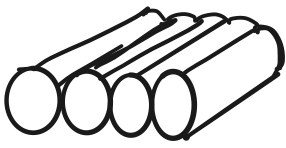


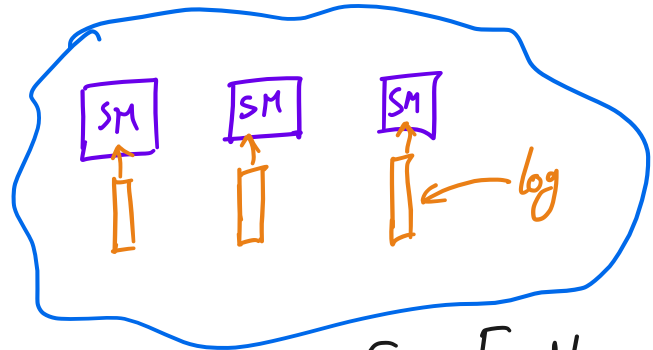
# Raft



Last week: RSMs: A framework for building fault tolerant distributed application

State Machine

Deterministic



Replicated For Fault Tolerance

Raft is one such protocol.

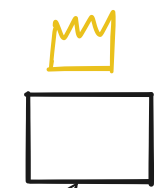
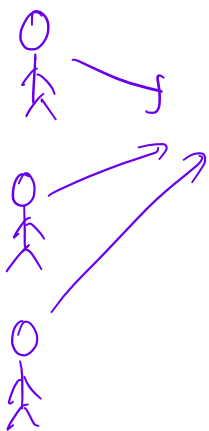
Assumptions:

Partial synchrony

Safety Guarantees:

Agreement

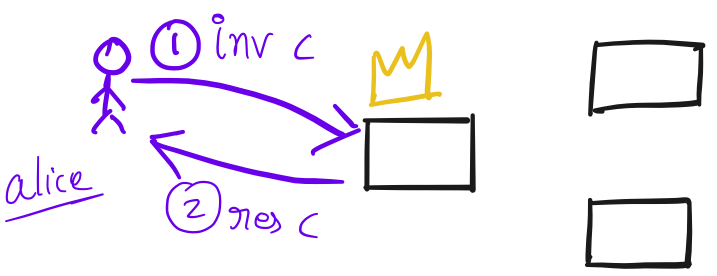
Validity



- Sequencing Commands
- Broadcasting to Replicate



Useful to Understand Requirements



① alice invokes command  $c$

② System responds to alice

Commit Point

$c$ 's effects are guaranteed to persist despite permissible failures

Commit Point in Raft

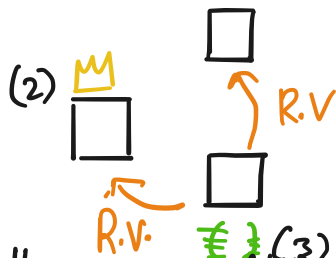
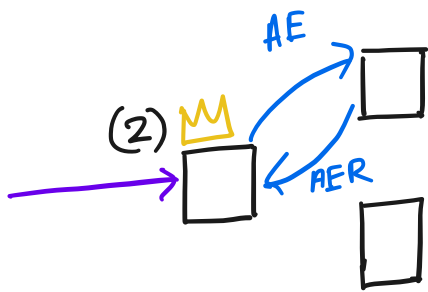
- After command  $c$  is committed at slot  $i$ , slot  $i$  always contains  $c$
  - All slots  $0..i-1$  contain committed commands
- Common to nearly all consensus protocols!
- Log Completeness
- Raft-Only

All of the portions of the Protocol are designed to maintain these properties.

- Replication (AE, AER)
- Leader Election (R.v.o, R.v.o.R.o)
- Reconfiguration

### Core Approach : Quorum Intersection

Any two changes to Rayt state (term, log, configuration) must have been signed off by a common node



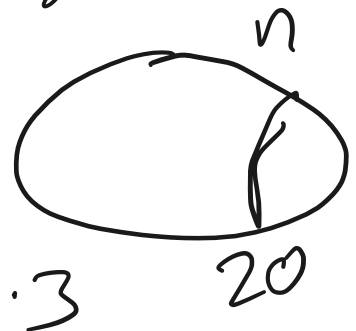
Common building block for all consensus protocols

Parameter:  $f$  = # of nodes that can fail.

