Raft

October 2025

Raft

- System for enforcing strong consistency (linearizability)
- Similar to Paxos and Viewstamped Replication, but much **simpler**
- Clear boundary between leader election and the log consensus
- Leader log is ground truth; log entries only flow in one direction (from leader to followers)

Recap: Raft Leader election

Everyone sets a randomized timer that expires in [T, 2T] (e.g. T = 150ms)

When timer expires, increment term and send a RequestVote to everyone

Retry this until either:

- 1. You get majority of votes (including yourself): become leader
- 2. You receive an RPC from a valid leader: become follower again

Conditions for granting vote

- 1. (Assignment 3) We did not vote for anyone else in this term
- 2. (Assignment 3) Candidate term must be >= ours
- 3. (Assignment 4) Candidate log is at least as *up-to-date* as ours
 - The log with higher term in the last entry is more up-to-date
 - b. If the last entry terms are the same, then the longer log is more up-to-date

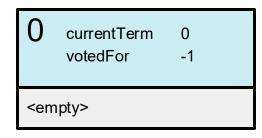
U	currentTerm	0	
	votedFor	-1	
	commitIndex	0	
	lastApplied	0	
	nextIndex	[]	
	matchIndex	[]	
(log entries here)			

Logs are 1-indexed

currentTermlatest term server has seenvotedForcandidate ID that received vote in current term,
or -1 if nonecommitIndexindex of highest log entry known to be committedlastAppliedindex of highest log entry applied to state machine

(Only on leader)

nextIndex	for each server, index of the next log entry to send to that server
matchIndex	for each server, index of highest log entry known to be replicated on the server



currentTerm latest term server has seen

votedFor candidate ID that received vote in current term,

or -1 if none

State required for election

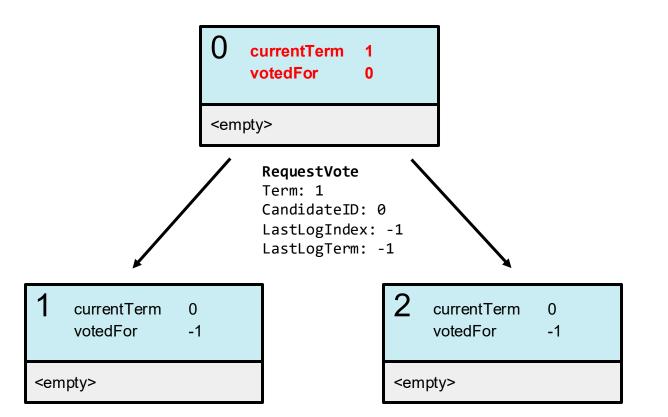
Scenario 1: During System Bootup

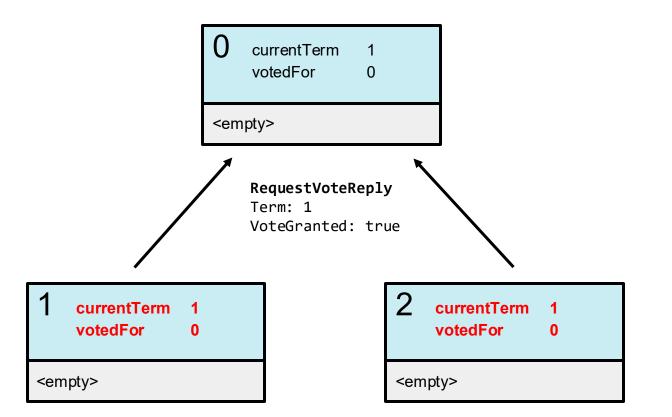
O currentTerm 0 votedFor -1

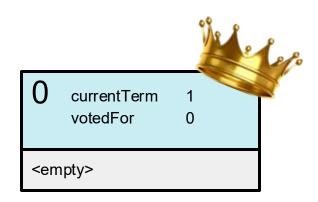
Timeout

1 currentTerm 0 votedFor -1 <empty>

2 currentTerm 0 votedFor -1 <empty>

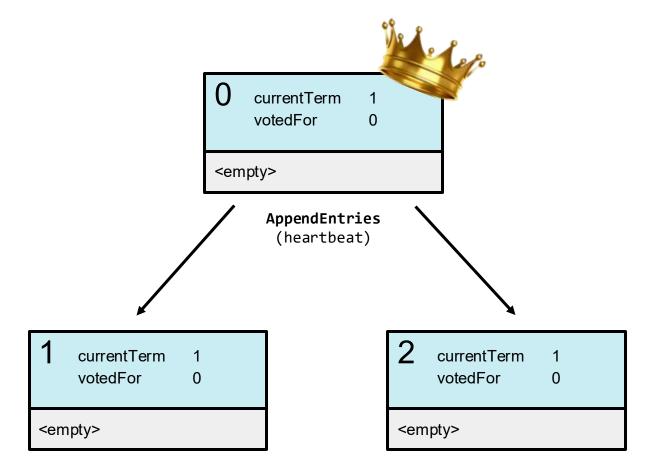






1 currentTerm 1 votedFor 0

2 currentTerm 1 votedFor 0 <empty>



Scenario 2: During Normal Execution

(suppose there are existing log entries...)

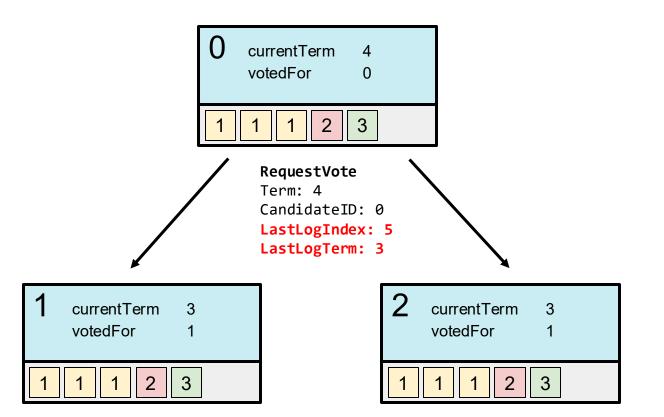
chance 2. Banng Monnai Excoation

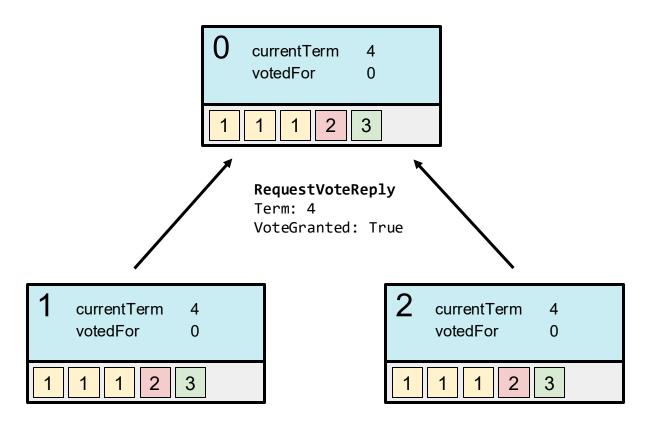
O currentTerm 3 votedFor 1

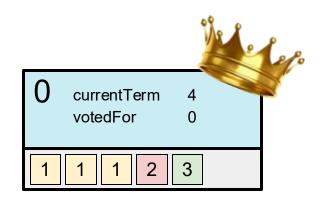
Timeout

1 currentTerm 3 votedFor 1

2 currentTerm 3 votedFor 1 1 1 2 3







1 currentTerm 4 votedFor 0

2 currentTerm 4 votedFor 0

Conditions for granting vote

- 1. We did not vote for anyone else in this term
- 2. Candidate term must be >= ours
- 3. Candidate log is at least as *up-to-date* as ours
 - a. The log with higher term in the last entry is more up-to-date
 - b. If the last entry terms are the same, then the longer log is more up-to-date

Which one is more *up-to-date*?

1 1 1 2 3

1 | 1 | 1 | 1 | 1 | 1 | 1

Which one is more *up-to-date*?

1 1 1 2 3

 1
 1
 1
 2
 3
 3
 3

Which one is more *up-to-date*?

