

MORE

Rsms

Where we are

- STATE MACHINES & STATE MACHINE REPLICATION

- RAFT

↳ OVERALL STRUCTURE

* LEADER SYNCHRONIZES LOG WITH REPLICAS
DECIDES WHEN COMMANDS ARE

COMMITTED

* AN ENTRY AT INDEX I ONCE COMMITTED

IS AT LOG INDEX I AT ALL FUTURE
LEADERS

* LOG COMPLETENESS.

⇒ PROTOCOL

✓ (i) COMMAND REPLICATION

✓ (ii) FAILURE DETECTION

... (iii) LEADER ELECTION / FAILURE HANDLING

(iv) RECONFIGURATION

LEADER ELECTION

- TERM

- REQUIREMENTS TO VOTE FOR A CANDIDATE

- REQUIREMENT FOR A CANDIDATE → LEADER

WHEN IS AN ENTRY COMMITTED?

- COMMAND REPLICATION?

INTERACTIONS WITH LEADER ELECTION

	0	1	2	3	4	5	6	7	8
$\text{M}_{2,4}$ A log	1	2							
B log	1	2							
C log	1	2							
D log	1	2							
M_3 E log	1	3							

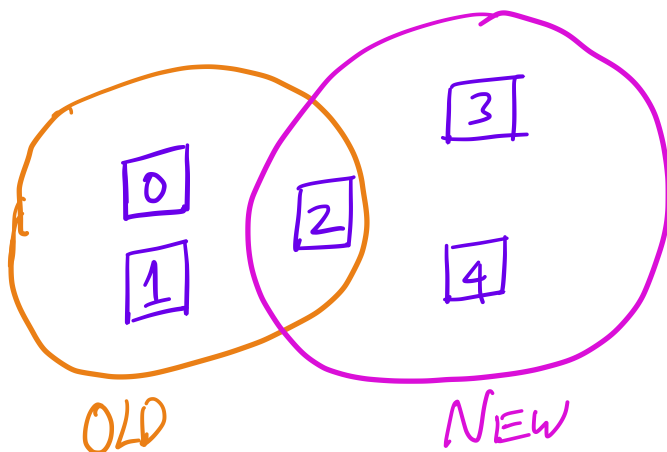
5 (LastLogTerm > NODE.LAST LOG TERM) OR

(LastLogTerm == NODE.LAST LOG TERM

AND LastLog Index > Node.Last Log INDEX)

RECONFIGURATION

PROBLEM: WANT TO ADD &/OR REMOVE NODES



Need to make sure no commands are lost

Avoid "SPLIT BRAIN"

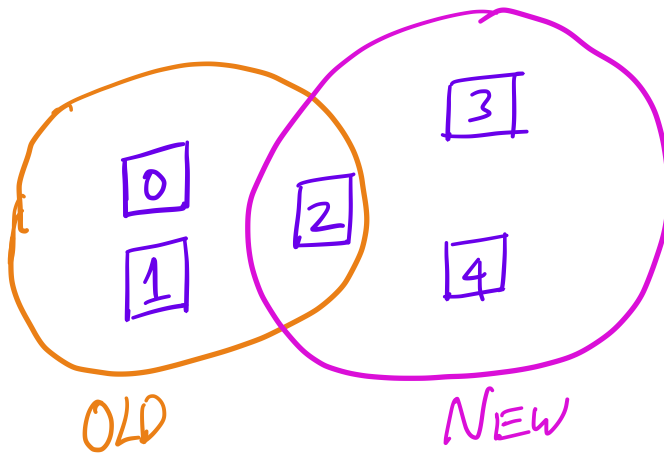
How:

- Store configuration in log
 ↑
 Set of active nodes

- Transition



Majority in OLD + Majority in New



MULTI PAXOS

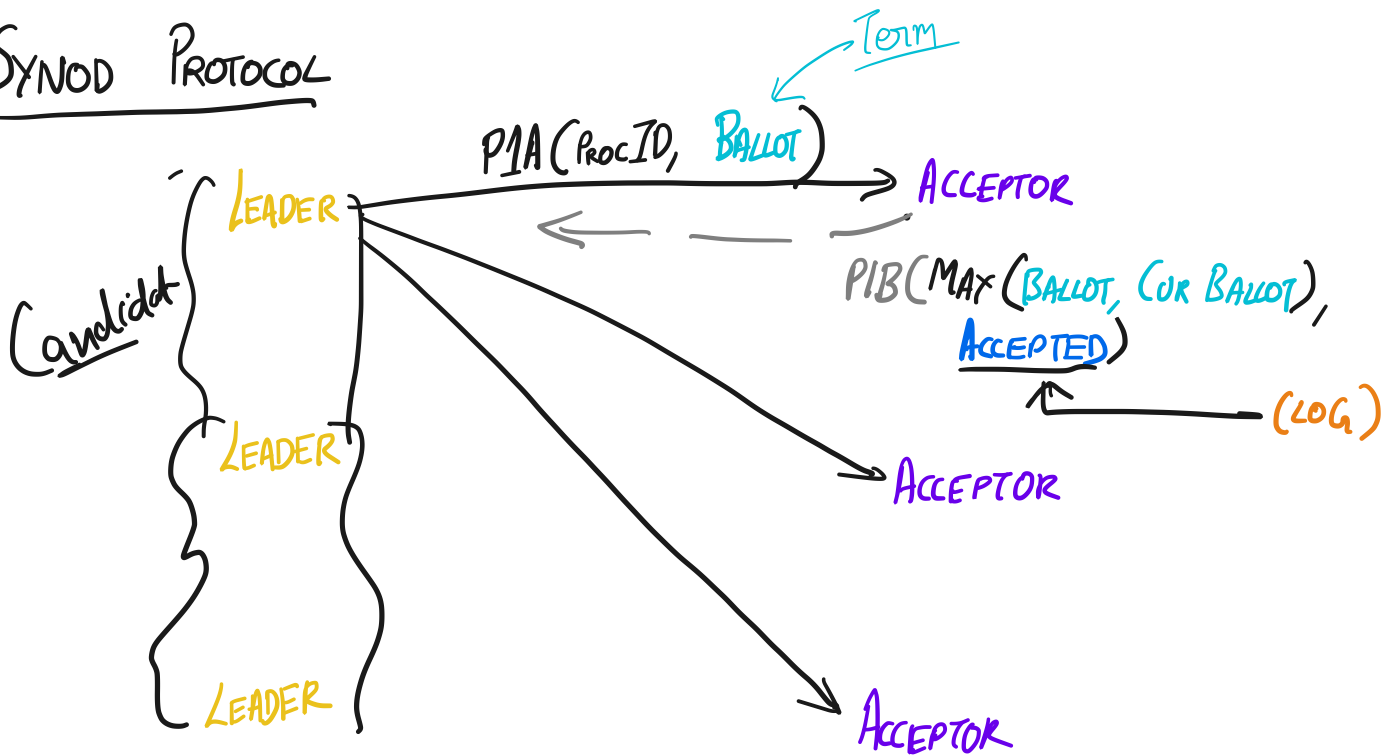
- CLIENTS / REPLICAS / LEADERS / ACCEPTORS

SIGH!

- REALLY NOT ALL THAT DIFFERENT FROM RAFT AS WE WILL SEE

REPLICA	LEADER	ACCEPTOR	
REPLICA	LEADER	ACCEPTOR	$f = 1$
			[NEED
REPLICA	LEADER	ACCEPTOR	$\geq 2f + 1$ acceptors
			$\geq f + 1$ replicas, leaders
]

SYNOD PROTOCOL



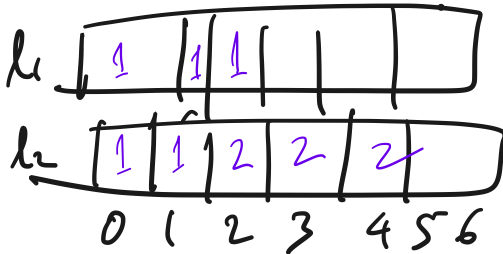
$\langle P1A, P1B \rangle \equiv$ LEADER ELECTION / VIEW CHANGE

DIFFERENCES

o No log completeness check.

o IMPLICATION: MUST ENSURE CHOSEN
LEADER HAS COMMITTED ENTRIES.

