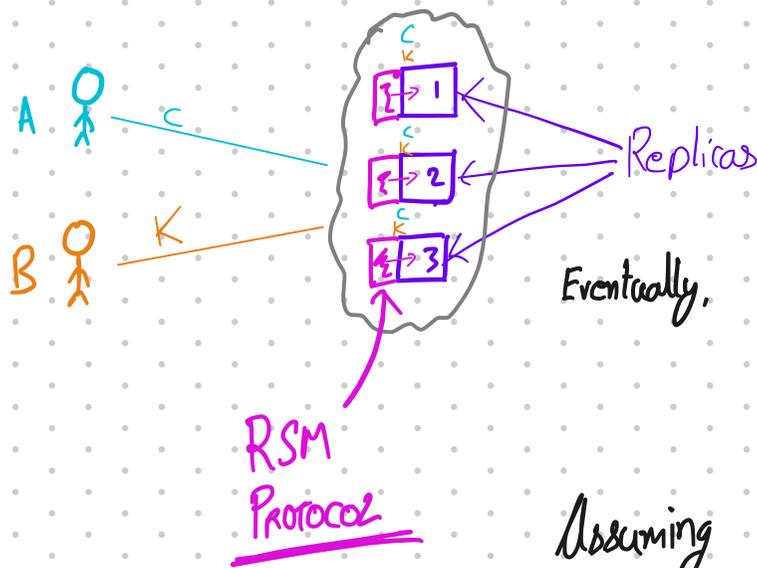


REPLICATED
 STATE
 MACHINES } RAFT



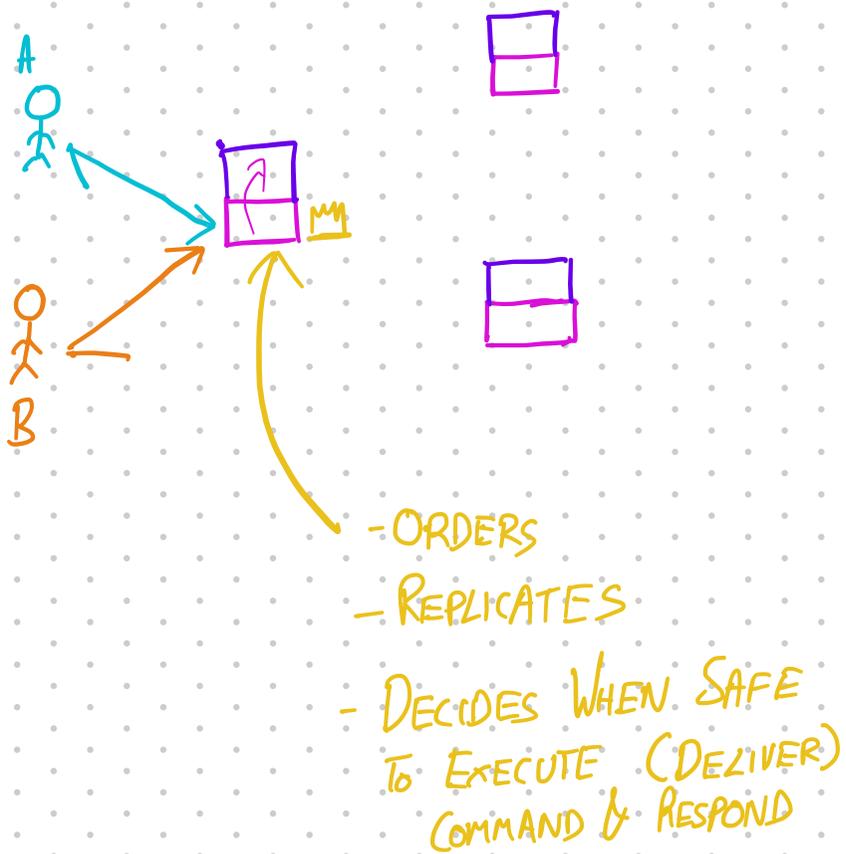
LAST CLASS:

- RSMs: A recipe for building fault tolerant systems



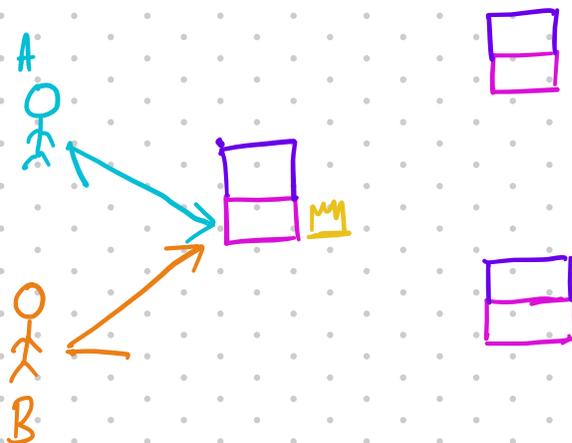
- Correctness requirements
 - Agreement
 - Ordering

- Leader based RSM



CORE CONCERNS

① WHEN IS IT SAFE TO EXECUTE/DELIVER?
[Commit Point]

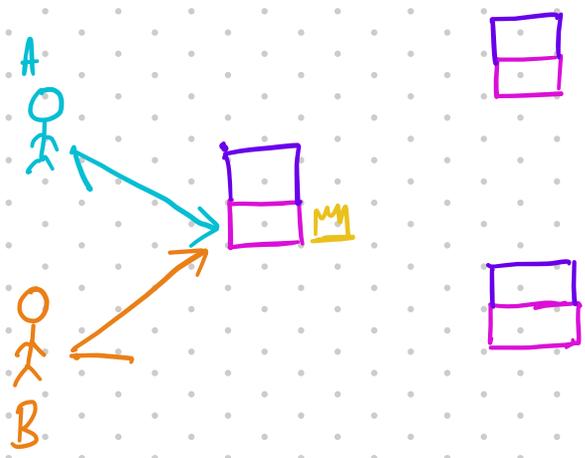


REQUIREMENTS :

③ How to detect leader failure &

④ who becomes leader after?