1 1 1 2 3

1 1 4

Why reject logs that are not *up-to-date*?

Leader log is always the ground truth

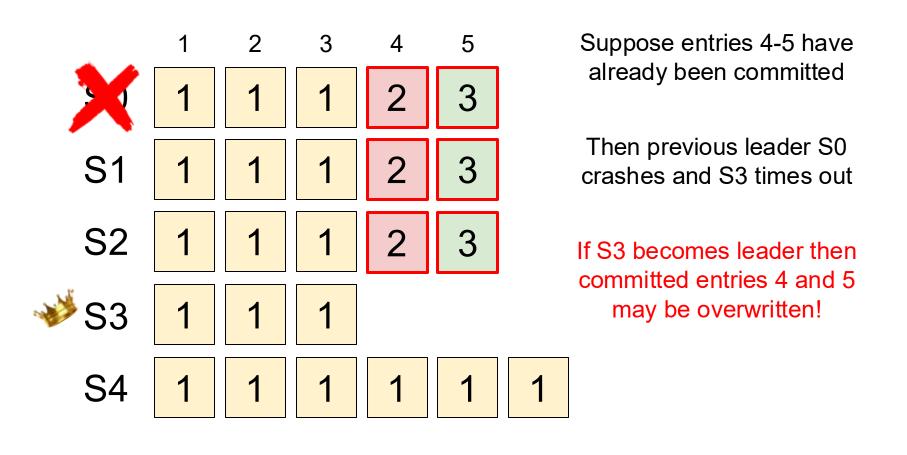
Once someone is elected leader, followers must throw away conflicting entries

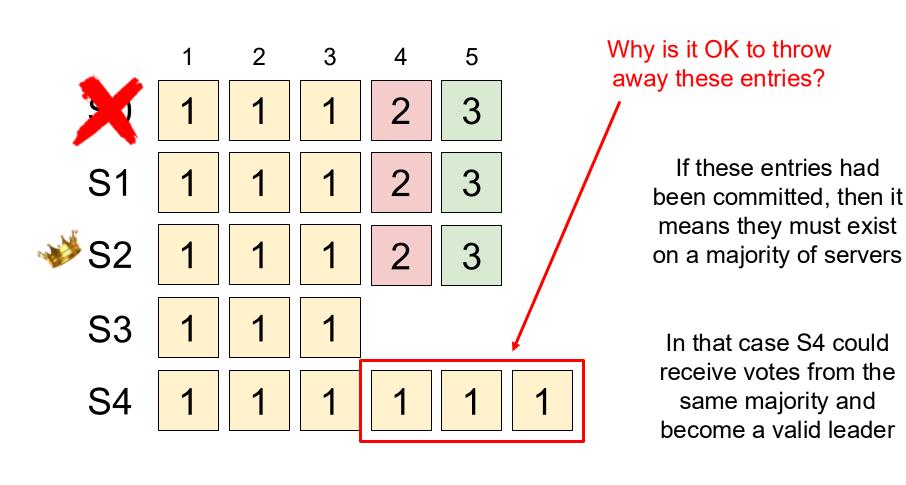
Must NOT throw away committed entries!

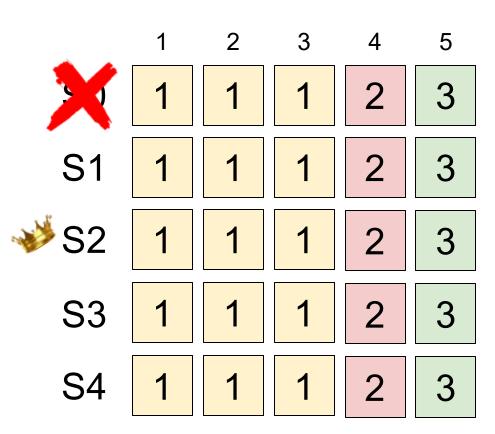
Note: Log doesn't need to be the MOST up-to-date among all servers

What if we accept logs that are not as

up-to-date as ours?







Raft Normal Operation

n	currentTerm	0		
	votedFor	-1		
	commitIndex	0		
	lastApplied	0		
	nextIndex	[]		
	matchIndex	[]		
<empty></empty>				

Logs are 1-indexed

 currentTerm
 latest term server has seen

 votedFor
 candidate ID that received vote in current term, or -1 if none

 commitIndex
 index of highest log entry known to be committed

 lastApplied
 index of highest log entry applied to state machine

(Only on leader)

nextIndex for each server, index of the next log entry to send

to that server

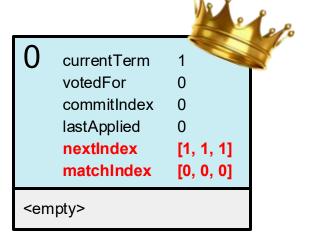
matchindex for each server, index of highest log entry known to

be replicated on the server

```
O currentTerm 0
votedFor -1
commitIndex 0
lastApplied 0
nextIndex []
matchIndex []
```

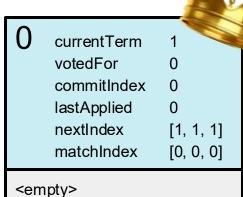
```
1 currentTerm 0
votedFor -1
commitIndex 0
lastApplied 0
nextIndex []
matchIndex []
<empty>
```

```
2 currentTerm 0
votedFor -1
commitIndex 0
lastApplied 0
nextIndex []
matchIndex []
<empty>
```



```
1 currentTerm 1
votedFor 0
commitIndex 0
lastApplied 0
nextIndex []
matchIndex []
```

```
2 currentTerm 1
votedFor 0
commitIndex 0
lastApplied 0
nextIndex []
matchIndex []
```



AppendEntries

Term: 1 LeaderID: 0

PrevLogIndex: 0

PrevLogTerm: -1 LeaderCommit: 0

currentTerm votedFor commitIndex **lastApplied** nextIndex [] matchIndex

<empty>

```
AppendEntries
```

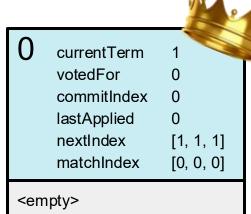
Term: 1 LeaderID: 0

PrevLogIndex: 0

PrevLogTerm: -1 LeaderCommit: 0 currentTerm votedFor 0 commitIndex 0 **lastApplied** 0

nextIndex [] matchIndex []

<empty>



AppendEntriesReply

Term: 1

Success: True

1 currentTerm 1 votedFor 0 commitIndex 0 lastApplied 0 nextIndex [] matchIndex []

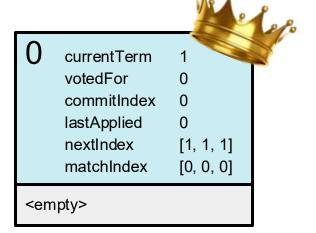
AppendEntriesReply

Term: 1

Success: True

2 currentTerm 1
votedFor 0
commitIndex 0
lastApplied 0
nextIndex []
matchIndex []

<empty>

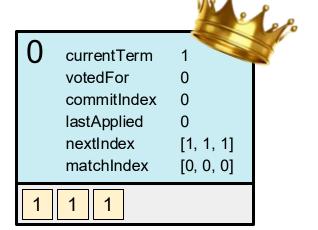


1	currentTerm	1		
	votedFor	0		
	commitIndex	0		
	lastApplied	0		
	nextIndex	[]		
	matchIndex	[]		
<empty></empty>				

Request 1

Client

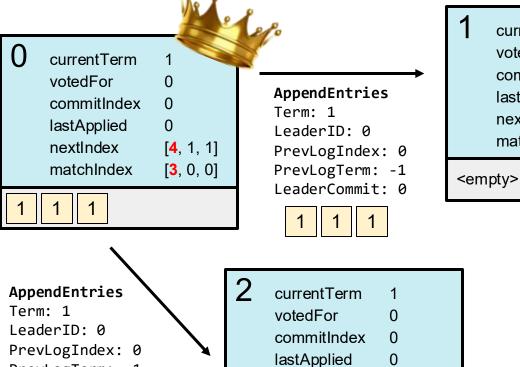
2 currentTerm 1
votedFor 0
commitIndex 0
lastApplied 0
nextIndex []
matchIndex []



