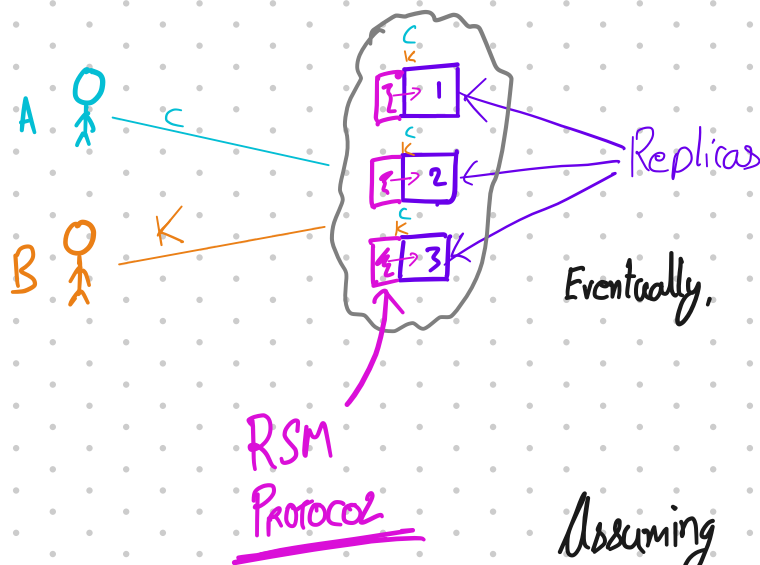


R
 STATE
 MACHINES } RAFT



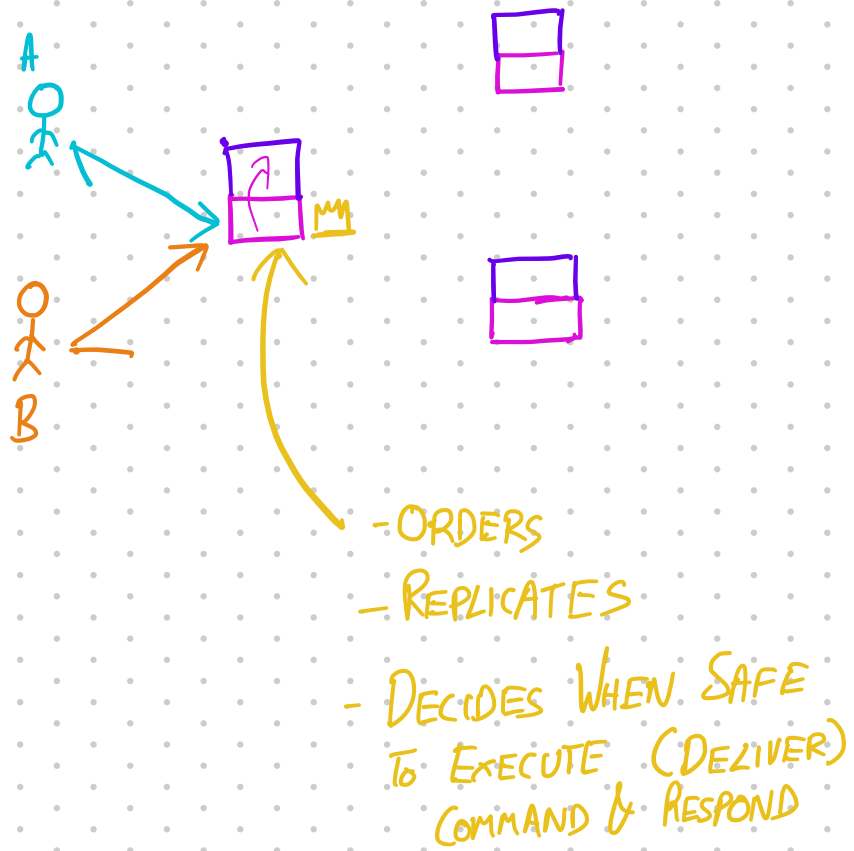
LAST CLASS:

- RSMs: A recipe for building fault tolerant systems



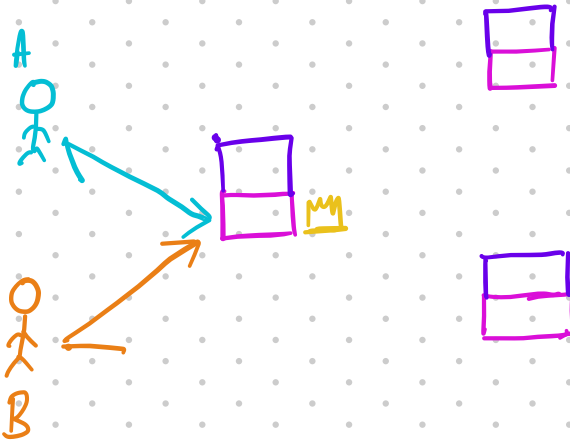
- Correctness requirements
 - Agreement
 - Ordering

- Leader based RSM



CORE CONCERNS

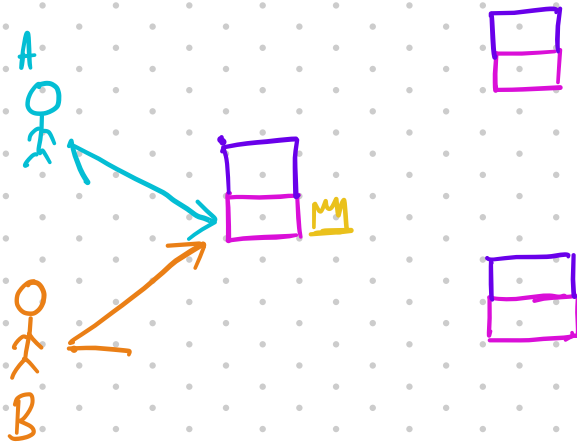
① WHEN IS IT SAFE TO EXECUTE/DELIVER?