② who becomes leader of nodes
[Election]



REQUIREMENTS:

TARGET CONSISTENCY MODEL:

PRACTICAL CONCERNS WE WILL IGNORE

- READS ARE EXPENSIVE

- CRASH RECOVERY

RAFT: AN RSM PROTOCOL

Failure model:

Fail-stop ;
(Really fail-recover, but needs more)

$$< \frac{N}{2}$$

PROCESSES CAN FAIL
(N: Number of nodes)

USEFUL FOR FIGURING OUT COMMIT POINT, etc.

CORE CONCEPT: QUORUM INTERSECTION.

Two operations: store (val) → bool
load () → v

Goal: If store succeeds any subsequent load returns val [Safety]

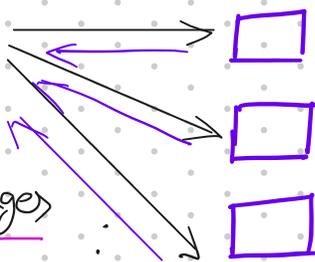
[Note: Not guaranteeing termination]
Failure model: at most f fail

< f fail
> f fail

Consider n=3 f=1

store(val)

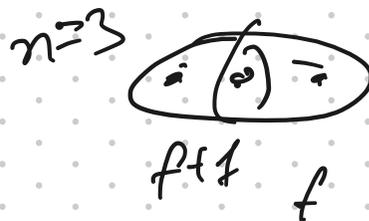
broadcast (store, val)
wait for ≥ f+1 messages
if |done| ≥ f+1
return true
else

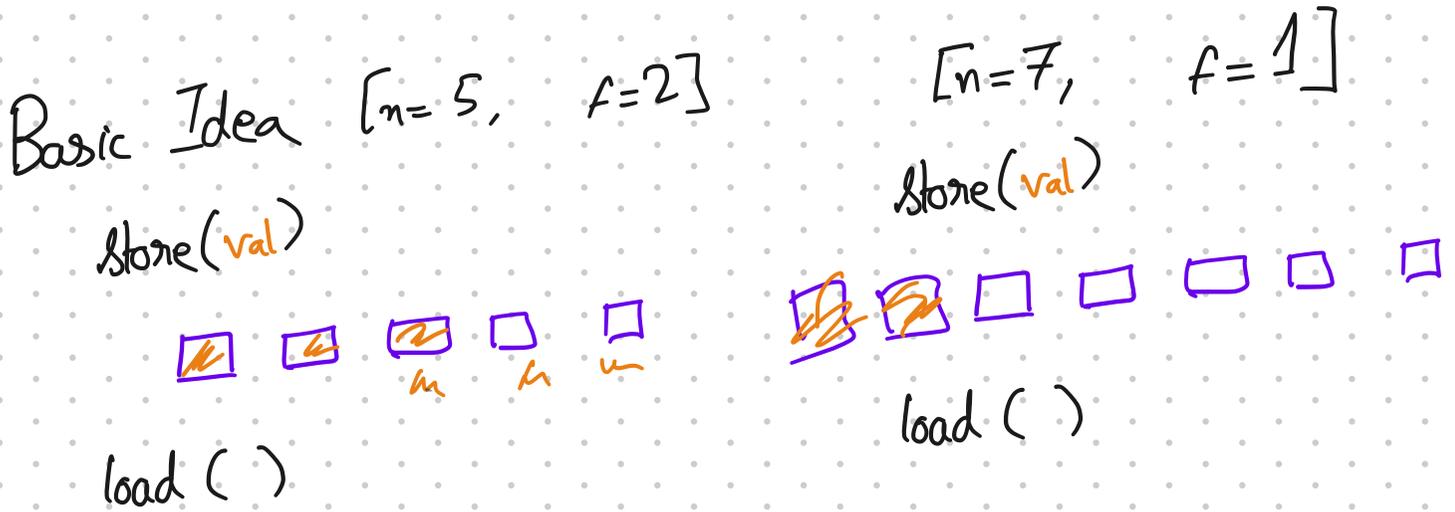
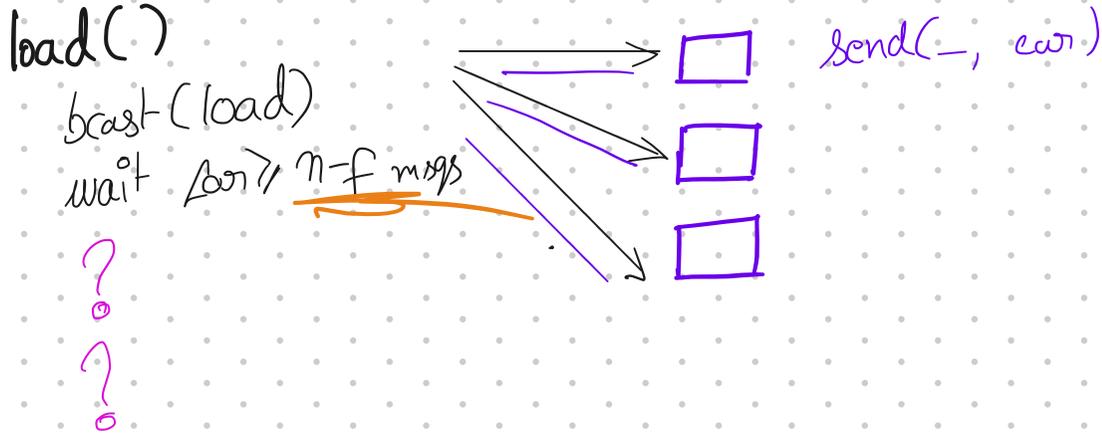