1 currentTerm 1 votedFor 0 commitIndex 0 lastApplied 0 nextIndex [] matchIndex []

Client Request 1
Request 2
Request 3

2 currentTerm 1
votedFor 0
commitIndex 0
lastApplied 0
nextIndex []
matchIndex []



nextIndex

<empty>

matchIndex

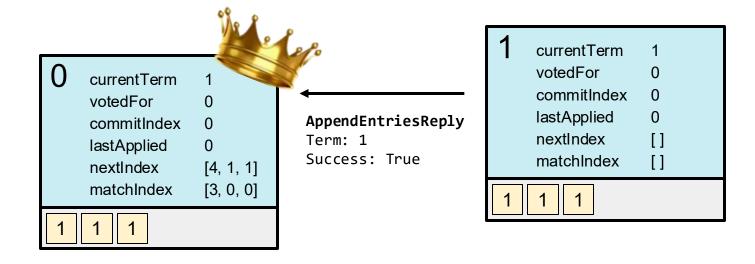
[]

[]

PrevLogTerm: -1

LeaderCommit: 0

1 currentTerm 1 votedFor 0 commitIndex 0 lastApplied 0 nextIndex [] matchIndex []

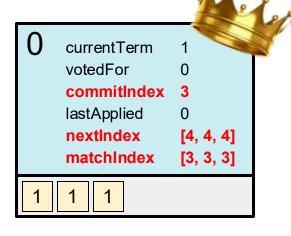


AppendEntriesReply

Term: 1

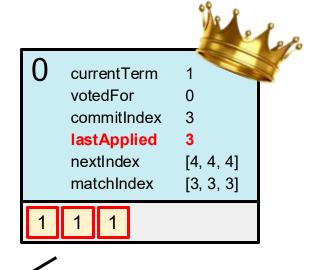
Success: True

2 currentTerm 1
votedFor 0
commitIndex 0
lastApplied 0
nextIndex []
matchIndex []



Entry 3 is now replicated on a majority, so we can commit it

while commitIndex > lastApplied,
apply commands to state machine



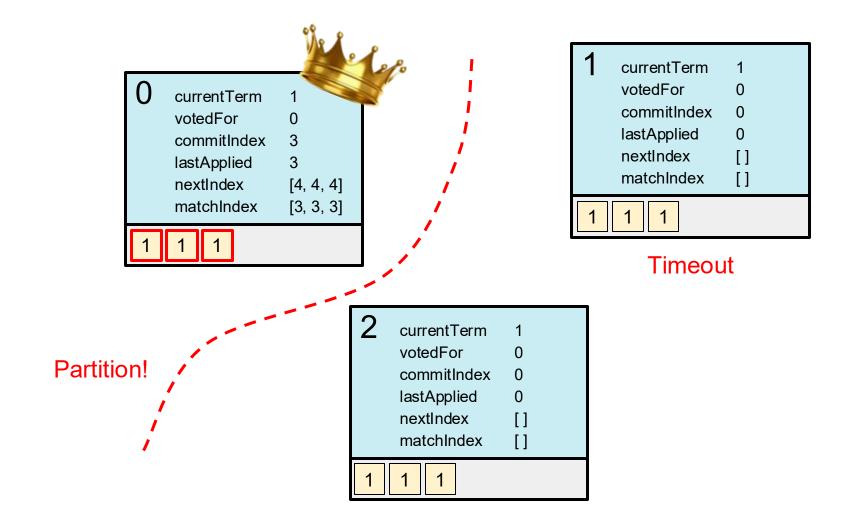
Once leader has applied an entry to state machine, it is safe to tell the client that the entry is committed