REQUIREMENTS :

② ... & RESPOND TO LEADER

(2) How to DETECT & RESPOND TO FAILURE.



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



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: Term.
> f fail:



Consider n=3 f=1

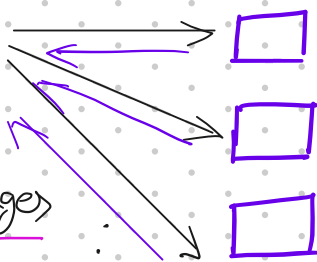
store(val)

broadcast (store, val)
wait for ≥ f+1 messages

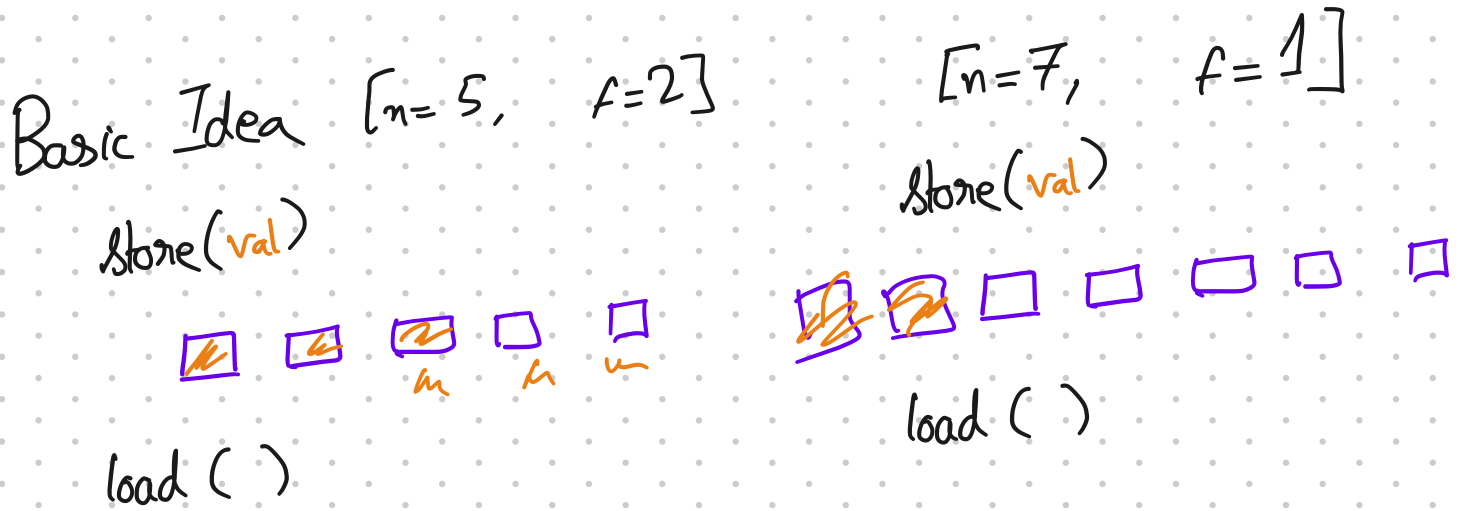
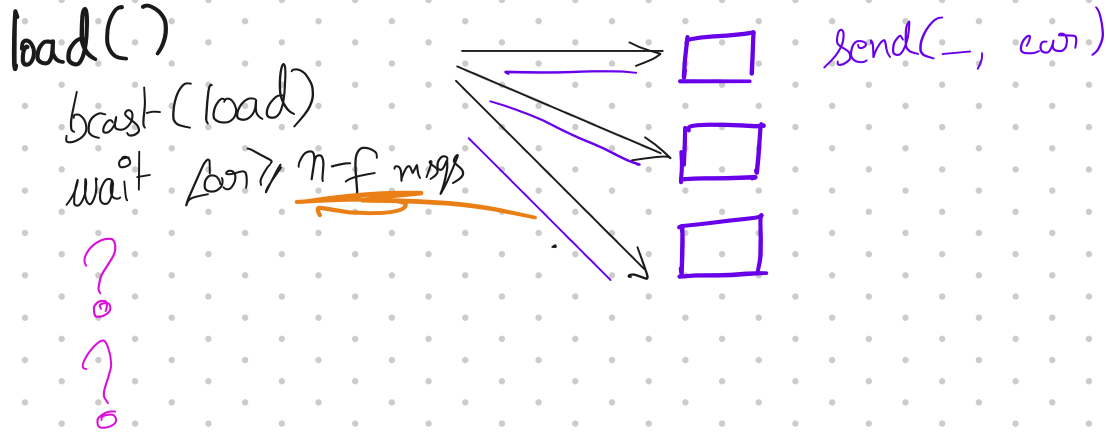
if |done| ≥ f+1
return true

else

REVISIT N. WEEK



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



Observe: $f < \frac{n}{2}$ & n -odd $\Rightarrow f+1 = n-f$

$$f = \frac{n-1}{2}$$

$$f+1 = \frac{n+1}{2} \quad \left| \quad n-f = \frac{n+1}{2} \right.$$

How does this all relate to Raft?

