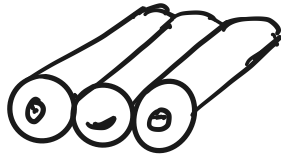


RAFT: An RSM Protocol



FINAL PROJECT PROPOSALS

- GROUP SIZE
- TOPICS
- PROPOSAL LENGTH

MIDTERM OBSERVATIONS

- o SUCCINCTNESS:

DON'T TELL ME ABOUT SAFETY & LIVENESS
IN A QUESTION ABOUT FAIRNESS

- o ON PROOFS

↳ CAN ASSUME ANYTHING COVERED IN CLASS
DON'T NEED TO PROVE FLP

→ BUT MUST SHOW WHY
WHY IS WFD CONSENSUS?

→ MUST BE LOGICALLY CONSISTENT

WHERE WE LAST LEFT OFF

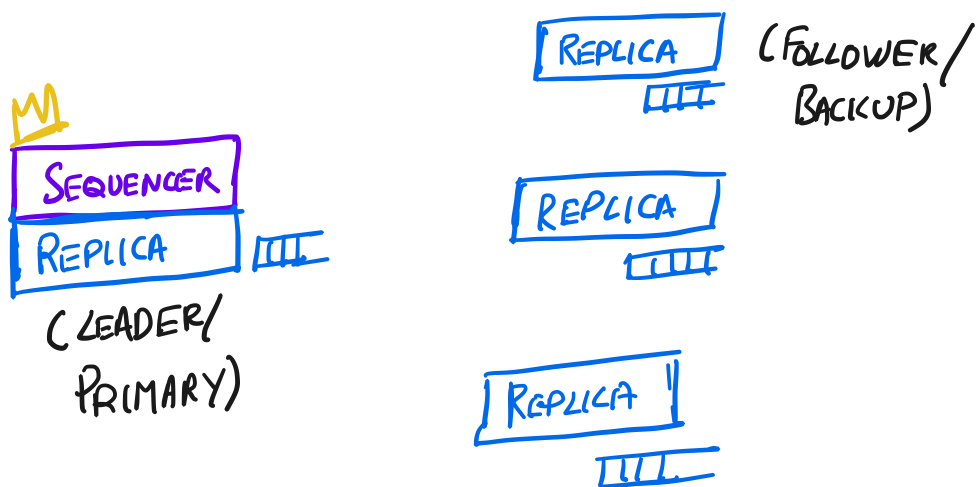
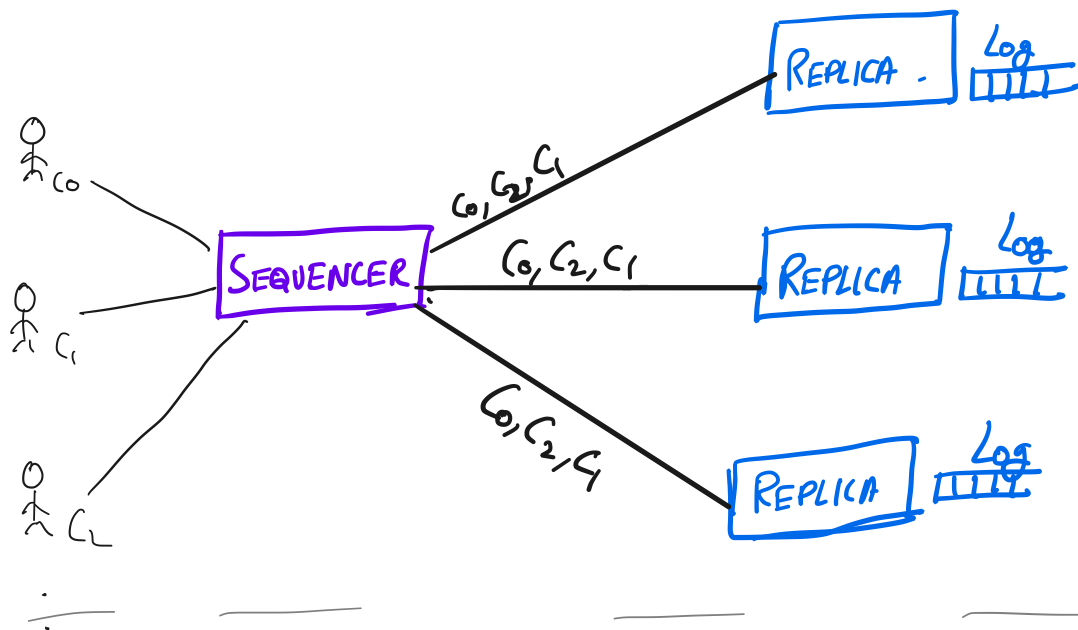
- o STATE MACHINES