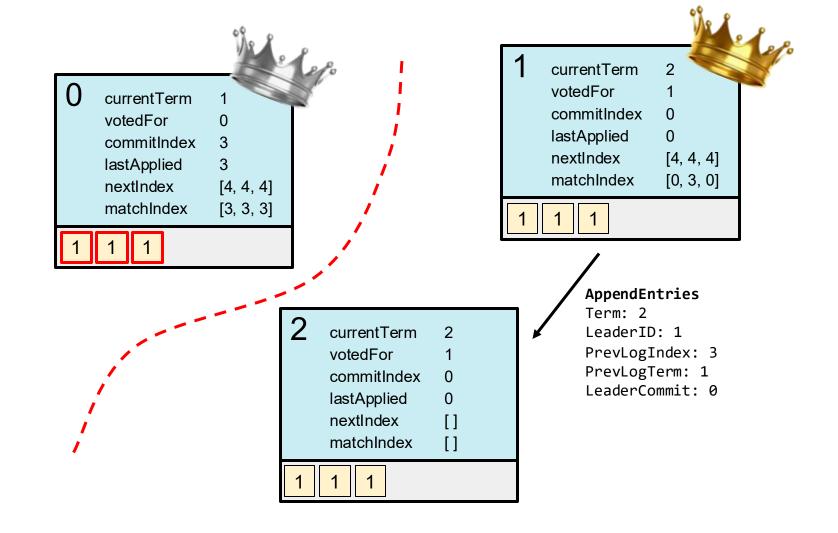
Response 1 2 3

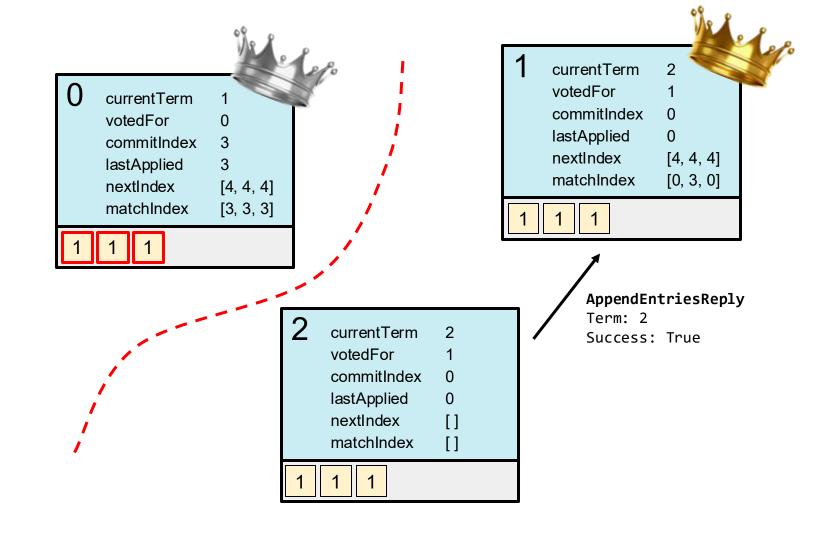
Client

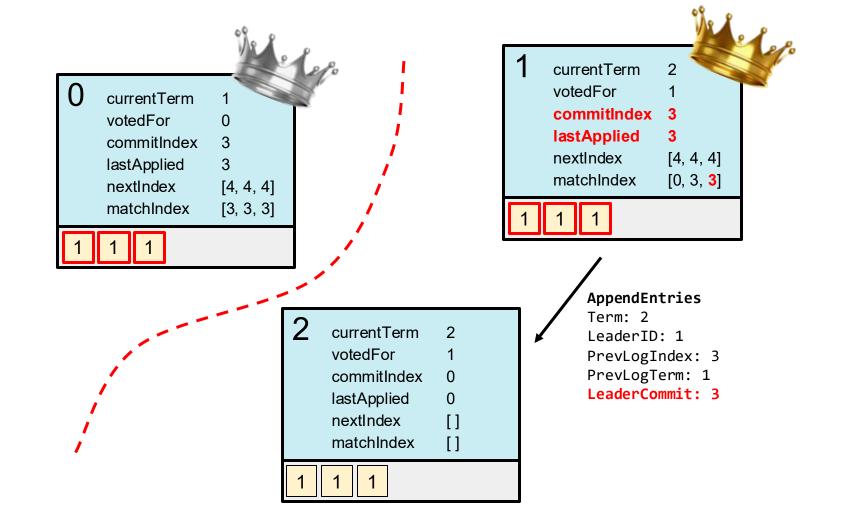
Raft

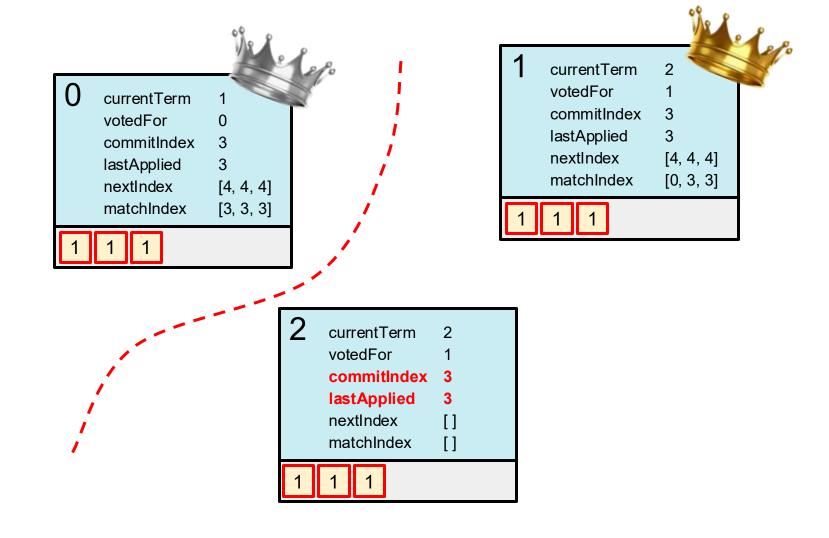
After new leader election

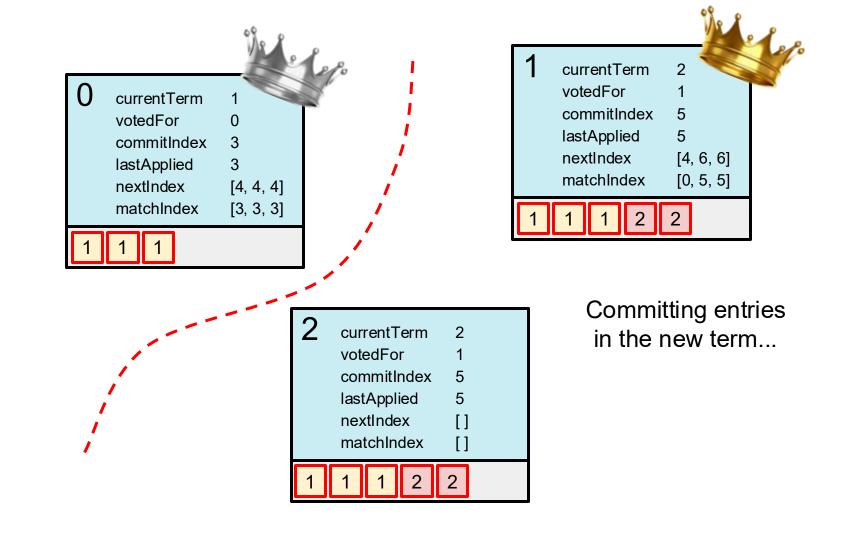




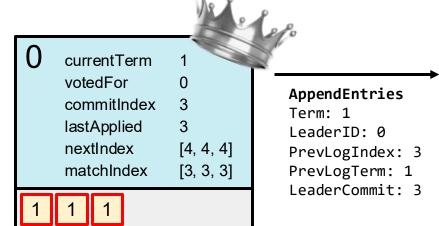


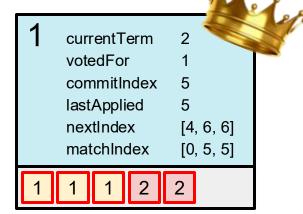






Later, the network partition is fixed ...



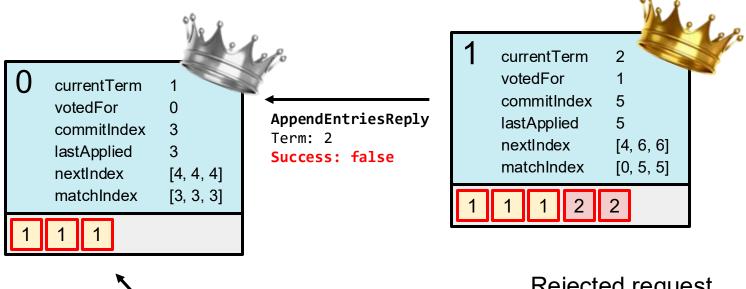


AppendEntries
Term: 1

LeaderID: 0

PrevLogIndex: 3
PrevLogTerm: 1
LeaderCommit: 3

2 currentTerm 2 votedFor 1 commitIndex 5 lastApplied 5 nextIndex [] matchIndex []



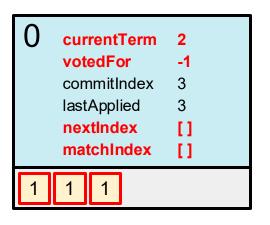
AppendEntriesReply

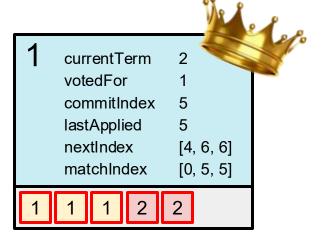
Term: 2

Success: false

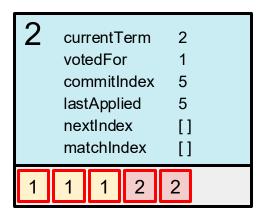
2 currentTerm 2 votedFor 1 commitIndex 5 lastApplied 5 nextIndex [] matchIndex []

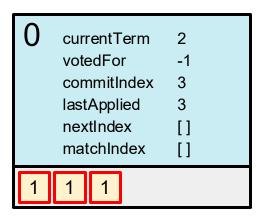
Rejected request because local term is higher (2 > 1)





Old leader is dethroned!





AppendEntries

Term: 2

LeaderID: 1

PrevLogIndex: 3
PrevLogTerm: 1

LeaderCommit: 5

2 2

 1
 currentTerm
 2

 votedFor
 1

 commitIndex
 5

 lastApplied
 5

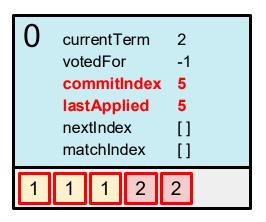
 nextIndex
 [4, 6, 6]

 matchIndex
 [0, 5, 5]

 1
 1
 2

 2
 2

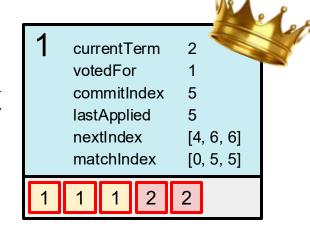
2 currentTerm 2 votedFor 1 commitIndex 5 lastApplied 5 nextIndex [] matchIndex []



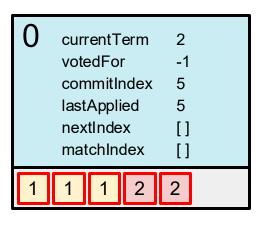
AppendEntriesReply

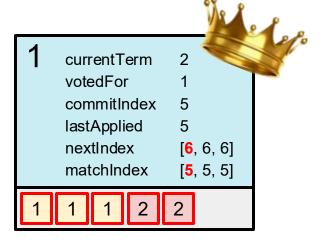
Term: 2

Success: true

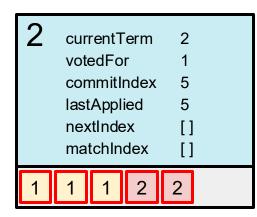


2 currentTerm 2 votedFor 1 commitIndex 5 lastApplied 5 nextIndex [] matchIndex []





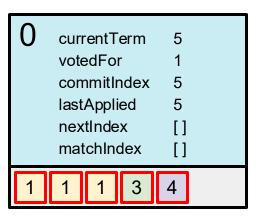
Everyone is on the same page again

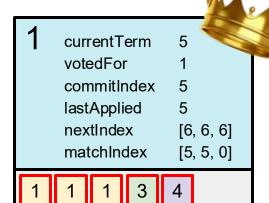


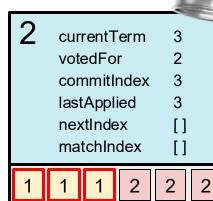
When log entries don't match...