Normal operation } Append Entry to add command:
 - Leader working as expected, etc. } \hookrightarrow 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

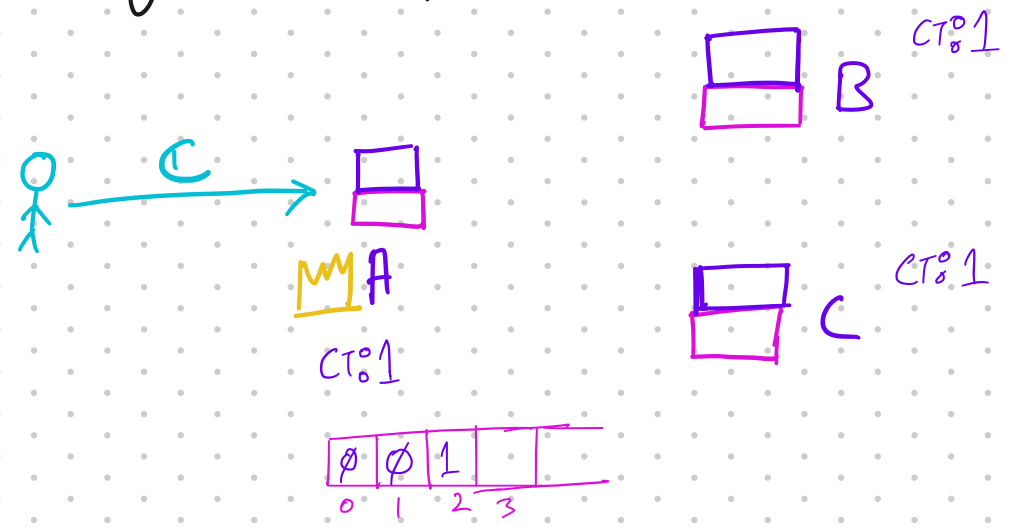
- currentTerm is increasing on each node.
- For any term at most 1 process/node is LEADER.
- Leader (for the current term*) has THE authoritative log.
- term, log index maps to a unique command

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

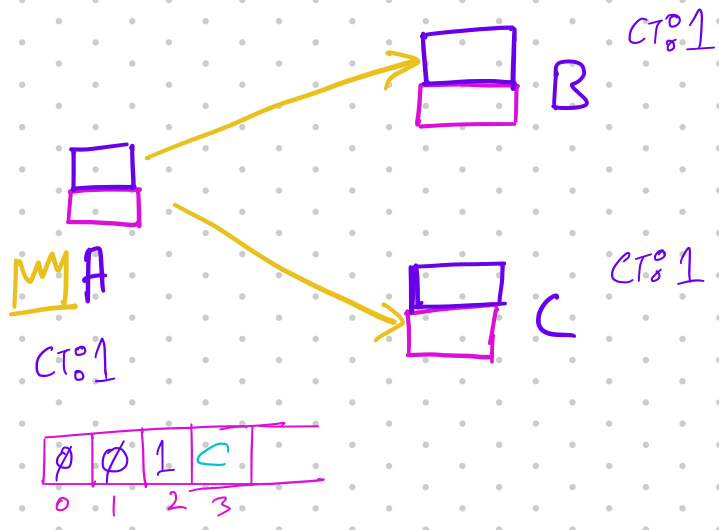
0	0	0	
---	---	---	--

- COMMIT POINT:

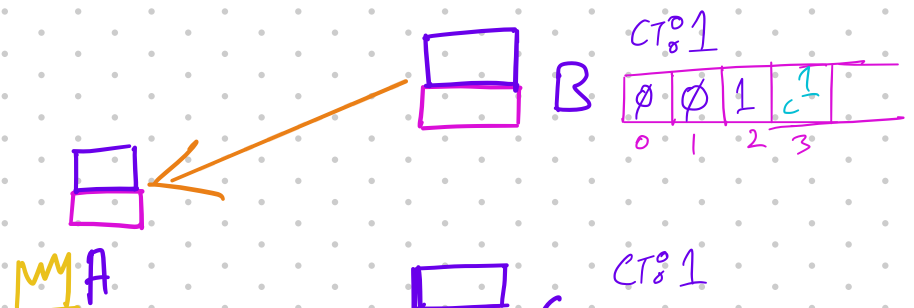
Walk through the basic protocol — see how invariants are preserved.



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



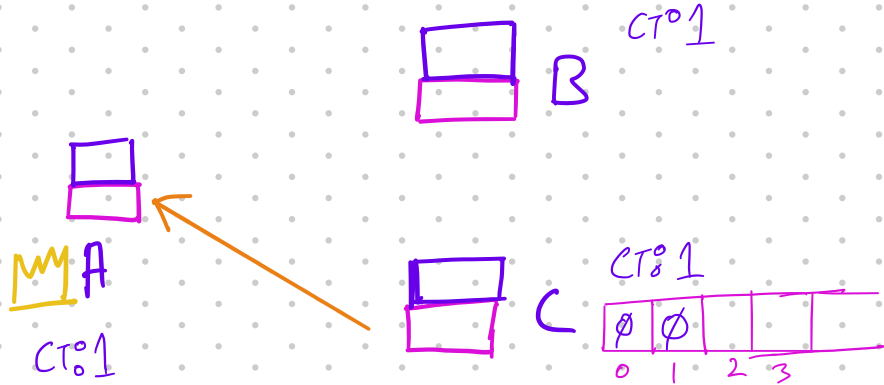
Append Entry (index=3, c, term=1, lastLogIndex=2, lastLogTerm=1)