- o REPLICATED STATE MACHINES

- o CONSENSUS

↳ IMPOSSIBILITY & WORKAROUNDS

TODAY: AN RSM PROTOCOL



INVARIANTS

- ① LEADER'S LOG IS 'AUTHORITATIVE'
 PROTOCOL MAKES SURE REPLICA LOGS == LEADER LOGS
- ② IF A MAJORITY OF REPLICAS HAVE LOG ENTRY E AT INDEX I , THEN ALL FUTURE LEADERS WILL HAVE LOG ENTRY $E @ I$

⇒ SAFE TO EXECUTE COMMAND AT INDEX I ONCE

(a) $I-1$ EXECUTED

(b) I REPLICATED BY MAJORITY

(i) LEADER ELECTION PROTOCOL ENSURES STABILITY OF A SUFFICIENTLY REPLICATED ENTRY

(ii) LEADER COUNTS # OF TIMES AN ENTRY IS REPLICATED & DECIDES WHEN COMMANDS ARE EXECUTED

(3) (RAFT SPECIFIC: LOG COMPLETENESS) IF INDEX I IS COMMITTED THEN INDEX $I-1$ MUST BE COMMITTED

SPLIT PROTOCOL INTO FOUR PORTIONS

(1) LOG SYNCHRONIZATION

HOW THE LEADER ADDS TO THE LOG

(2) FAILURE DETECTION

HOW FOLLOWERS DETECT LEADER FAILURES

(3) FAILURE RECOVERY

HOW A NEW LEADER IS CHOSEN

(4) REPLICATION

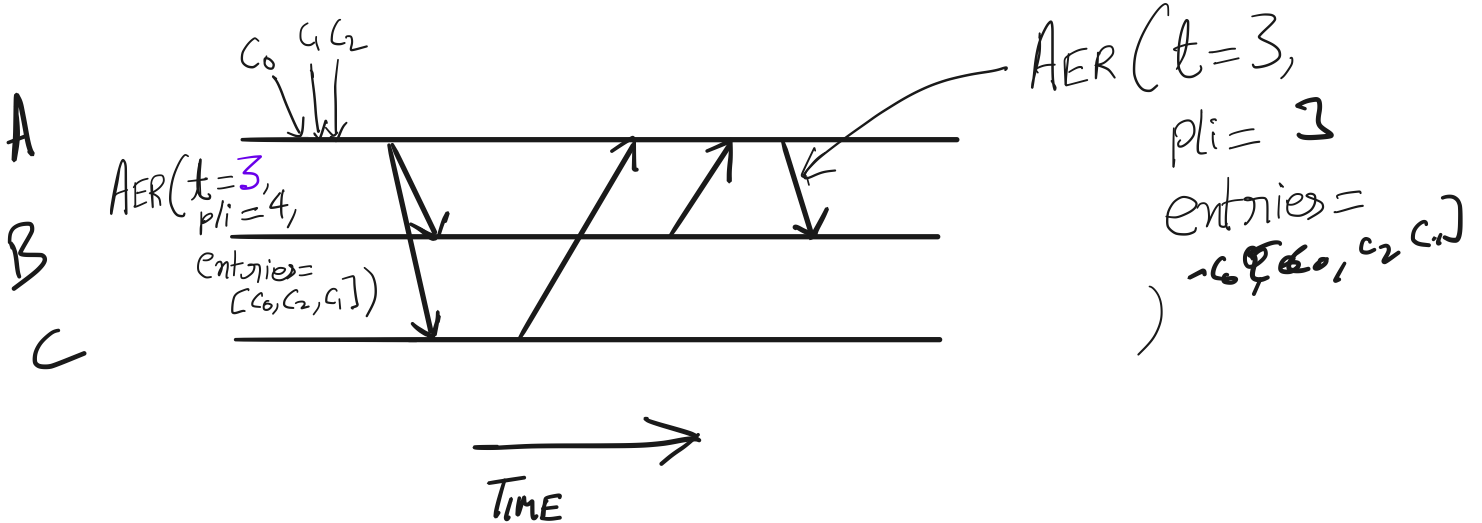
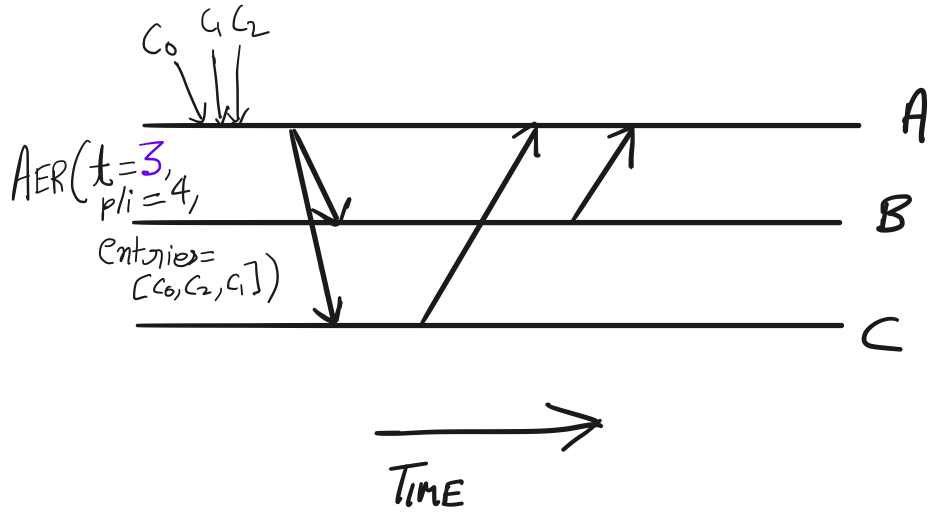
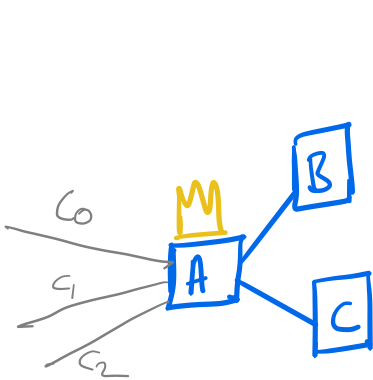
④ RECONFIGURATION

ADDING NEW MACHINES & REMOVING EXISTING ONES

LOG REPLICATION



	0	1	2	3	4	5	6	7	8
A log	2	2	2	3					
B log	2	2	2						
C log	2	2	2	3	c ₀	c ₂	c ₁		



no ||| A execute c₀?

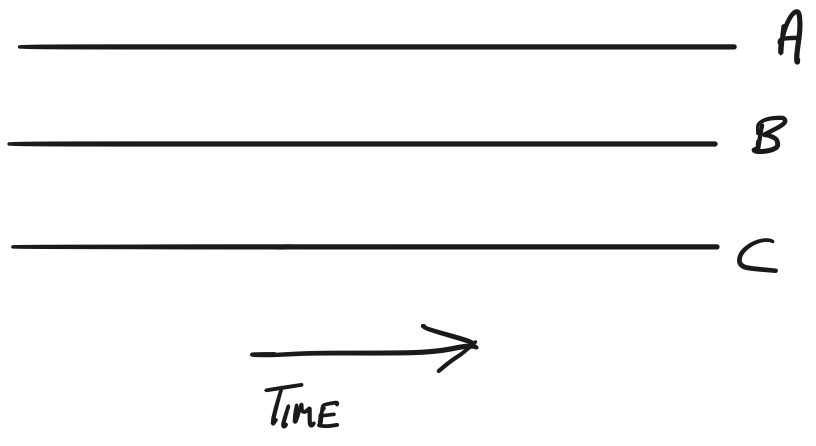
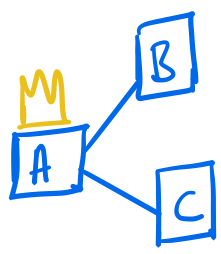
Q: When can A execute C?