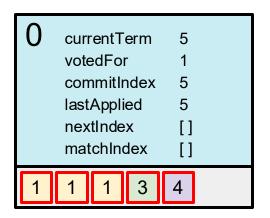
When log entries don't match...

- The leader will find the latest log entry in the follower where the two logs agree
- At the follower:
 - Everything after that entry will be deleted
 - The leader's log starting from that entry will be replicated on the follower



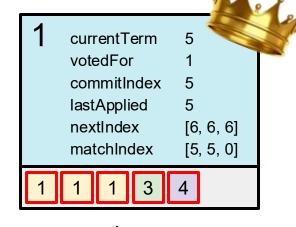






prevLogIndex = 5
 S1 log[5] = 4
 S2 log[5] = 2

Mismatch!



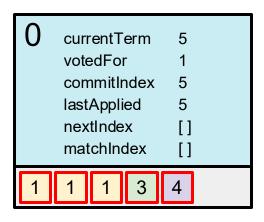
2 currentTerm 3 votedFor 2 commitIndex 3 lastApplied 3 nextIndex [] matchIndex []

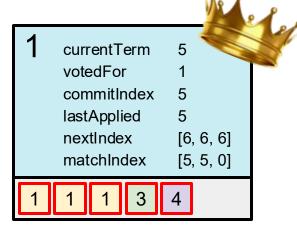
AppendEntries

Term: 5

LeaderID: 1

PrevLogIndex: 5
PrevLogTerm: 4
LeaderCommit: 5



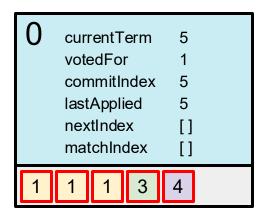


2 currentTerm 5
votedFor -1
commitIndex 3
lastApplied 3
nextIndex []
matchIndex []

AppendEntriesReply

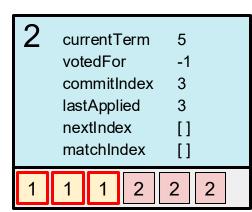
Term: 5

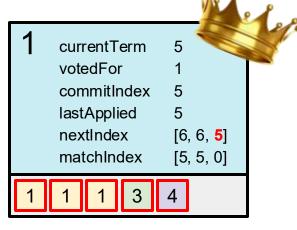
Success: False



prevLogIndex = 4
 S1 log[4] = 3
 S2 log[4] = 2

Mismatch!





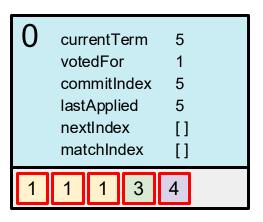
AppendEntries

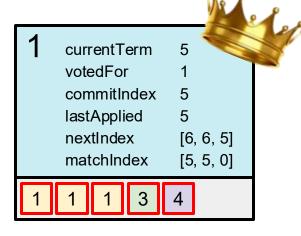
Term: 5

LeaderID: 1
PrevLogIndex: 4

PrevLogTerm: 3
LeaderCommit: 5

4



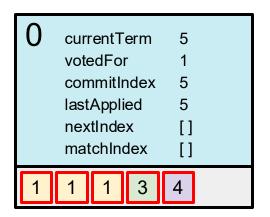


2 currentTerm 5
votedFor -1
commitIndex 3
lastApplied 3
nextIndex []
matchIndex []

AppendEntriesReply

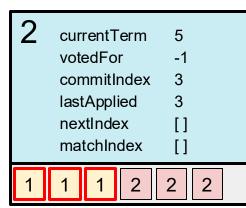
Term: 5

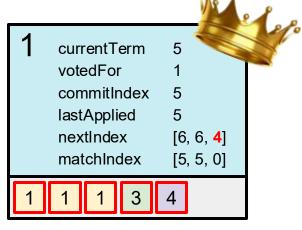
Success: False



```
prevLogIndex = 3
    S1 log[3] = 1
    S2 log[3] = 1
```

Match!



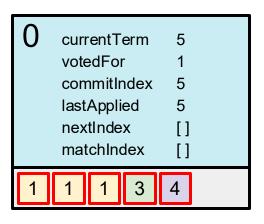


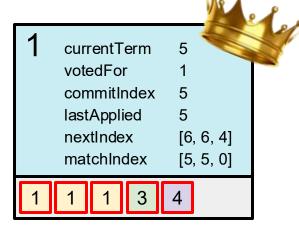
AppendEntries

Term: 5 LeaderID: 1

PrevLogIndex: 3
PrevLogTerm: 1
LeaderCommit: 5

3 4

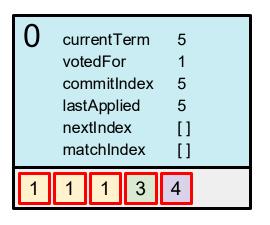


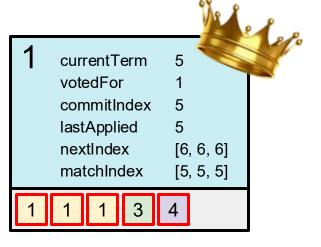


2 currentTerm 5 votedFor -1 commitIndex 5 lastApplied 5 nextIndex [] matchIndex []

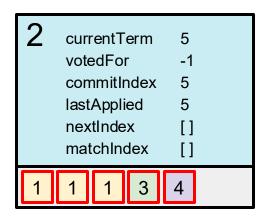
AppendEntriesReply

Term: 5
Success: True





Everyone is on the same page again



number of messages?

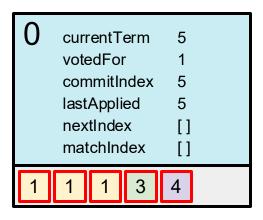
Optimization to reduce

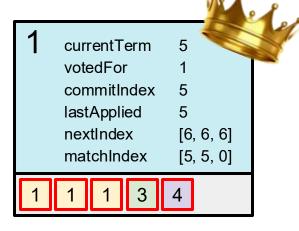
Key Idea

- Reduce the number of rejected AppendEntries RPCs
- One RPC per conflicting term, rather than one RPC per conflicting entry

Detailed Algorithm:

- When rejecting an AppendEntries request, the follower can include the term
 of the conflicting entry and the first index it stores for that term.
- With this information, the leader can decrement nextIndex to bypass all of the conflicting entries in that term.
- See page 7-8 in <u>Raft (extended version)</u>



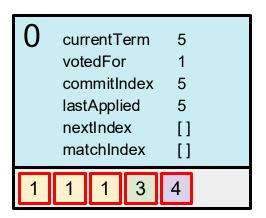


2 currentTerm 3 votedFor 2 commitIndex 3 lastApplied 3 nextIndex [] matchIndex []

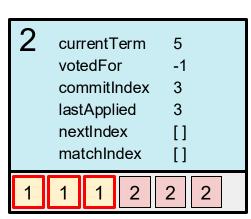
 ${\bf AppendEntries}$

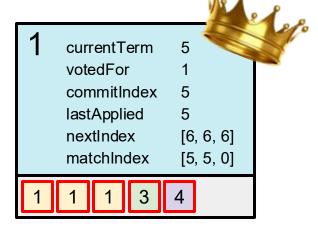
Term: 5 LeaderID: 1

PrevLogIndex: 5
PrevLogTerm: 4
LeaderCommit: 5



Specify the conflicting term and the first index of this term





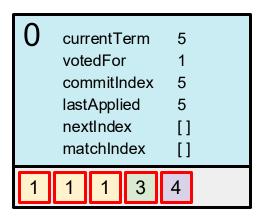
AppendEntriesReply

Term: 5

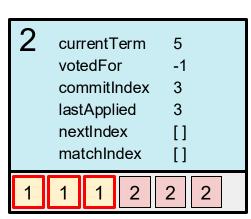
Success: False

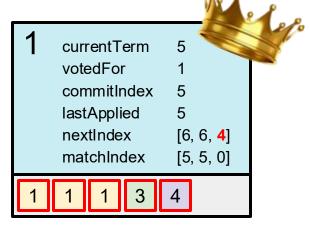
ConflictTerm: 2

ConflictFirstIndex: 4



Leader sends its log entries that are different from the follower's starting the specified conflicting term