if (cur = L)
cur = val
send(-, done)
else
send(-, no)
(cur, val)

else

REVISIT N. WEEK





Important: Goal is ensuring quorum intersection

$n = 2f + 1$; etc. are all just ways to achieve this.

How does this all relate to Raft?

Normal operation } Append Entry to add command:
 - Leader working as expected, etc. } \rightarrow Leader: store(log idx, cmd)

Leader Election? (Request-Vote to become leader)

Leader election
 ↳ Candidate: store (term, node ID)
 ↳ Leader is value in last stored term

Back to Raft as Raft

TERM, COMMAND

Index

	0	1	2	3	4
1	TERM CMD				
2					
3					
4					
5					

Logs

State @ each node

- Log
- currentTerm
- votedFor
- type
 - Follower
 - Candidate
 - Leader
- commit index
- last applied

(Figure 2 in the paper is your friend)

Invariants

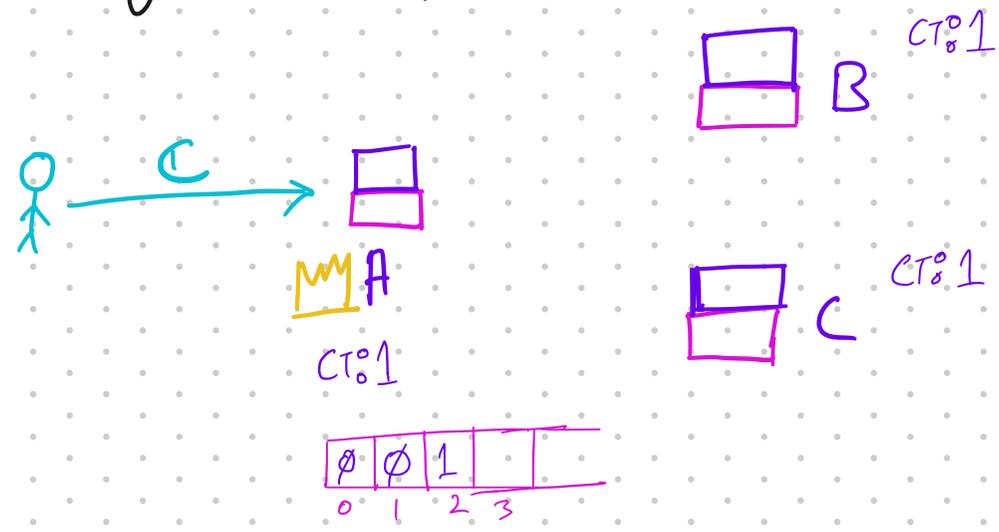
- currentTerm is increasing on each node.
- For any term at most 1 process/node is LEADER.
- Leader (for the current term*) has THE authoritative log.
- term, log index maps to a unique command

0	0	1	
---	---	---	--

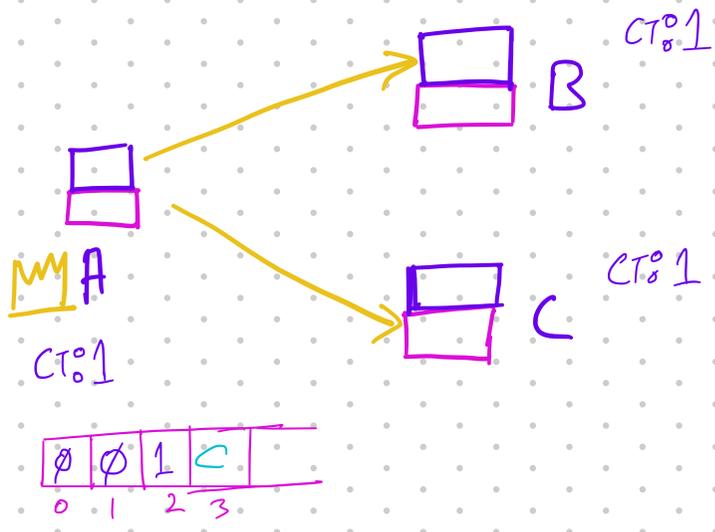
0	0	0	
---	---	---	--

- COMMIT POINT:

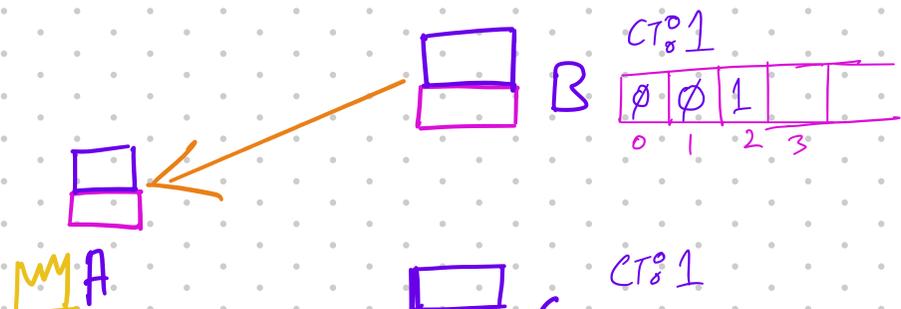
Walk through the basic protocol — see how invariants are preserved.

