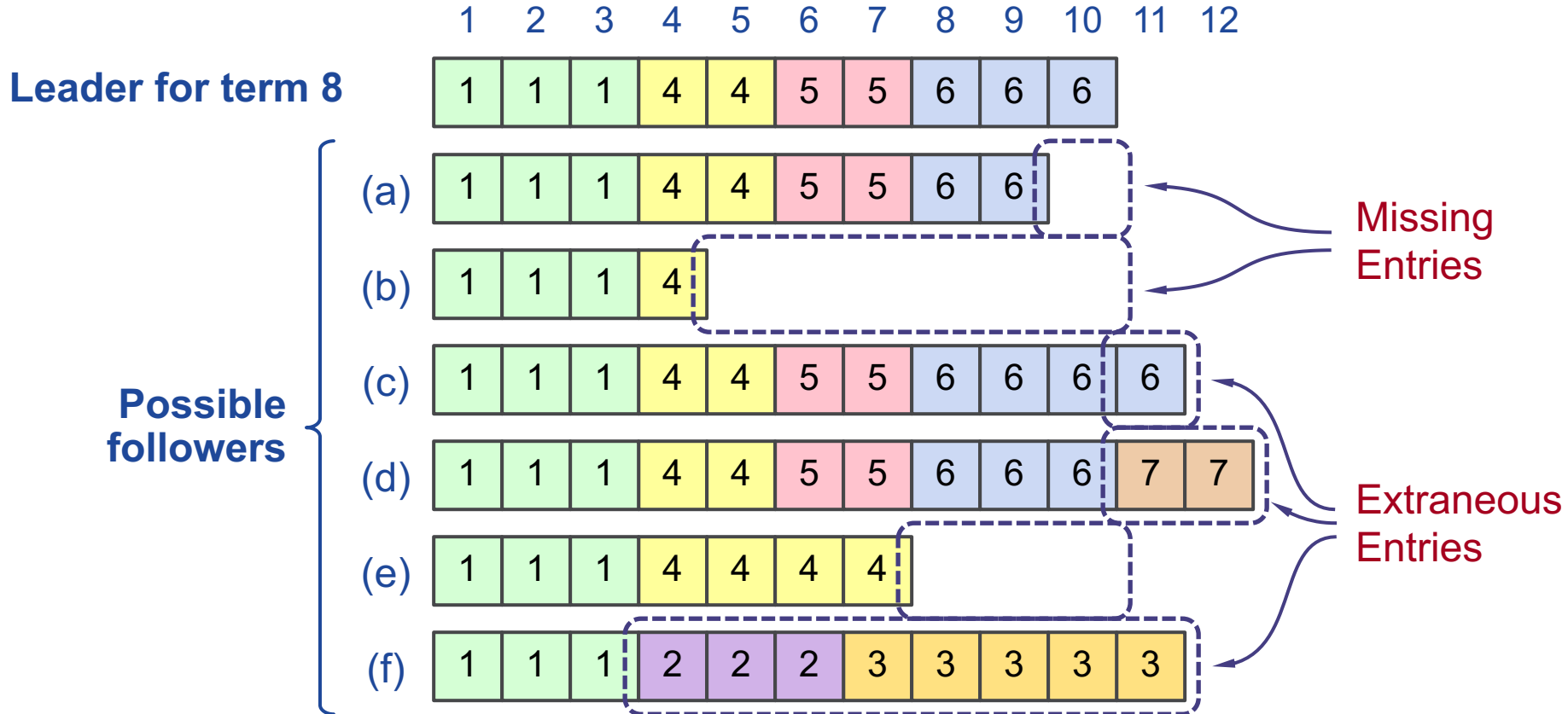
- Case #2: Leader trying to finish committing entry from earlier
- Entry 3 not safely committed:
 - s_5 can be elected as leader for term 5 (how?)
 - If elected, it will overwrite entry 3 on s_1 , s_2 , and s_3