After C

Q: When can A respond to C?

After replication/after C

② FAILURE DETECTION



Problem: Leaders initiate all communication.

Silence → Failure

silence → failure
OR
No requests

Solution: Don't let leaders be silent for too long \Rightarrow Heartbeats.

Q: Heartbeat frequency?

Q: How long should a follower wait before assuming leader has failed? (ELECTION TIMEOUT)

Randomized Timeouts

Q: How accurate is this process? Can a follower incorrectly suspect a live leader?

Yes

Q: Leader stability?

Failure Recovery / Leader Election

* IF A MAJORITY OF REPLICAS HAVE LOG ENTRY E AT INDEX I , THEN ALL FUTURE LEADERS WILL HAVE LOG ENTRY E @ I

REQUIREMENTS

(R1) ANY FOLLOWER CAN INITIATE LEADER ELECTION
RANDOMIZED ELECTION TIMEOUT

(R2) * MUST HOLD FOR ANY FOLLOWER THAT BECOMES A LEADER

o CANDIDATE MUST COMPARE LOG WITH A MAJORITY OF REPLICAS

(R3) AT MOST 1 LEADER AT A TIME.

R3° TERM NUMBERS.

- REPLICAS REMEMBER HIGHEST TERM NUMBER THEY HAVE SENT OR RECEIVED
- REJECT ANY MESSAGES WITH LOWER TERM
- **OBSERVATION** ◦ REPLICAS LOG CAN ONLY BE CHANGED BY LEADER WITH CORRECT TERM.
- MUST ENSURE AT MOST ONE LEADER PER TERM

R1 + R2 + AT MOST ONE LEADER PER TERM

ANY FOLLOWER CAN BECOME CANDIDATE
(ON ELECTION TIMEOUT)

① → INCREMENTS TERM NUMBER
(TERM += 1)

IMPLICATIONS?

② BROADCASTS REQUEST VOTE MESSAGE TO ALL REPLICAS

REQUEST VOTE (TERM, ID, LAST LOG INDEX,

LAST LOG TERM

CLAIM: (INDEX, TERM) COMPLETELY IDENTIFIES A LOG ENTRY. WHY?

③ WHEN A REPLICHA RECEIVES REQ VOTE

① UPDATE TERM IF NECESSARY

② Check if already voted for this term (why)

$\xrightarrow{\text{If Voted}}$ DO NOT GRANT VOTE

Log Matching

③ { if (lastLogTerm > NODE.LAST LOG TERM) OR
(lastLogTerm == NODE.LAST LOG TERM
AND lastLog Index >= Node.Last Log INDEX)

Vote yes

else
vote no

③ Write log response from a

waits for responses from
majority. If a majority
vote yes \Rightarrow candidate
becomes leader.

Requirements

1. At most one leader per term.

Why?

2. If a majority of replicas have log entry E at
index I , then all future leaders
will have log entry E @ I

Claim \circ ($lastLogTerm > Node.LastLogTerm$) OR

($lastLogTerm == Node.LastLogTerm$

AND $lastLogIndex \geq Node.LastLogIndex$)

\Rightarrow

IF A MAJORITY OF REPLICAS HAVE LOG ENTRY E AT INDEX I, THEN ALL FUTURE LEADERS WILL HAVE LOG ENTRY E @ I

Why?

	0	1	2	3	4	5	6	7	8
A log									
B log									
C log									

Other issues with leader election

- o Entries from previous terms.

	0	1	2	3	4	5	6	7	8
$\mathbb{M}_{2,4}$ A log	1	2							
B log	1	2							
C log	1	2							
D log	1								
\mathbb{M}_3 E log	1	3							

Problem: Larger term dominates.

④ Reconfiguration

- Add or remove nodes.

Problem: What constitutes a majority?

Solution: Require majority in both old & new configuration

C_{old} → C_{old,new} → C_{new}

System
as described
so far

System where
any entry must
be committed in
both C_{old} & C_{new}
before being
executed

System
as described
so far

Commit configuration changes using Raft