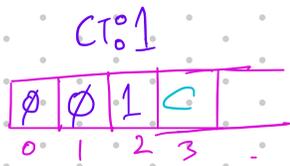


Append Entry (index=3, c, term=1, lastLogIndex=2, lastLogTerm=1)

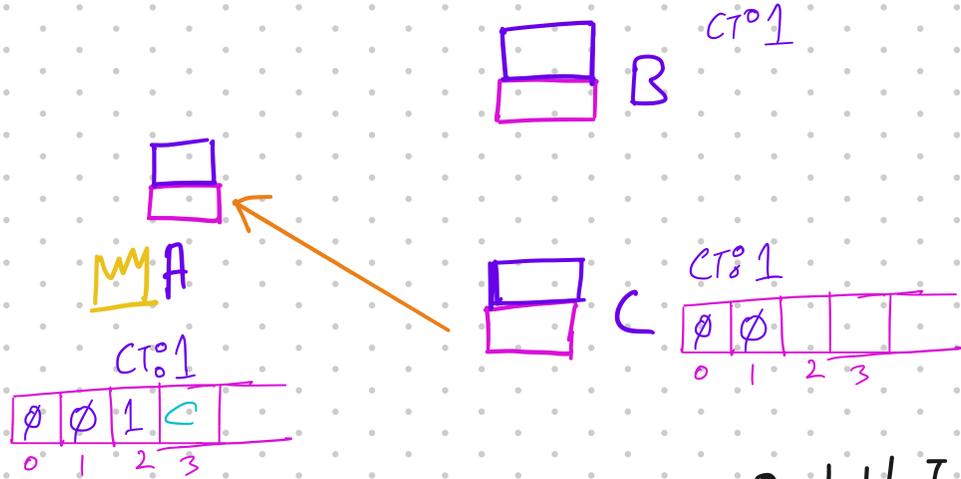


Append Entry (index=3, c, term=1, lastLogIndex=2, lastLogTerm=1)

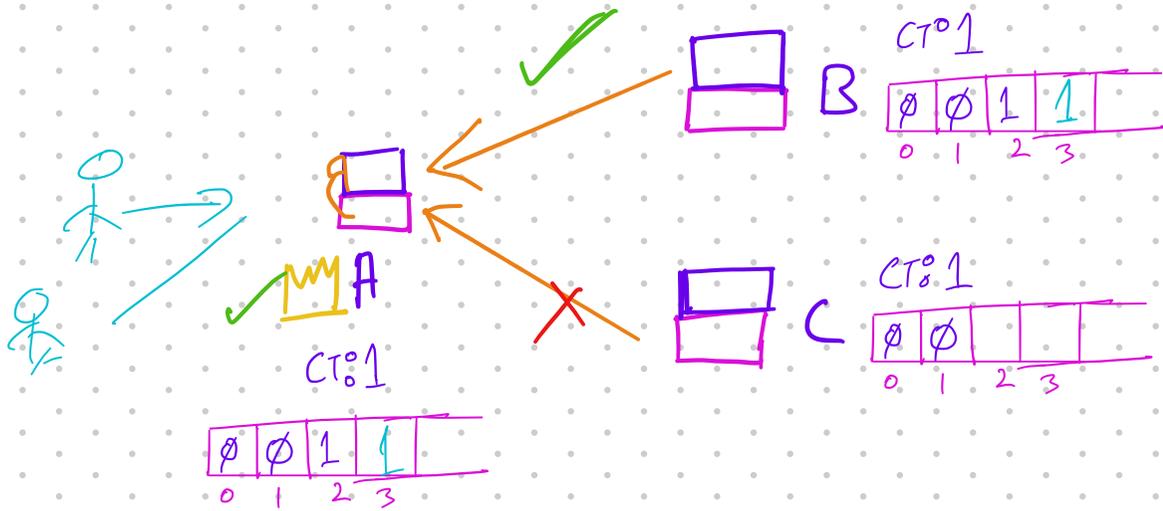




Append Entry (index=3, C, term=1, lastLogIndex=2, lastLogTerm=1)

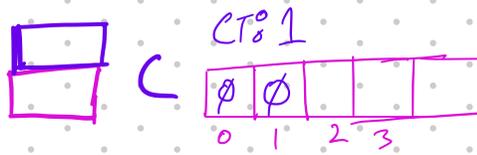
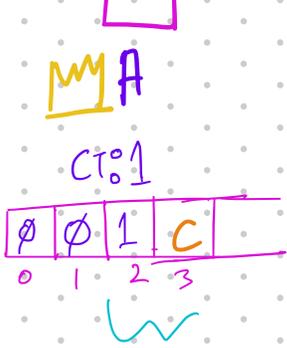