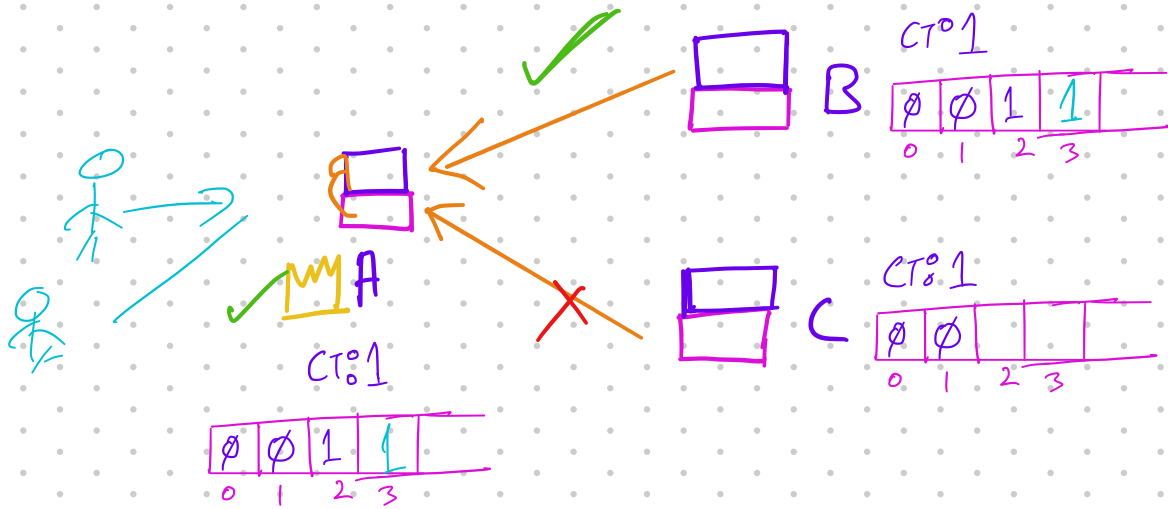
CT:1



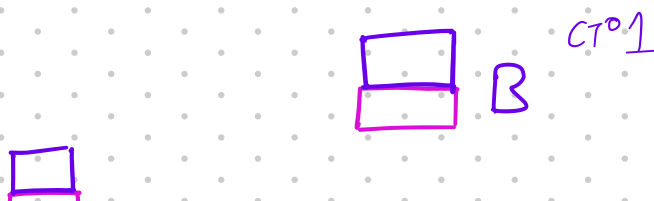
Append Entry (index=3, c, term=1, last Log Index=2, last Log Term=1)

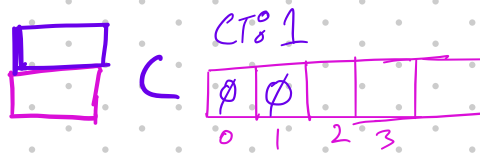
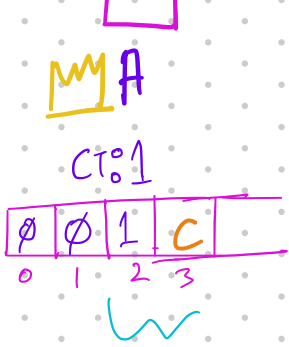


Append Entry (index=3, c, term=1, last Log Index=2, last Log Term=1)

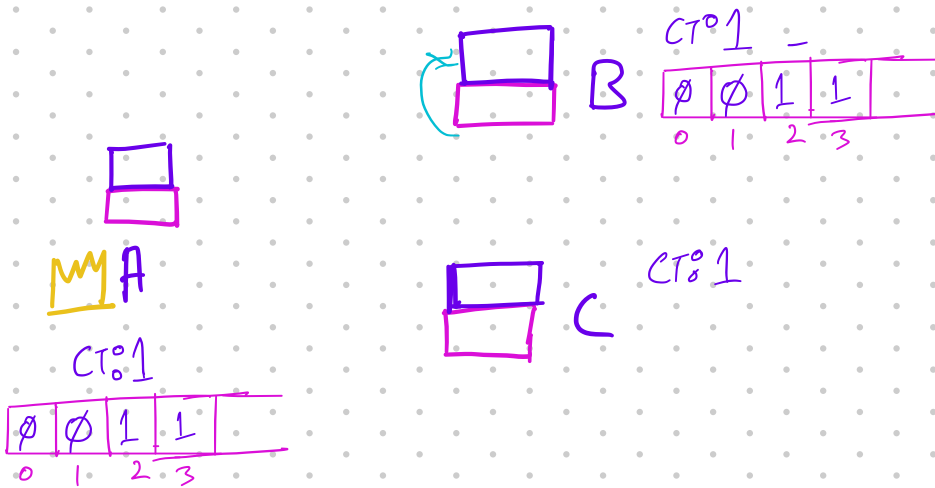


What happens with C?