Min quorum size?

$$n - f$$

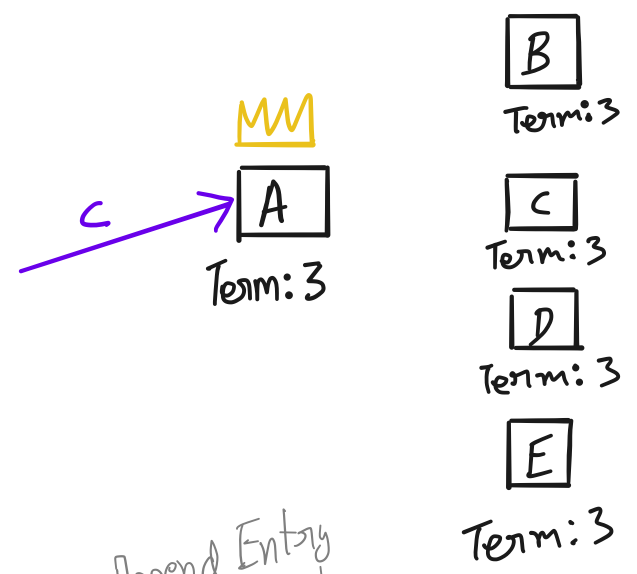
$$n - f + 1$$



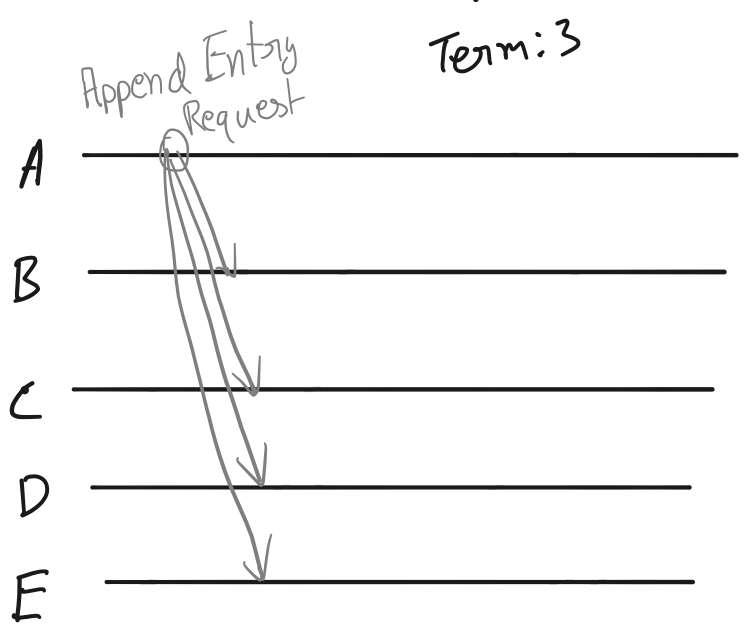
# Invariants

- Term:  $\mathbb{Z} : \emptyset$  or 1 leader. (Leader election ensures this)
- (term, index)  $\leftrightarrow$  command or  $\perp$
- In steady state, leader's log is authoritative.

# Replication/Log Synchronization



	0	1	2	3	4	5	6	7	8
A	To	T <sub>3</sub>							
B	To	T <sub>3</sub>							
C	To								
D	To								
E									



Append Entry Contents

- Term 3
- Command(s) c

Entries

- Last log index } Why? 1
- Last log term } 3
- CI

$\rightarrow [3, 1], c$

Processing at B



A	T <sub>0</sub>	T <sub>3</sub>	C						
B	T <sub>0</sub>	T <sub>3</sub>	C						
C	T <sub>0</sub>								
D	T <sub>0</sub>	T <sub>1</sub>							
E									

Processing at C

Processing responses at A (M)

- When is a command committed?

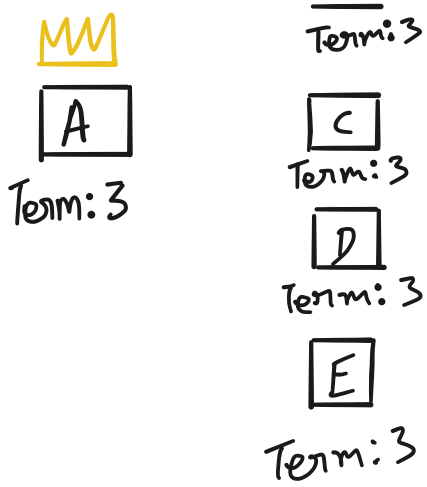
$$n = 2f + 1?$$

- When can A execute the command?

- When can D execute the command?

B

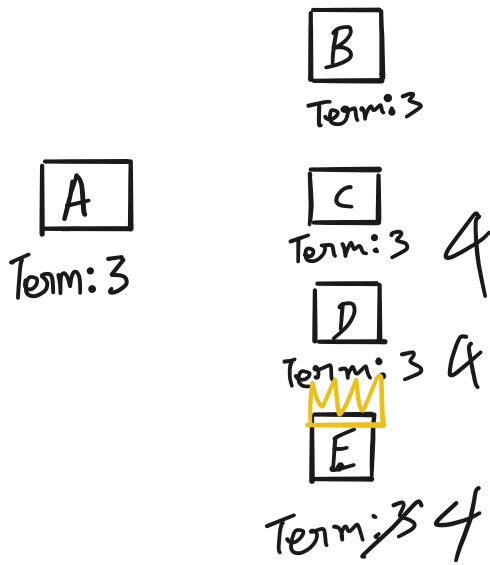
	0	1	2	3	4	5	6	7	8
A	T <sub>0</sub>	T <sub>1</sub>							



A	T <sub>0</sub>	T <sub>0</sub>							
B	T <sub>0</sub>	T <sub>0</sub>							
C	T <sub>0</sub>	T <sub>0</sub>							
D	T <sub>0</sub>								
E									

When can A execute command at Index 1?