Append Entry (index=3, C, term=1, lastLogIndex=2, lastLogTerm=1)

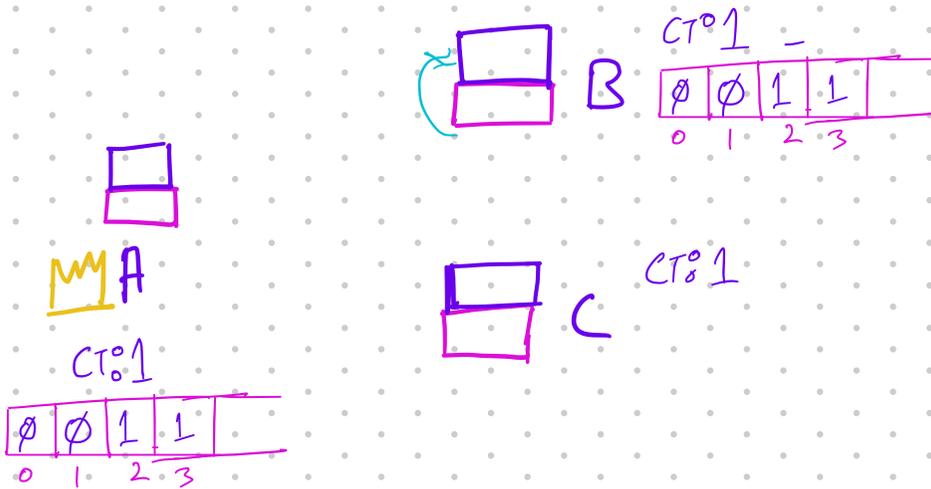


What happens with C?