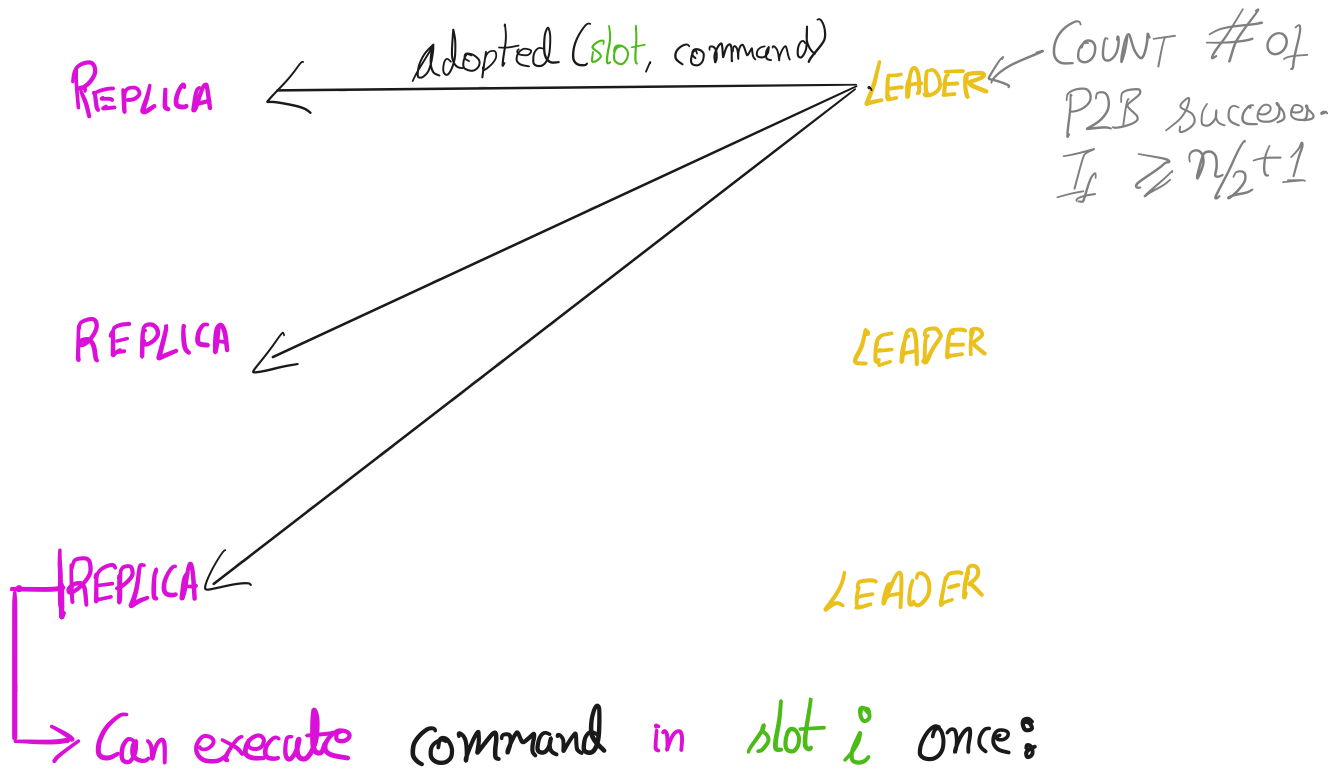
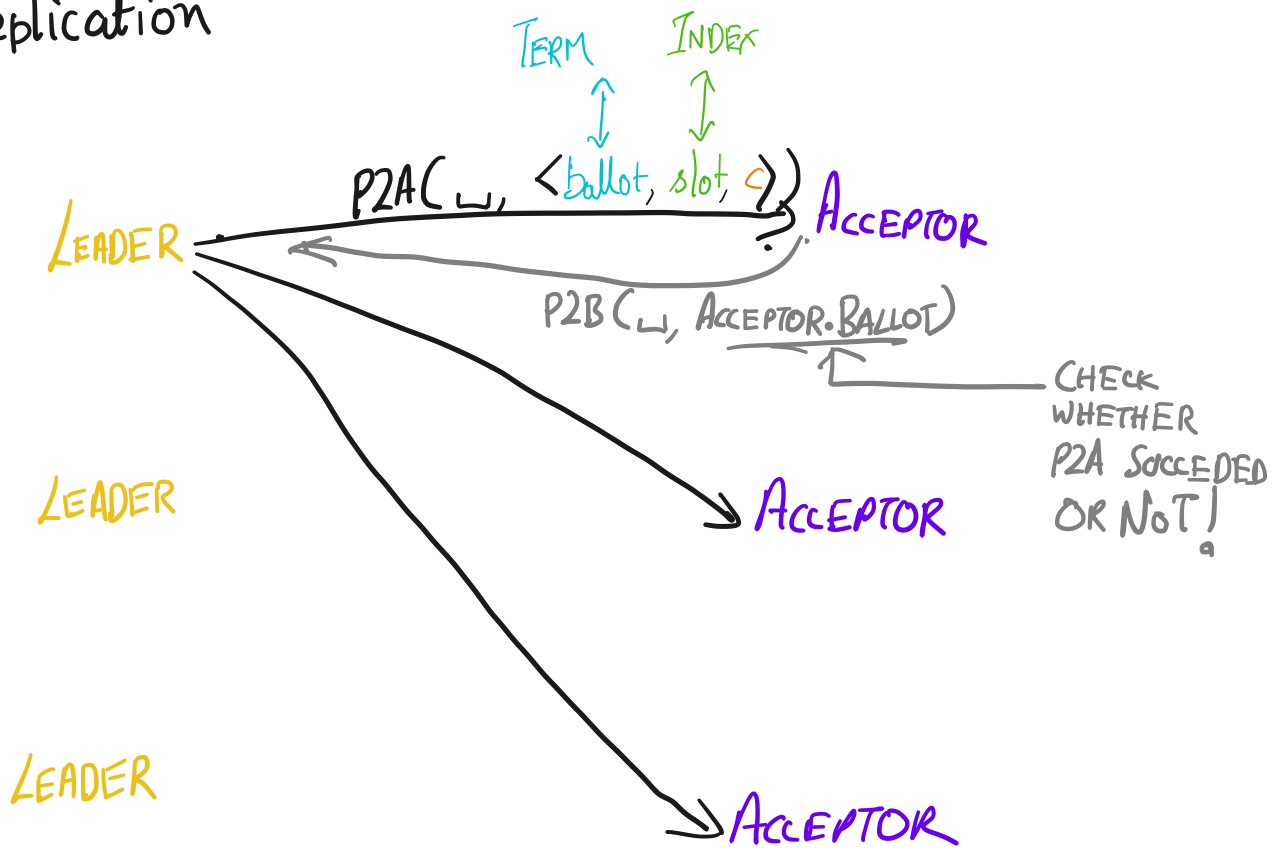
How?