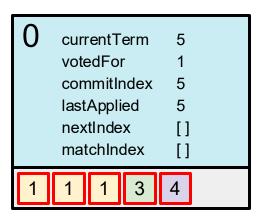


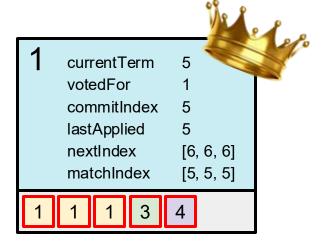
AppendEntries

Term: 5

LeaderID: 1
PrevLogIndex: 3
PrevLogTerm: 1
LeaderCommit: 5

3 4





2 currentTerm 5
votedFor -1
commitIndex 5
lastApplied 5
nextIndex []
matchIndex []

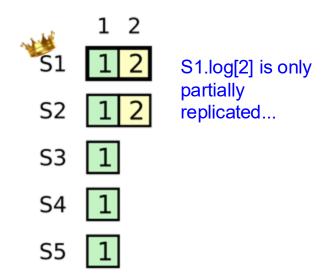
Key Idea:
Decrement nextIndex
one term at a time

Conditions for committing an entry

- 1. The entry exists on a majority AND was appended to leader in the current term
- 2. Or, the entry precedes another entry that is committed

Caveat for committing old entries

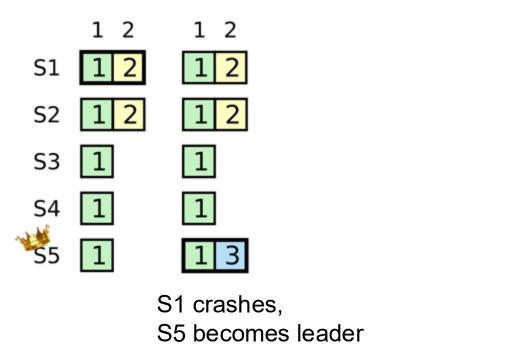
Can't assume an old entry has been committed *even if* it exists on a majority



S1 is the leader

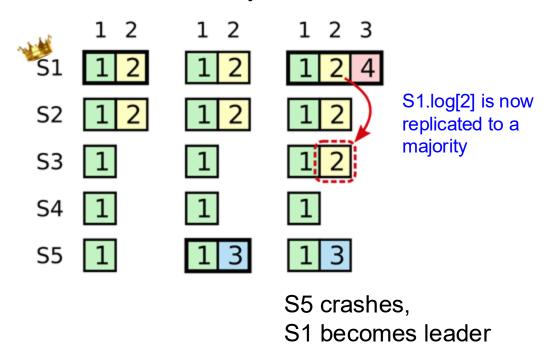
Caveat for committing old entries

Can't assume an old entry has been committed *even if* it exists on a majority



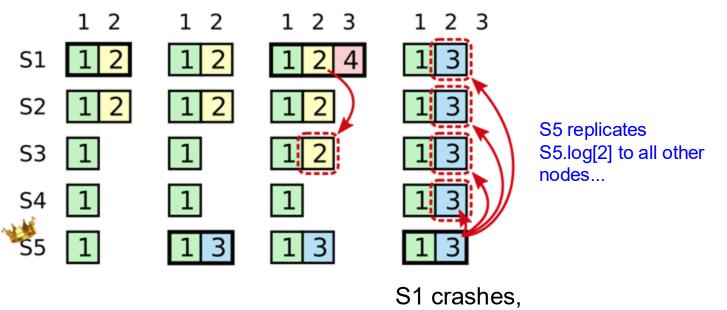
Caveat for committing old entries

Can't assume an old entry has been committed *even if* it exists on a majority



Caveat for committing old entries

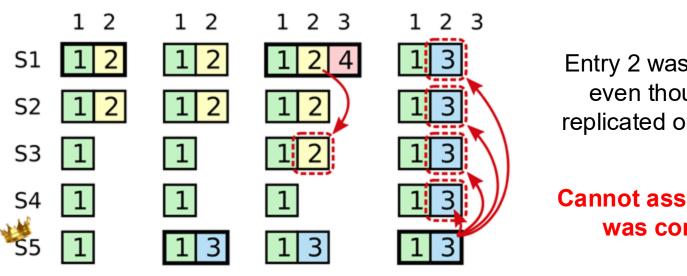
Can't assume an old entry has been committed even if it exists on a majority



S1 crasnes, S5 becomes leader

Caveat for committing old entries

Can't assume an old entry has been committed even if it exists on a majority

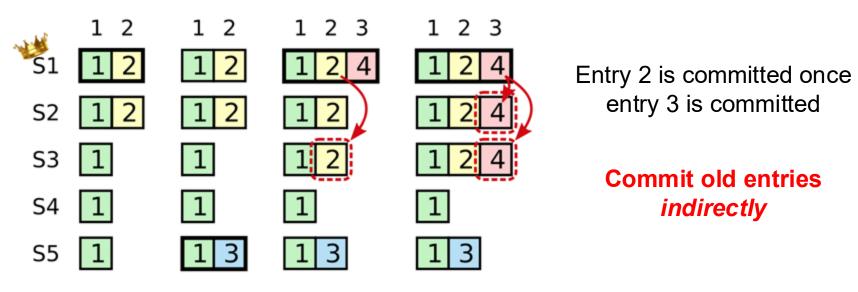


Entry 2 was overwritten even though it was replicated on a majority!

Cannot assume entry 2 was committed

Caveat for committing old entries

Can't assume an old entry has been committed even if it exists on a majority

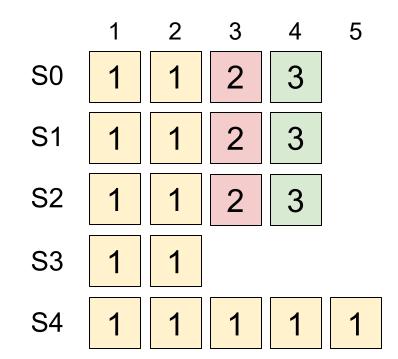


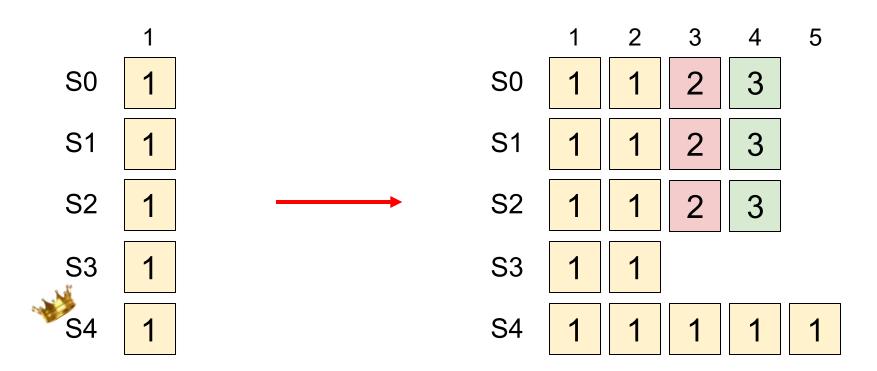
S1 commits entry 3

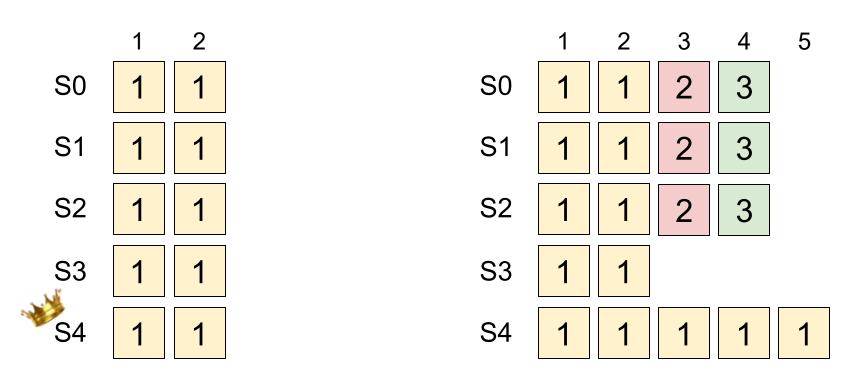
Exercise...