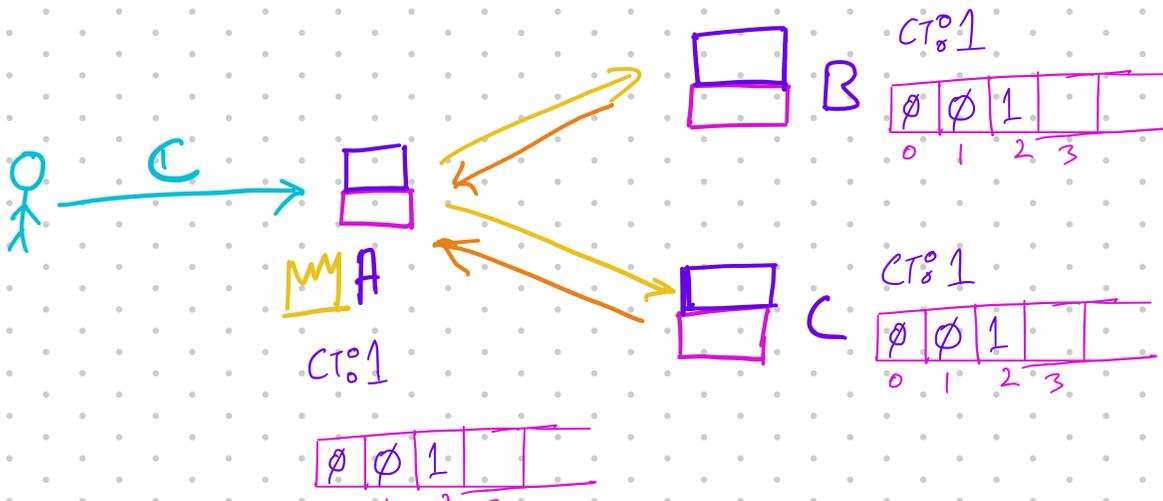


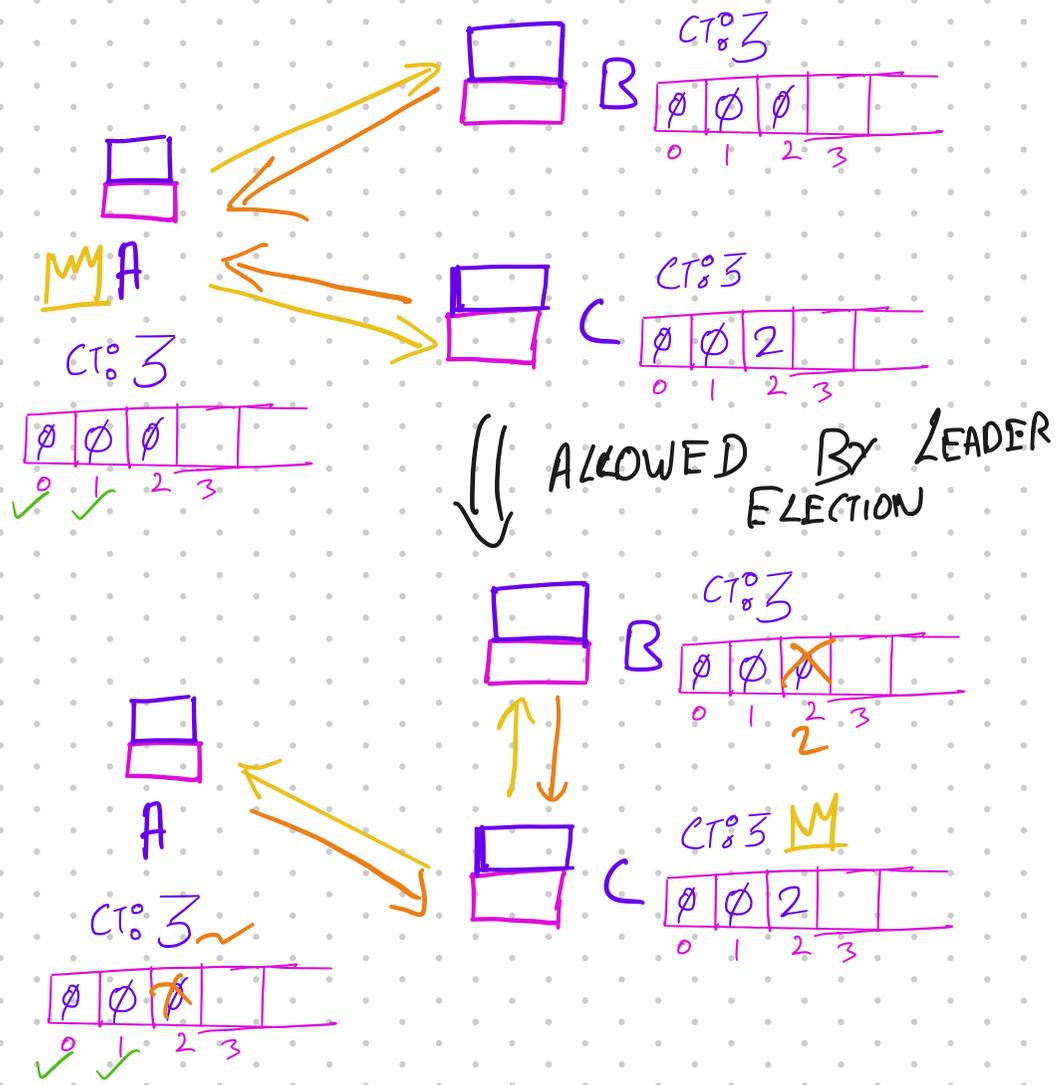


When does B deliver the command?



COMMIT POINTS





- Leader (for the current term*) has THE authoritative log.

WHEN IS AN ENTRY COMMITTED?

LEADER ELECTION

Invariant - For any term at most 1 process/node is LEADER.

① Trigger election.

Desirable Avoid leader election when possible (until leader fails)

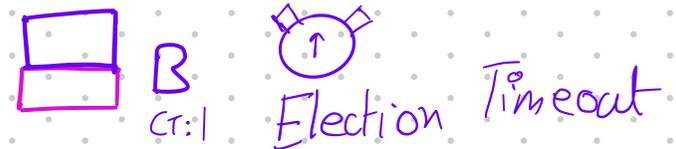
} Async model
⇒ Impossible

Why? Better performance

Instead tradeoff

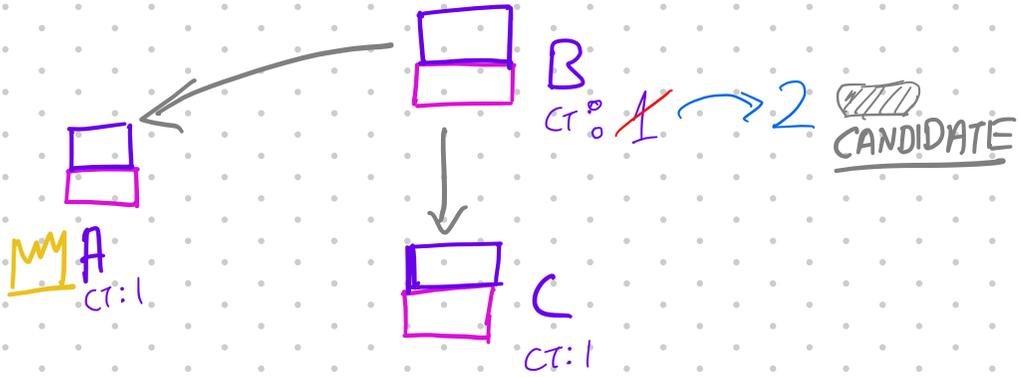
Freq. of unnecessary election

Time to recover from leader failure.



HB timeout < Election timeout?

② When election timeout expires at node B



Request Vote (ID=B, term=2)

INV - Leader (for the current term*) has THE authoritative log.

⇒ Any node elected in term 2 must have all committed entries.

Two paths

→ Only elect leaders with all committed entries [Raft]

→ Fix up log after election [Next week, etc.]

PROBLEM: ONLY THE LEADER KNOWS WHAT ENTRIES ARE COMMITTED

SOLUTION: QUORUM INTERSECTION

↳ Any committed entry was replicated to a quorum.
[$\frac{n}{2} + 1$ nodes]

⇒ Known to at least one process in any other quorum.