o WHEN IS P1A TRIGGERED &
LEADER STABILITY?

SCOUTS

Phase 2: Replication



- * Command in $i-1$ adopted & executed &
- * Command in slot i ADOPTED

ALTERNATE VIEW

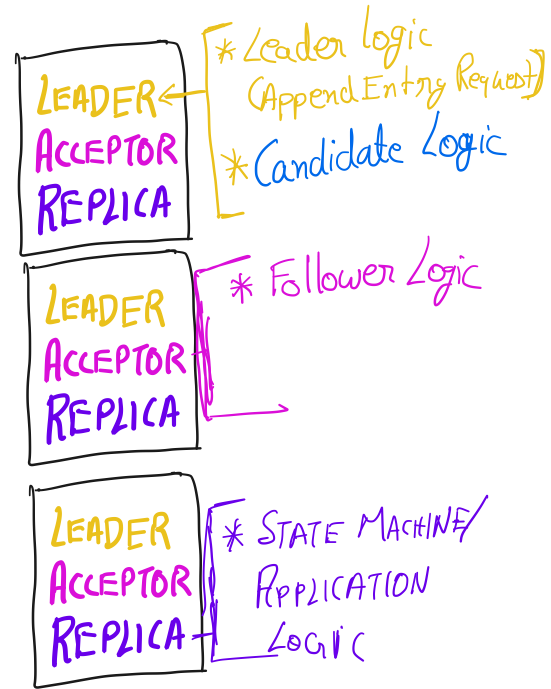
M-PAXOS

REPLICA LEADER ACCEPTOR

REPLICA LEADER ACCEPTOR

REPLICA LEADER ACCEPTOR

RAFT



SOME THINGS LEFT UNDERSPECIFIED

* FAILURE DETECTION LOGIC

MY CONTENTION: RAFT IS DESIGNED FOR A PARTICULAR USE

MULTI PAXOS PROVIDES A DESIGN PATTERN THAT CAN BE APPLIED TO DIFFERENT DEPLOYMENTS

EXAMPLES (MORE IN THE PAPER)

DISK PAXOS: ACCEPTORS ARE DISKS

* WHY?

* WHAT CHANGES?

ACCEPTOR INTERFACE: READ BLOCK, WRITE BLOCK

NOT ACTIVE

MENCLUS/EPAXOS: CONSENSUS IN THE WIDE AREA

CONSENSUS IN PRACTICE: LOCK SERVICES.