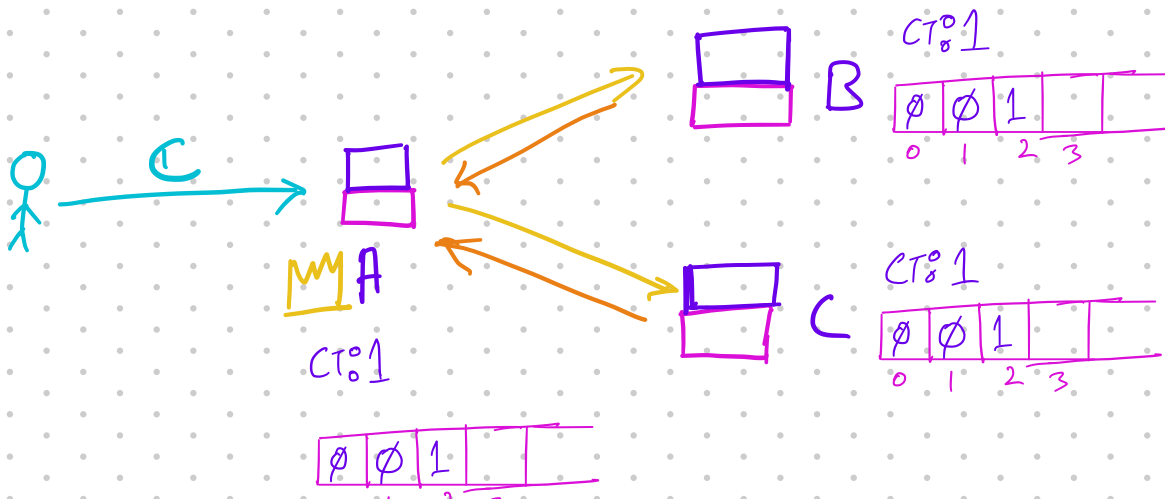


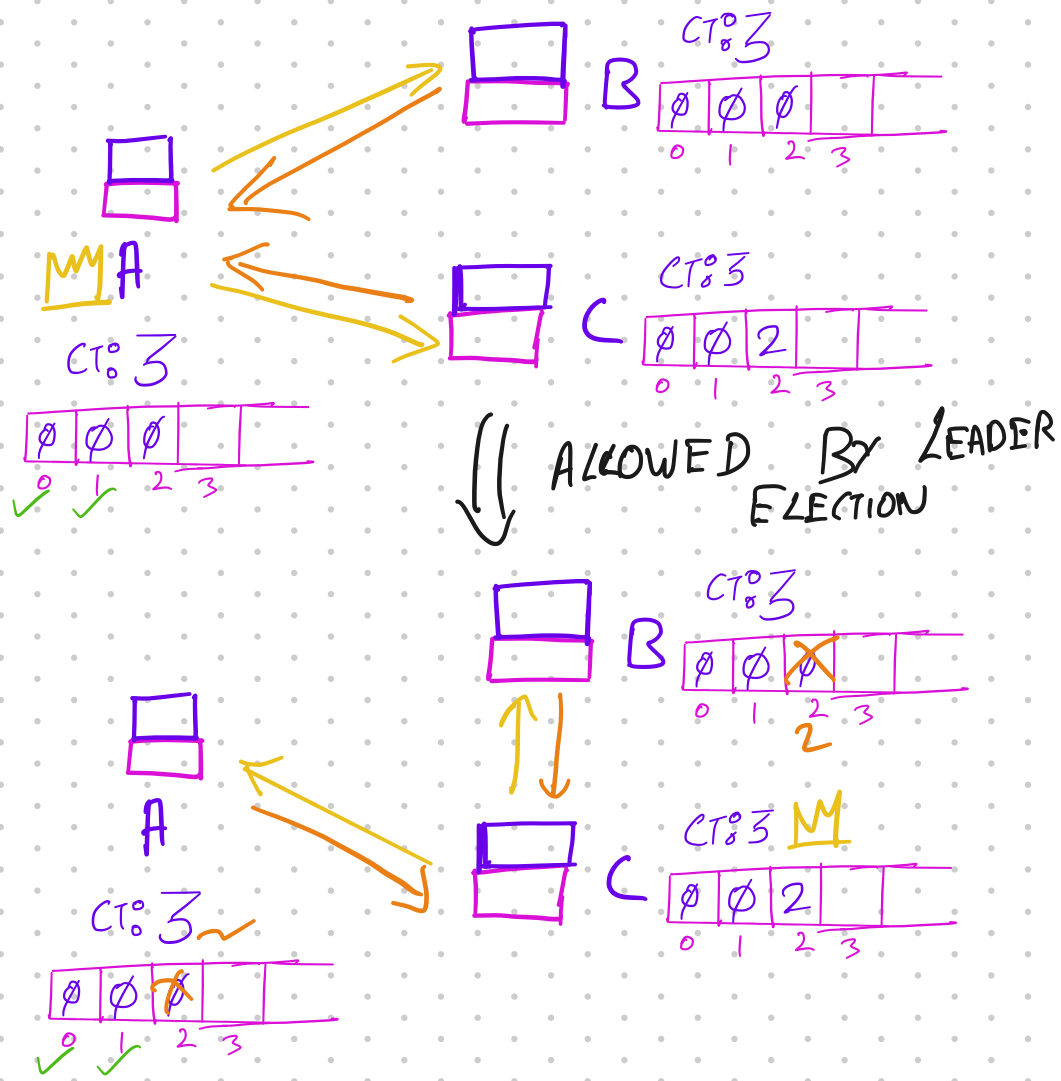


When does B deliver the command?



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.