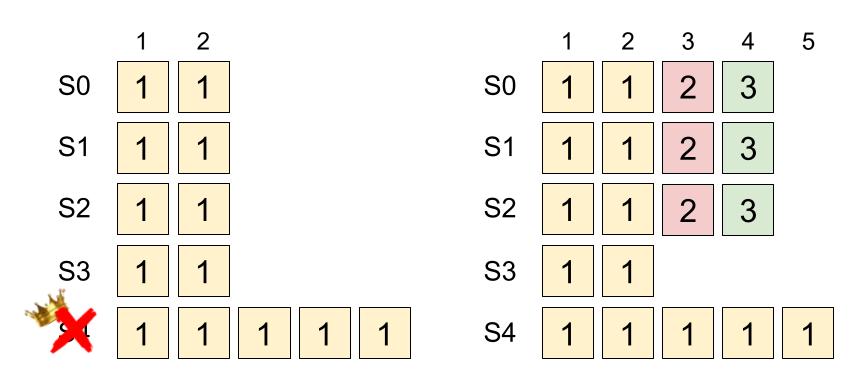
Rules for deciding which log is more up-to-date:

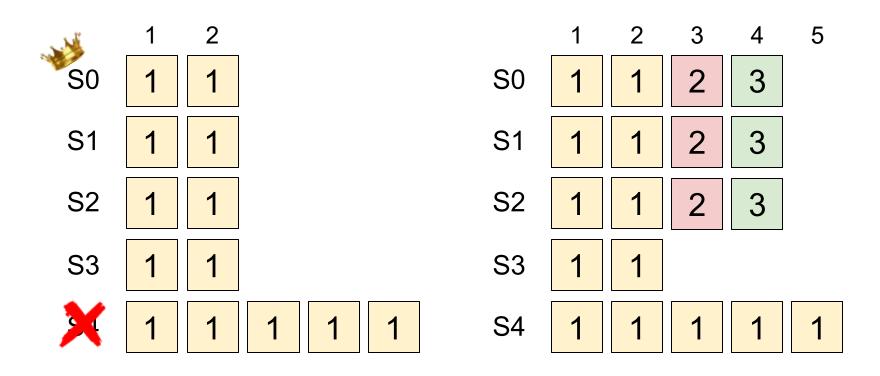
- Compare index and term of last entries in the logs
- If the terms are different: log with later term is more up-todate
- If the terms are the same: longer log is more up-to-date

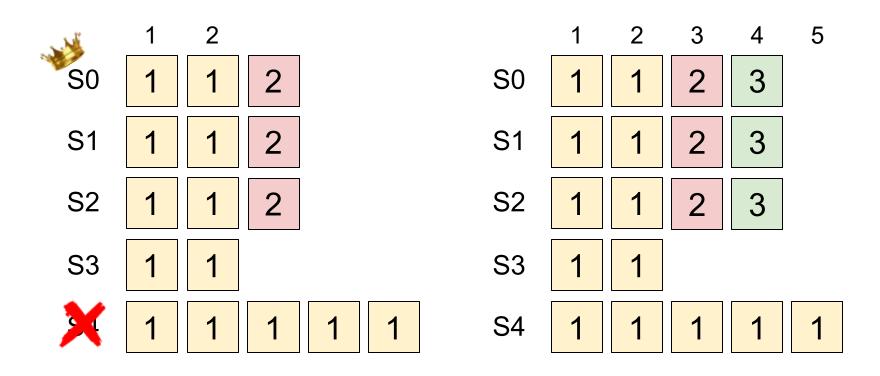


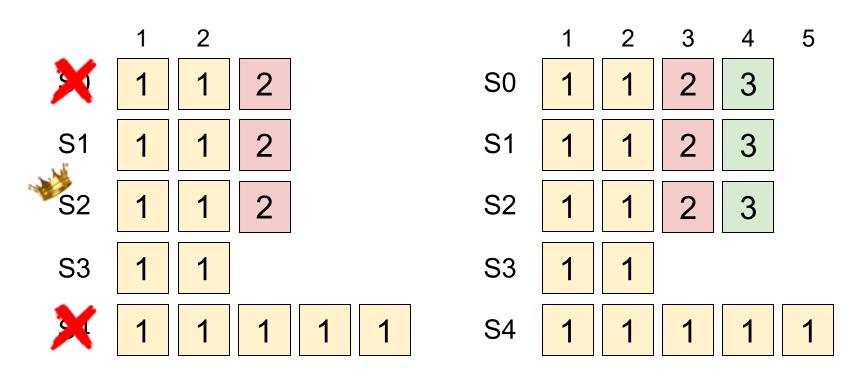


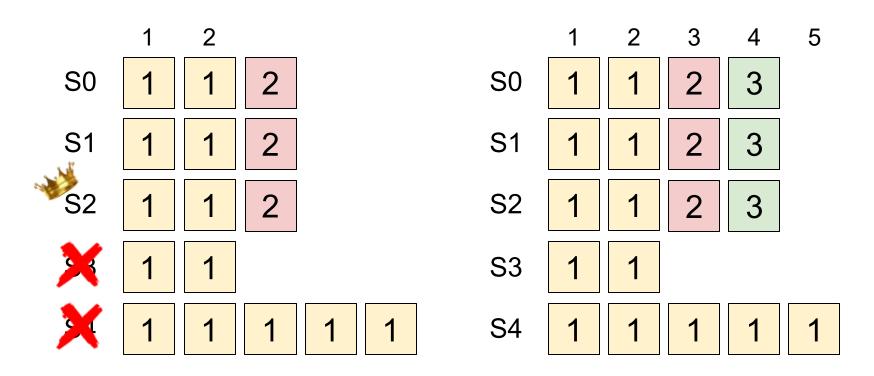


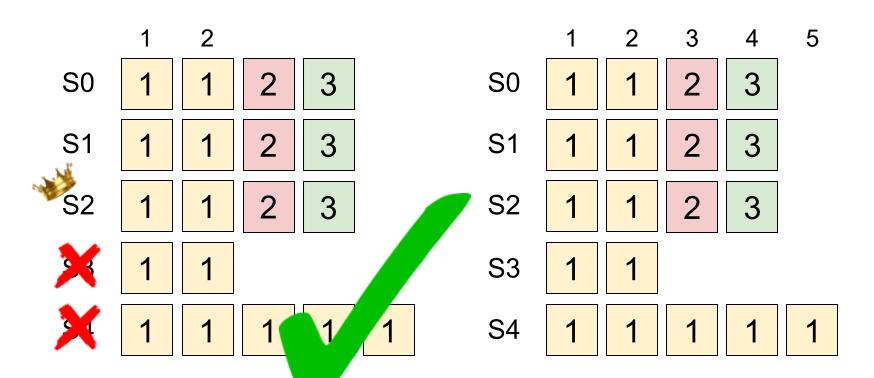


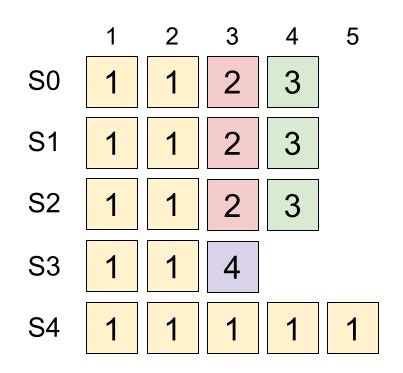






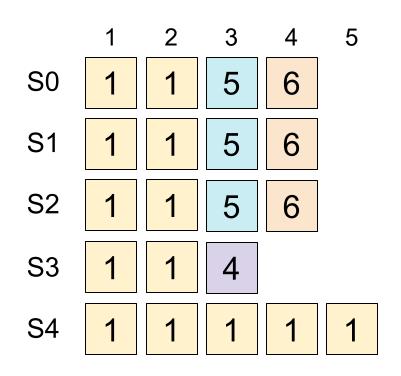






NO

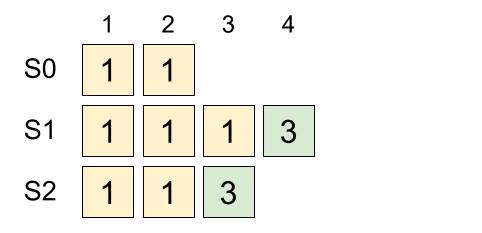
S3 cannot become leader in term 4 (Who's going to vote for him?)