Why?



	0	1	2	3	4	5	6	7	8
A	T <sub>0</sub>	T <sub>0</sub> '							
B	T <sub>0</sub>	T <sub>0</sub>							
C	T <sub>0</sub>	T <sub>0</sub>							
D	T <sub>0</sub>								
E	T <sub>0</sub>	T <sub>1</sub>							

Addendum: When is a command committed?

Leader Election

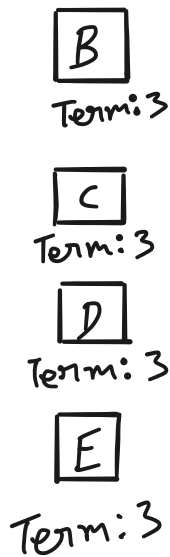
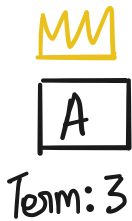
Goal: Allow progress if leader fails

# Conflict

- Stable leadership is necessary for progress. Why?

Avoid leader election when unnecessary

- Want to detect that a leader has failed soon. Why?



- Solution

- Heartbeats: no-op  
append entry requests

- Handling at Leader

- Handling at follower.

# Election Timeout

- Relation to Heartbeats?





A

Term: 3

Term: 3

C

Term: 3

D

Term: 3

E

Term: 3

## Leader Election

Requirement 1<sup>o</sup>: New leader's log should contain all committed log entries

Slight refinement of

After command  $c$  is committed at slot  $i$ , slot  $i$  always contains  $c$

Requirement 2<sup>o</sup>: At most one leader at a time.

Protocol assumes there is no preference for what node becomes leader



MM  
 A  
 Term: 3

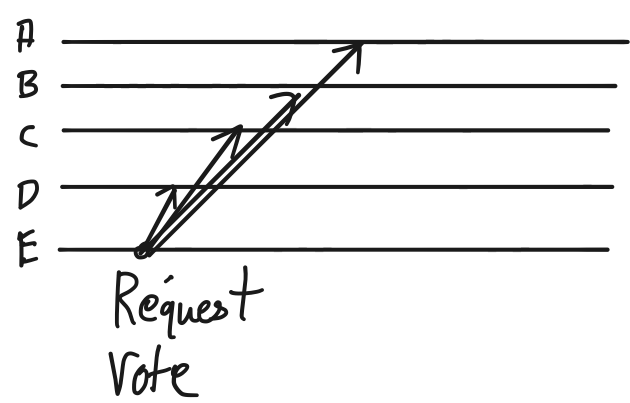
B  
 Term: 3/4

C  
 Term: 3

D  
 Term: 3

←→  
 E  
 Term: 4  
Candidate

	0	1	2	3	4	5	6	7	8
A	T <sub>0</sub>	T <sub>0</sub>							
B	T <sub>0</sub>	T <sub>0</sub>							
C	T <sub>0</sub>								
D	T <sub>0</sub>								
E	T <sub>0</sub>	T <sub>1</sub>							



Information Req Vote Req  
 term: 4  
 ID: E  
 Last Log Index: 1  
 Last Log Term: 1

What does a node (B) do on receiving R.V

B agrees to at most leader

(a) Update term (if received term is higher)  
 ↳ Consequence

(b) Check if B has already voted in term

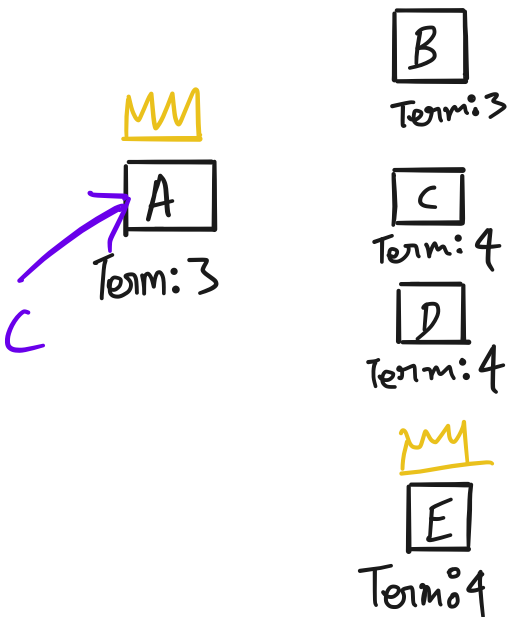
(c) Check who has more up-to-date log

Longer

Commit from a later term.