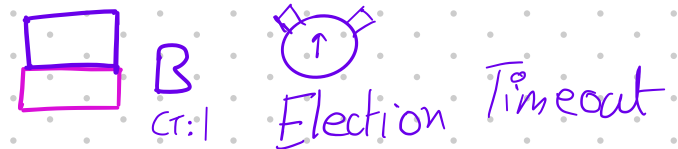
Desirable Avoid leader election when possible (until leader fails) } Async model
⇒ Impossible

Why? Better performance

Instead tradeoff

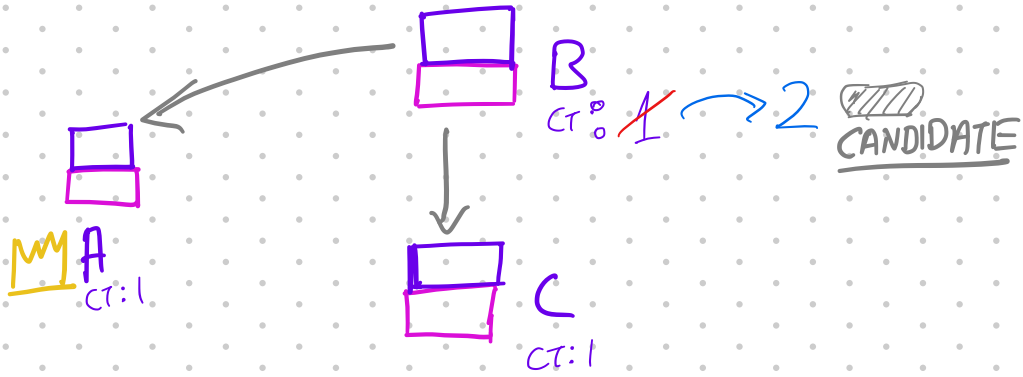
Freq. of unnecessary election

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.

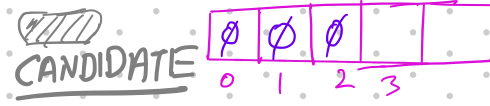
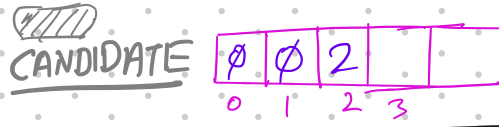
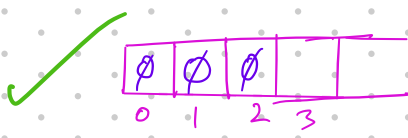
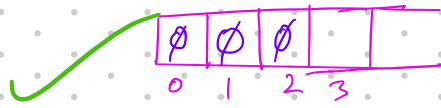
- Candidate needs vote from quorum to become leader.

How recently was log updated -

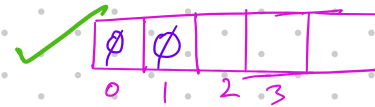
Nodes vote for a candidate only if

OR

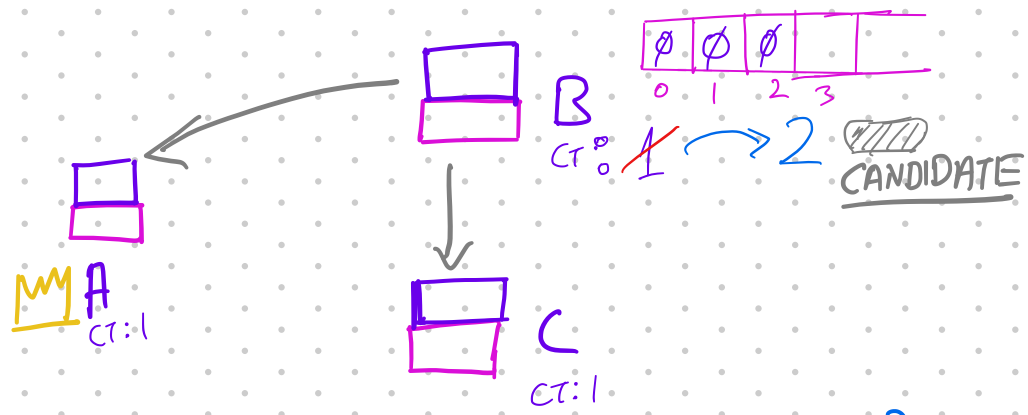
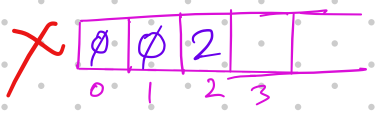
- (a) Candidate has more entries in its log
- (b) Candidate has the same number of entries but candidate's last entry term \geq node's last entry term.