Yes

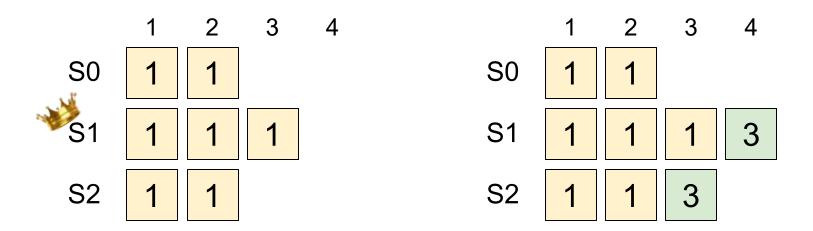
What happened to terms 2 and 3?

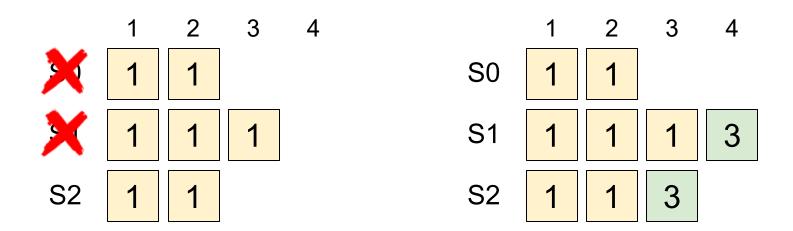
- 1. Split vote: no one became leader
- 2. Partitions: no one became leader
- 3. Simply no requests in these terms



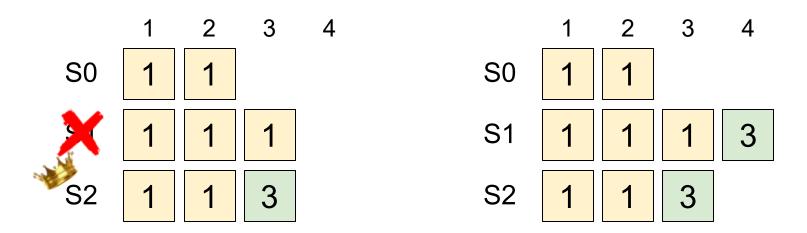
NO!

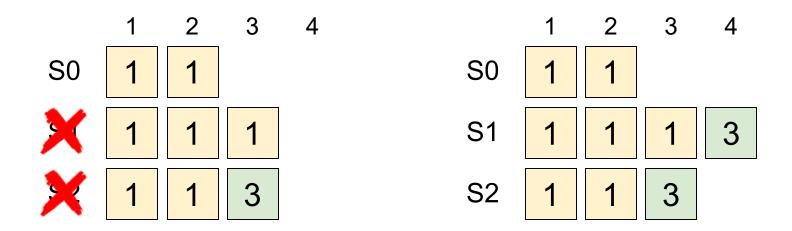
Let's try tracing the steps...

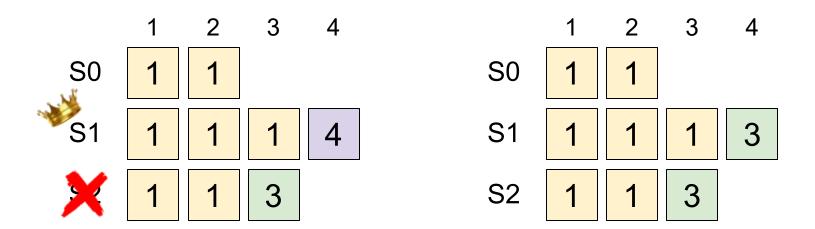




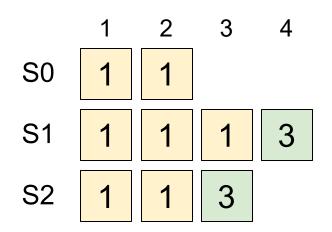
No one becomes leader in term 2...







S0 previously voted for S2 in term 3 S0 can only vote for S1 for term 4!



The two entries in term 3 are in different positions

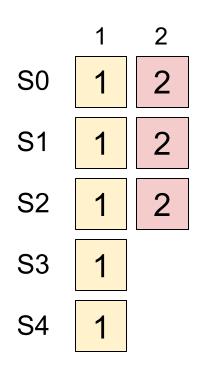
S1 and S2 could not have written these entries without being leaders

But they can't both be leaders in the same term!

Conditions for committing an entry

- 1. The entry exists on a majority AND was appended to leader in the current term
- 2. Or, the entry precedes another entry that is committed

Q5: Is entry 2 (term 2) guaranteed to be committed?

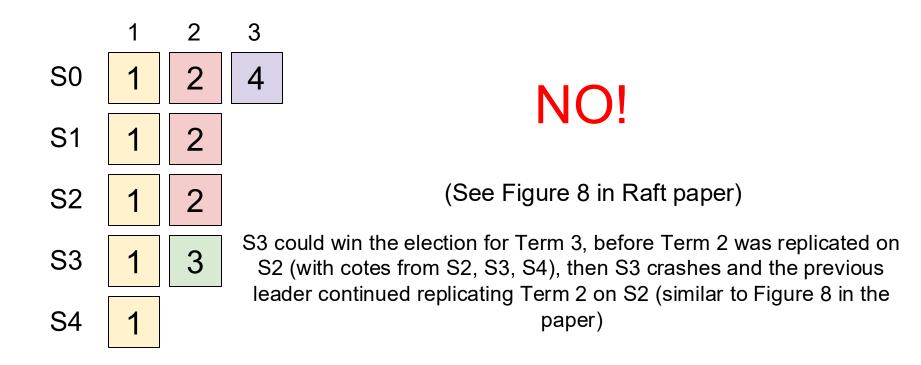


Yes!

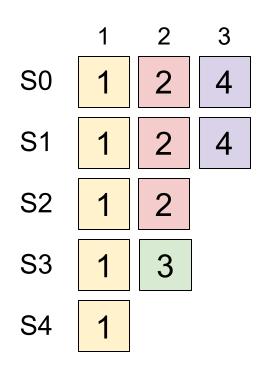
Entry 2 is on a majority of nodes

No one else has a more *up-to-date* log

Q6: Is entry 2 (term 2) guaranteed to be committed?



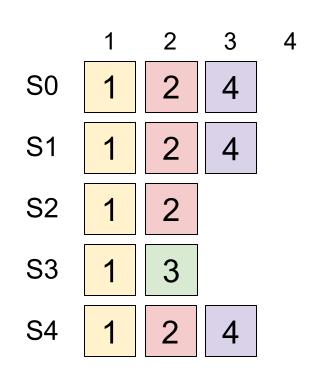
Q7: Is entry 2 (term 2) guaranteed to be committed?



NO!

S3 could still become leader if S0 crashes (votes from S2, S3 and S4)

Q8: Is entry 2 (term 2) guaranteed to be committed?



Yes!

Entry 4 is guaranteed to be committed because no one else has a more *up-to-date* log, *and* majority has entry 4

All entries before entry 4 are safe