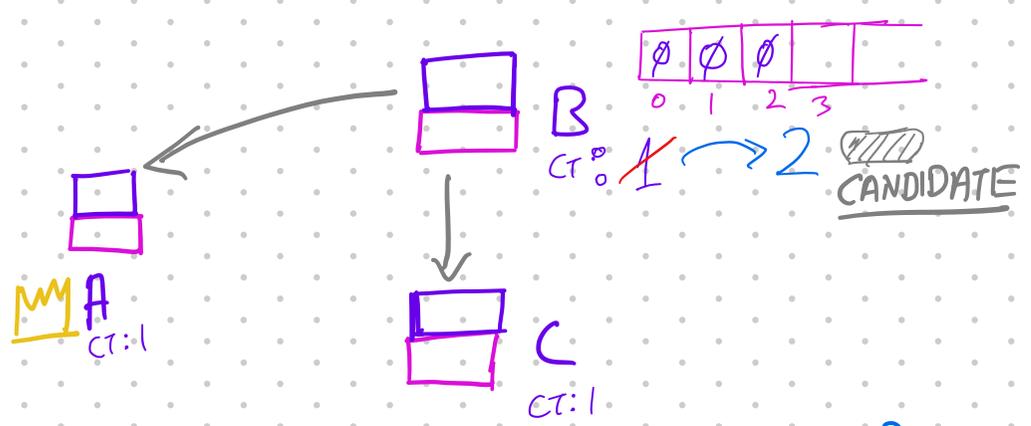
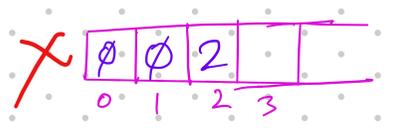
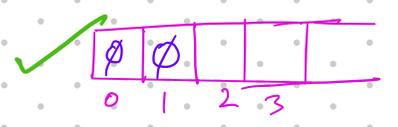
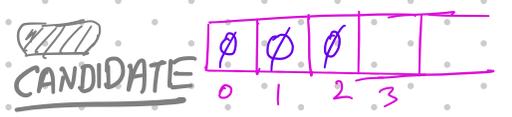
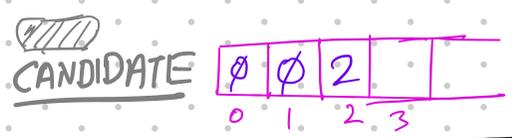
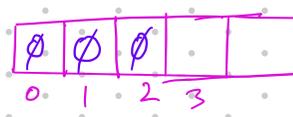
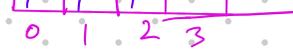
- Candidate needs vote from quorum to become leader.

How recently was log updated - Nodes vote for a candidate only if

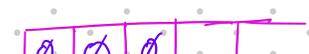
OR

- (a) Candidate has more entries in its log & last log entry term ^{higher}
- (b) Candidate has the same number of entries but candidate's last entry term \geq node's last entry term.

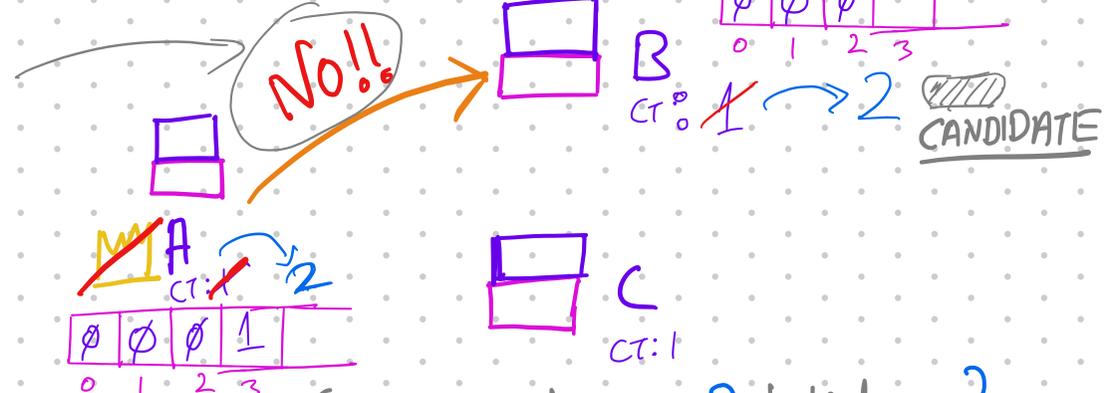




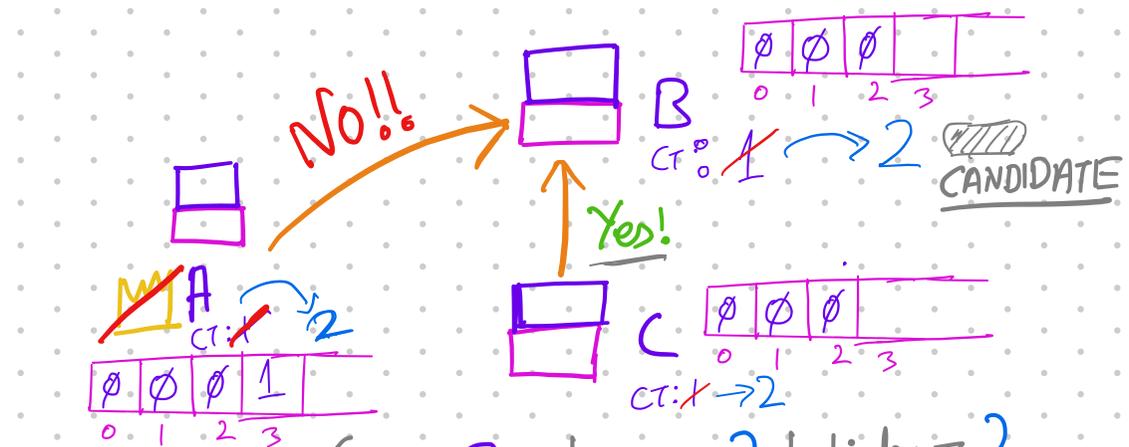
Request Vote (ID=B, term=2, last index = 2, last IndexTerm = ∅)