0 0 0 0 0
 0 0 0



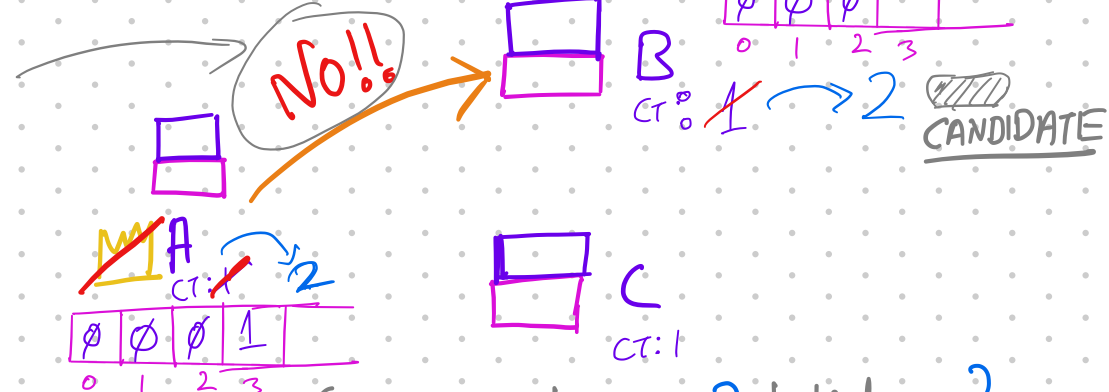
0 2



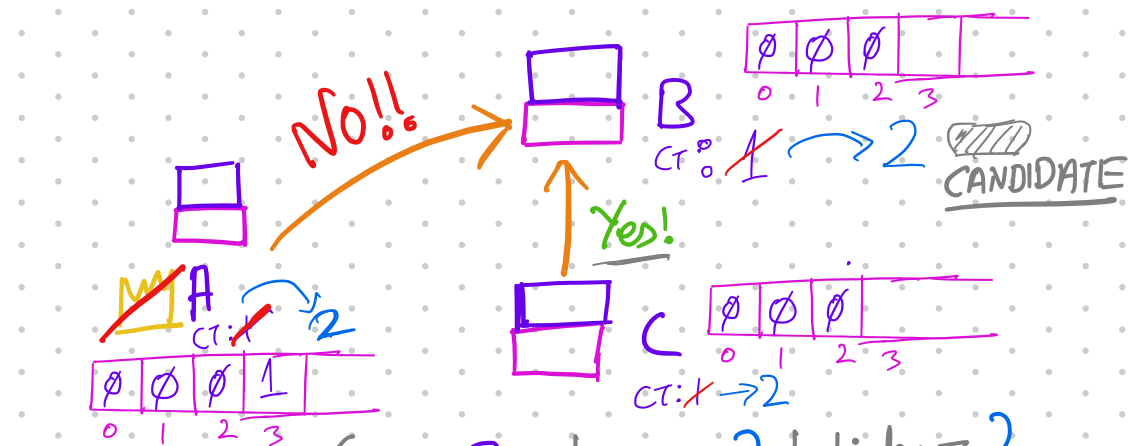
Request Vote (ID=B, term=2, last index=2,
 last IndexTerm = ∅)



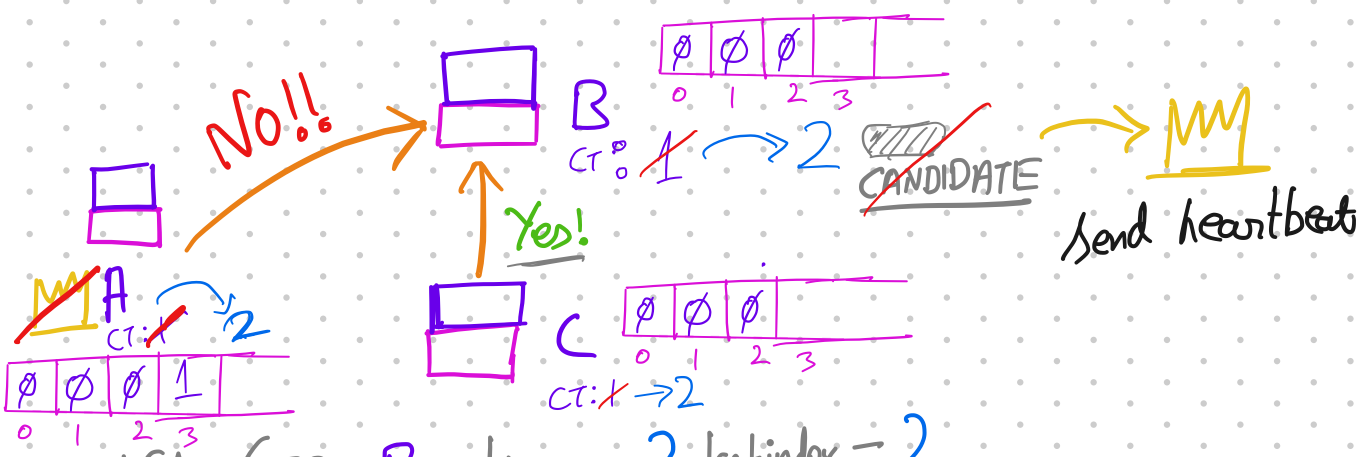
Why?



Request Vote (ID=B, term=2, lastIndex=2, lastIndexTerm=∅)

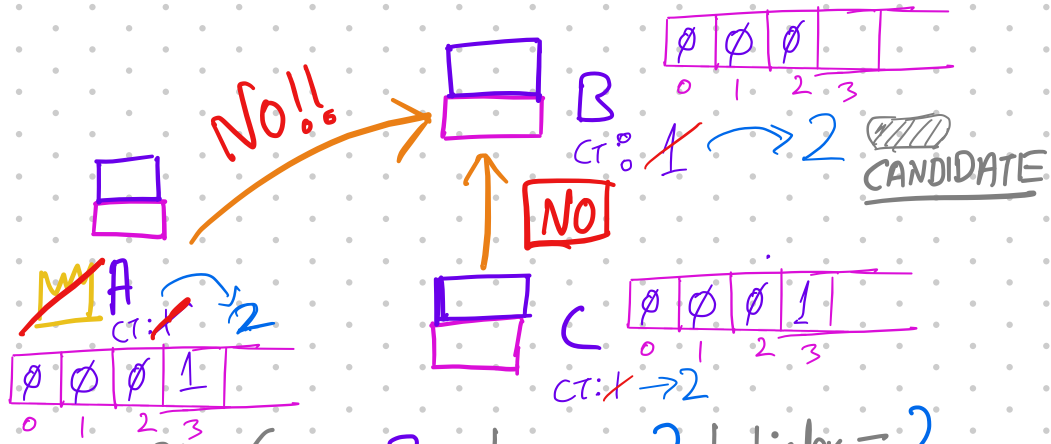


Request Vote (ID=B, term=2, lastIndex=2, lastIndexTerm=∅)