New Commitment Rules



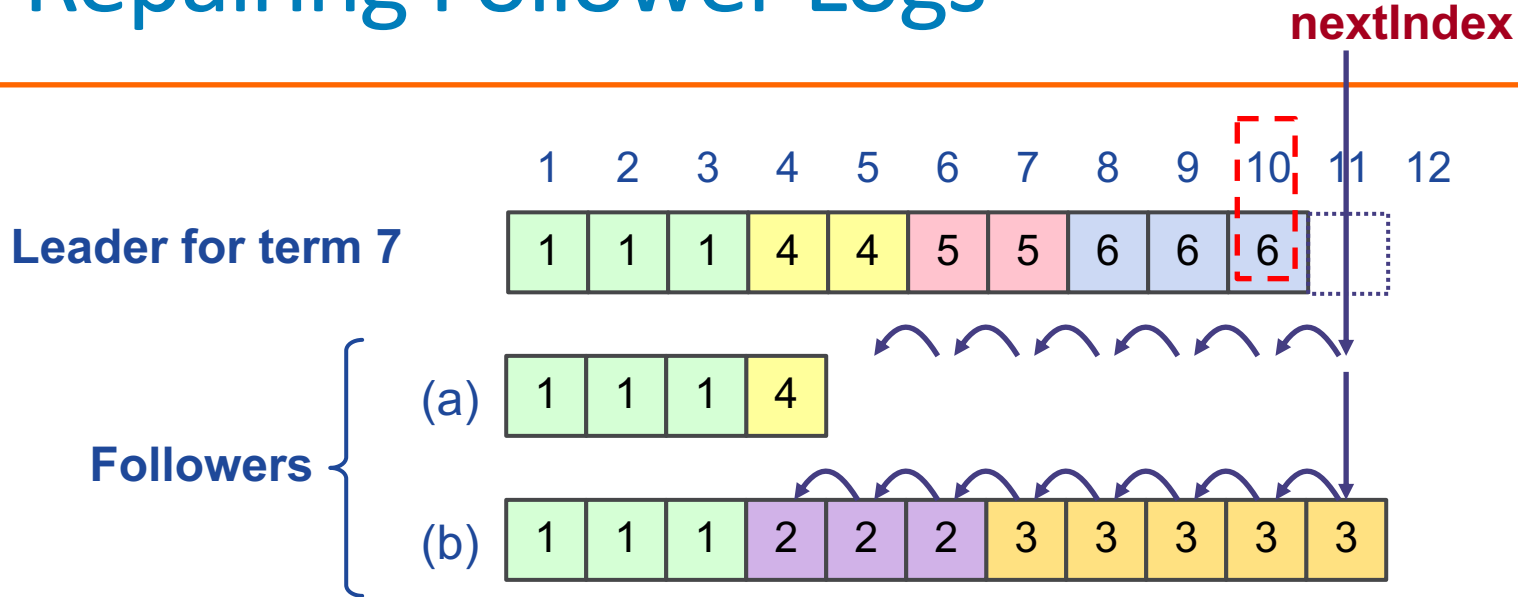
- For leader to decide entry is committed:
 1. Entry stored on a majority
 2. ≥ 1 new entry from leader's term also on majority

Challenge: Log Inconsistencies



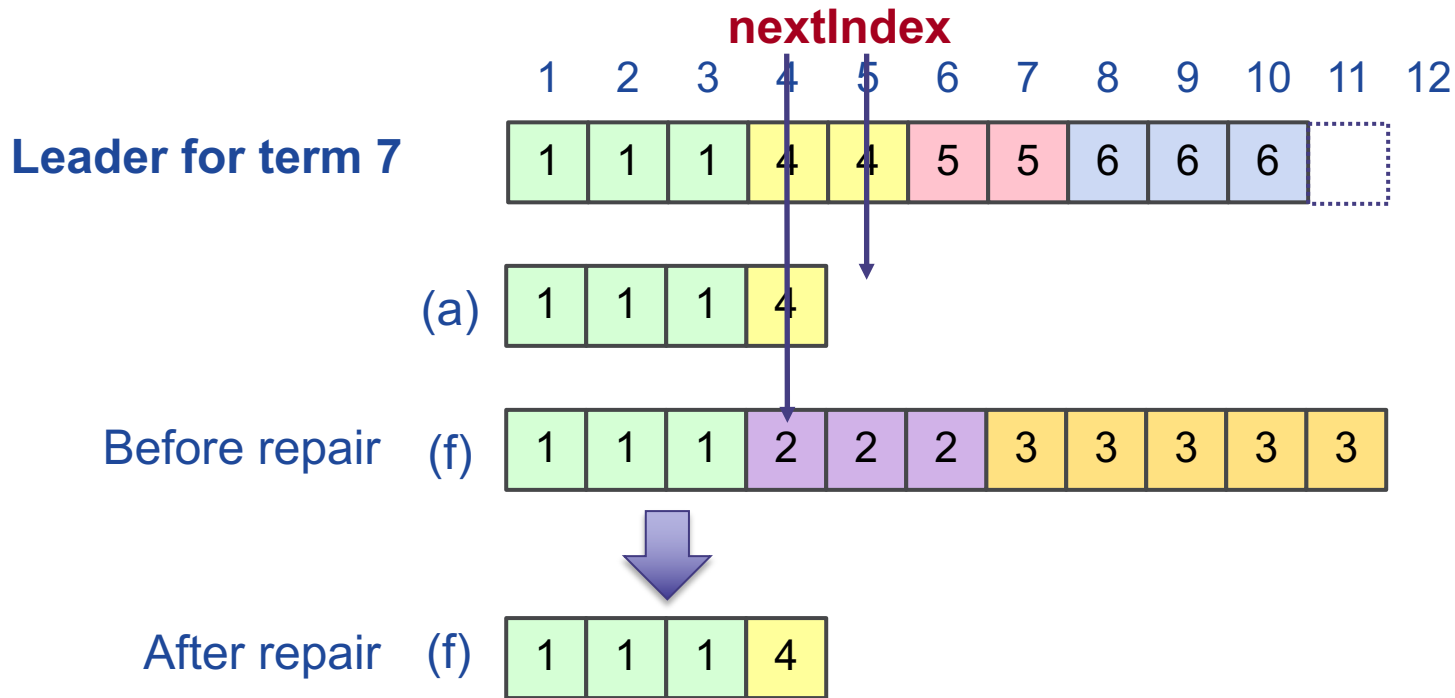
Leader changes can result in log inconsistencies

Repairing Follower Logs



- **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



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

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

WEB SIMULATOR/DEMO

<https://raft.github.io/raftscope/index.html>

UC San Diego