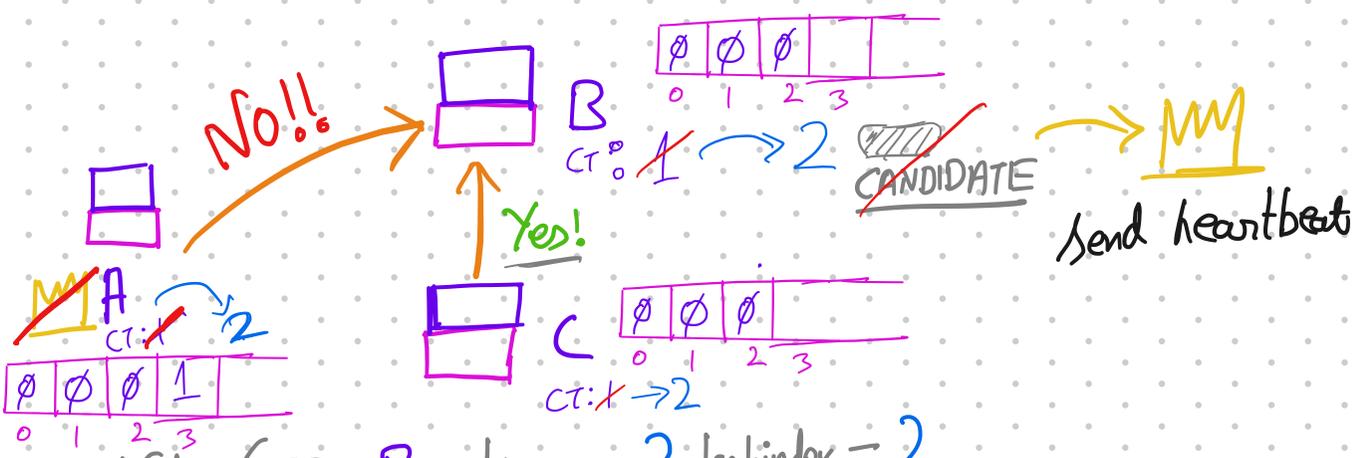
Why?



Request Vote ($ID=B$, $term=2$, $lastindex=2$,
 $lastIndexTerm = \emptyset$)

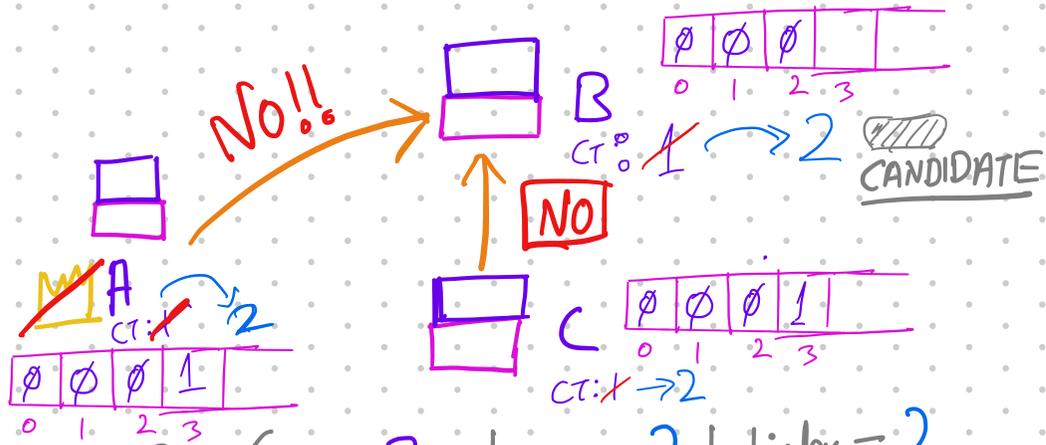


Request Vote ($ID=B$, $term=2$, $lastindex=2$,
 $lastIndexTerm = \emptyset$)



Request Vote ($ID=B$, $term=2$, $lastindex=2$,
 $lastIndexTerm = \emptyset$)

ALT HISTORY

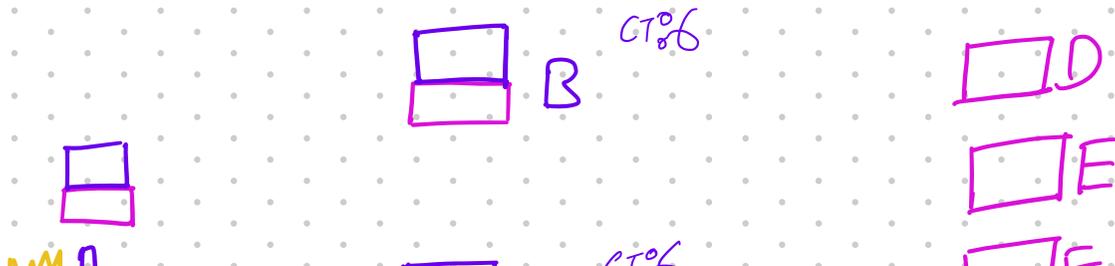


Request Vote (ID=B, term=2, lastindex=2, lastIndexTerm= \emptyset)

What now???

Joint quorums & reconfig.

PROBLEM



W A
CT:6

C C
CT:6

F
G

$N/2 + 1 = 4$ ∴ D, E, F & G form a quorum

But - do not have current log
(on any other state)

How to have them safely join

OLD + NEW ∴ Operations on two Raft instances!