Request Vote (ID=B, term=2, lastIndex=2, lastIndexTerm=∅)

ALT HISTORY



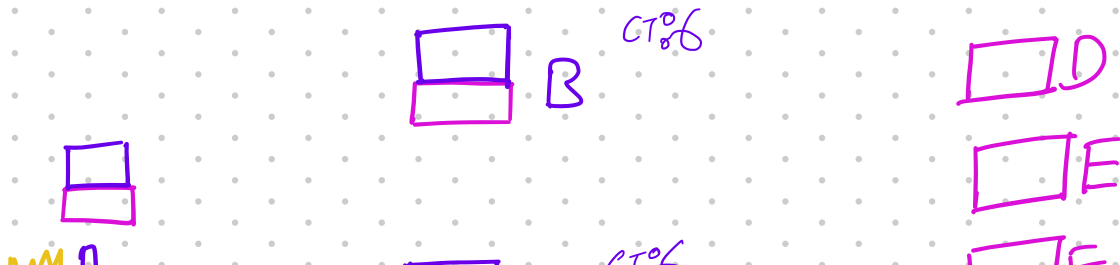
Request Vote (ID=B, term=2, last index=2, lastIndexTerm= \emptyset)

What now???

000

Joint quorums & reconfig.

PROBLEM



W A
CT:6

W C
CT:6

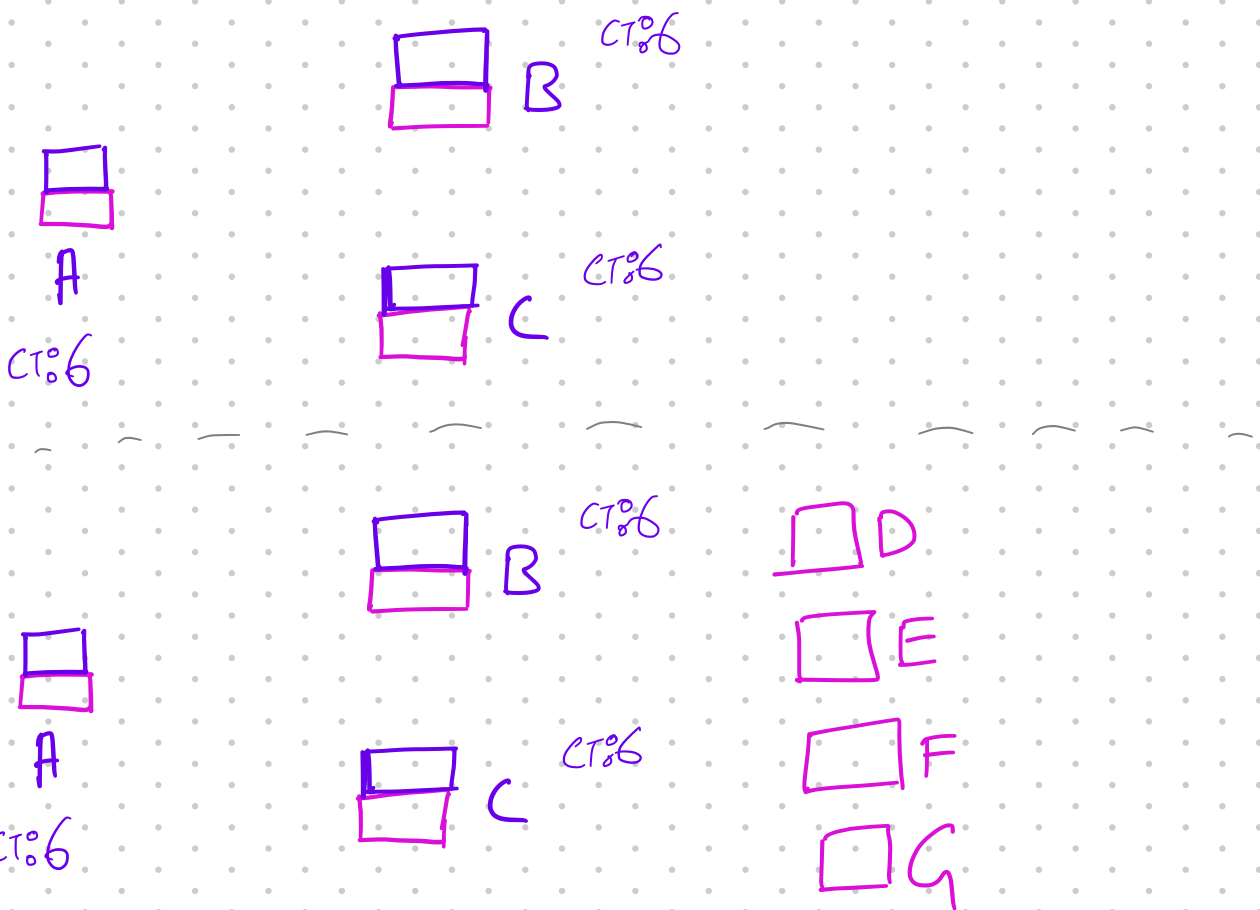
W F
W G

$N/2 + 1 = 4$ ∴ 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!



00
00

00
00

00
M002

00 (2) (3)

AER (rd=4, [],
ll=3, llt=2)