When does a candidate  $\rightarrow$  leader?

When does leader in term  $(t-1) \rightarrow$  FOLLOWER (in term  $t$ )?

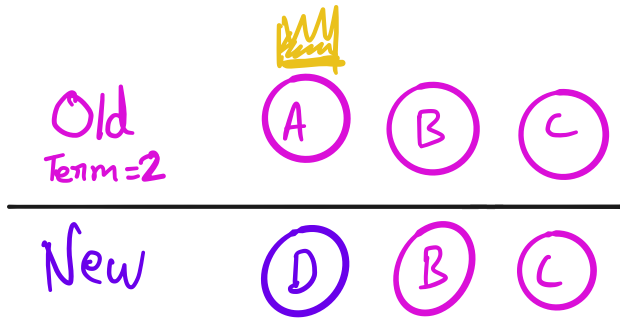


Re-configuration

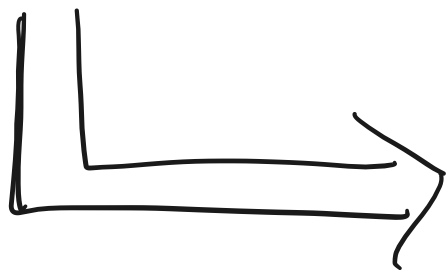
... participating

- Goal: Change the set of nodes participating

- Why hard? Need to pick a slightly complex setting



	0	1	2	3	4	5
A	t2					
B	t2					
C						
D						



Must ensure that D or C do not get elected. Why?

But A cannot continue as leader either. Why?

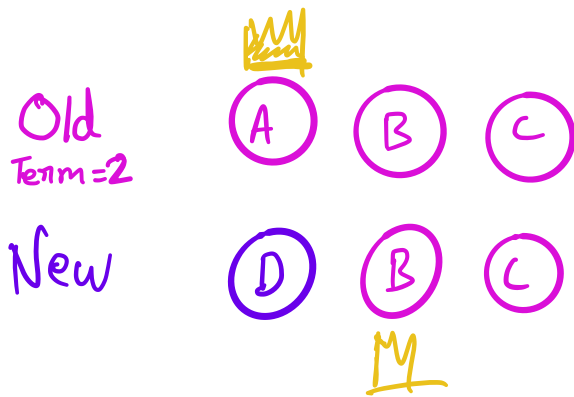
Things we need to fix

- (a) Synchronize logs so that a quorum of nodes in the new config contain committed entries.

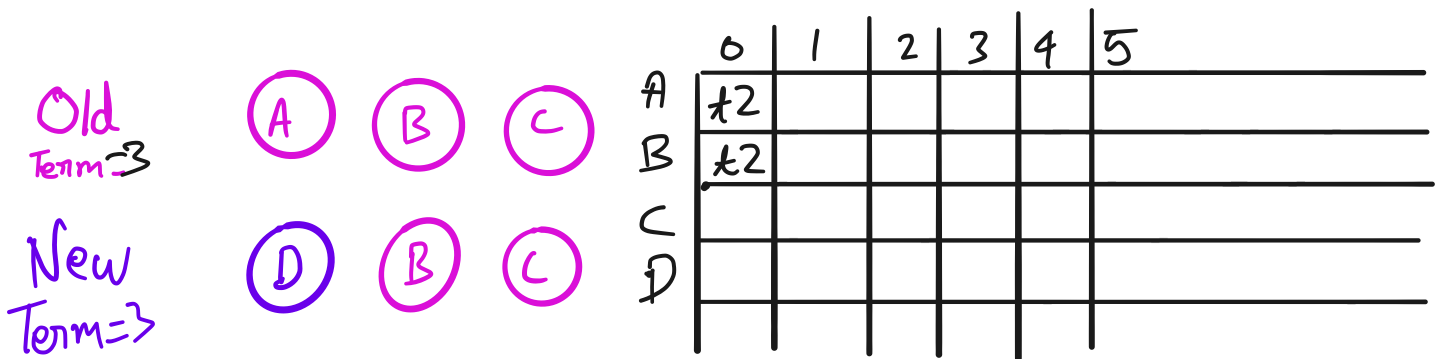
Observation: AE/AER synchronizes logs for a quorum of nodes so they match the leaders log.

Idea: ① Run leader election to find a node that

Joint Quorum {  
 (a) Can commit to nodes in the new configuration  
 (b) Has all committed entries



② Have the new leader commit an entry



③ Project?

Notes on practice

- Leader election

- Reconfiguration