

Distributed Systems ◦ Randomized Consensus

Plan for today

- Some notes on the exam
- Ben-On + Rabin
- Some notes on Lab 4

Exam

◦ Asynchronous

- No bounds on message delays
- No bounds on processing delays
- No global clock

→ Processes P_1, P_2 might not agree on how much time has passed

Note ◦

- Does not mean that processes "have no clock"

clock

receive do

} Logical
time?

... → doct+1
end } Clock!

Core problem: Cannot distinguish b/w delay & failure

• Partial synchrony

After some time (Δ) message and/or processing delays are bounded.

Before \leftrightarrow still asynchronous.

Observe: partially synchronous \rightarrow asynchronous

Cannot result in safety violations.

Why?

Hopefully clearer after next week.

- Π /Set of processes.

in the async/part. sync model

- Impossible to design protocols_n that assume knowledge of what processes are alive.

Why?

Generally $|P|-1$ on $n-1$

Number of processes in configuration

Back to Scheduled Programming

- Focus since RSM lecture: Consensus

Problem: FLP impossibility (next week!)

No deterministic, fault tolerant consensus
protocol in the asynch. model

Consensus: Agreement Validity Termination

How have we gone about this so far?

asynch \rightarrow Partially synchronous

Today: Deterministic \rightarrow Randomized.

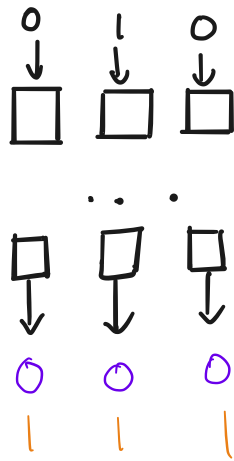
Why? ① Intellectually interesting

② Fewer assumptions

↳ Don't need to set HoBo timer;
Election timer; etc.

Ben-Or

Binary consensus



Convenient problem
to study

- Small set of possible
values

Also features in FLP

- Validity

$\emptyset \ \emptyset \ \emptyset \ \dots \ \emptyset \rightarrow \emptyset$

$1 \ 1 \ 1 \ \dots \ 1 \rightarrow 1$

- Agreement

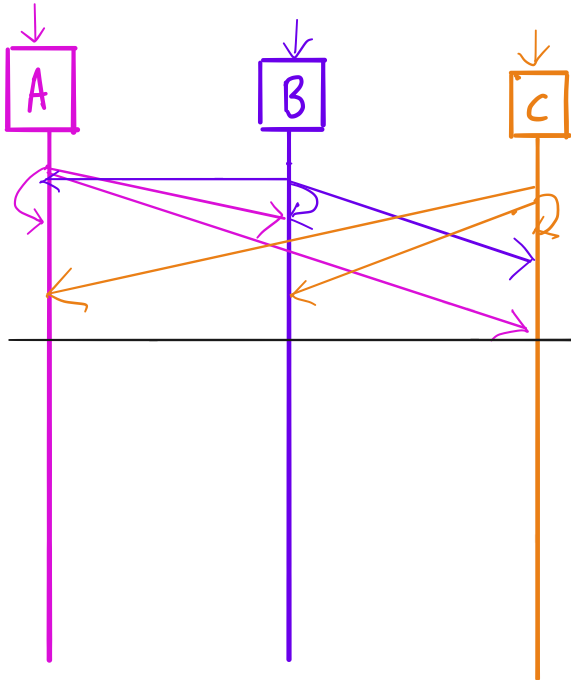
- Probabilistic Termination $p=1$ algorithm terminates

Core mechanism for ensuring agreement:

interaction

- Randomization to break ties

Round 1



$\langle \text{state}, 1, p \rangle$ ($\langle 1, 1, p \rangle$)

wait for $n-f$ messages.

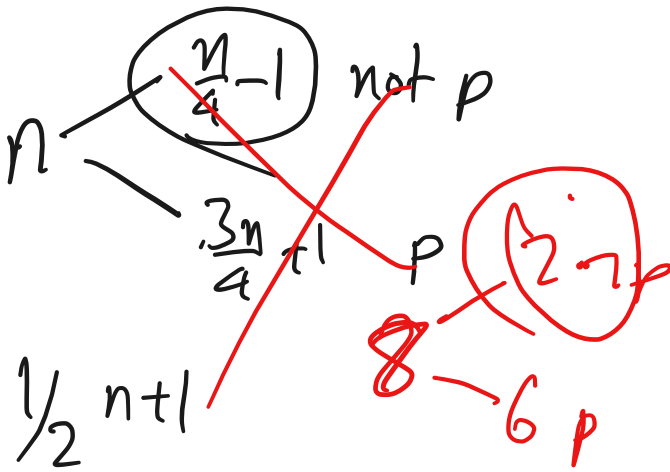
if $\geq \lfloor \frac{n}{2} \rfloor + 1$ messages agree on p
 vote = p else vote = ?

Note

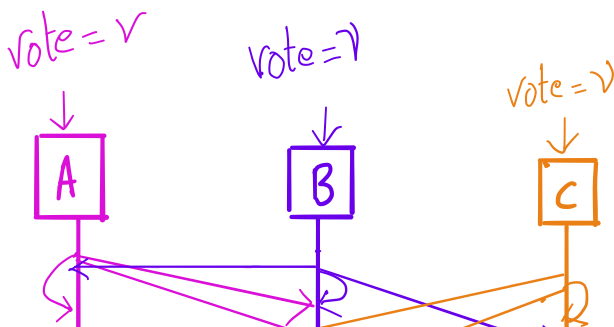
(a) If all nodes have the same input p then all nodes vote = p

(b) What if $\lfloor \frac{n}{2} \rfloor + 1$ nodes have input p ?

(c) What if $\lfloor \frac{3n}{4} \rfloor + 1$ nodes have input p ?

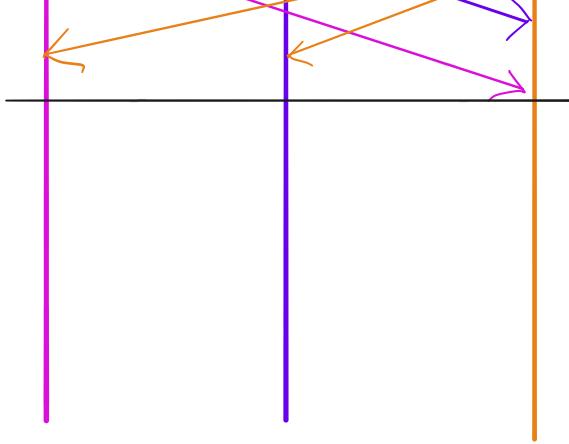


Round 1



$\langle \text{vote}, 1, v \rangle$ ($\langle 2, 1, v \rangle$)

- Reminder $v = \{1 \text{ or } 0\}$ on?



If $v=1$ for B
 ↳ Possible votes for A or C ?

- Wait for $(n-f)$ votes.

If $\geq \lfloor \frac{n}{2} \rfloor + 1$ votes for
 non-? value v

- Decide v

If ≥ 1 vote for non-?
 value v

- set proposal = v

If only ?-votes

- set proposal =
 Coin toss

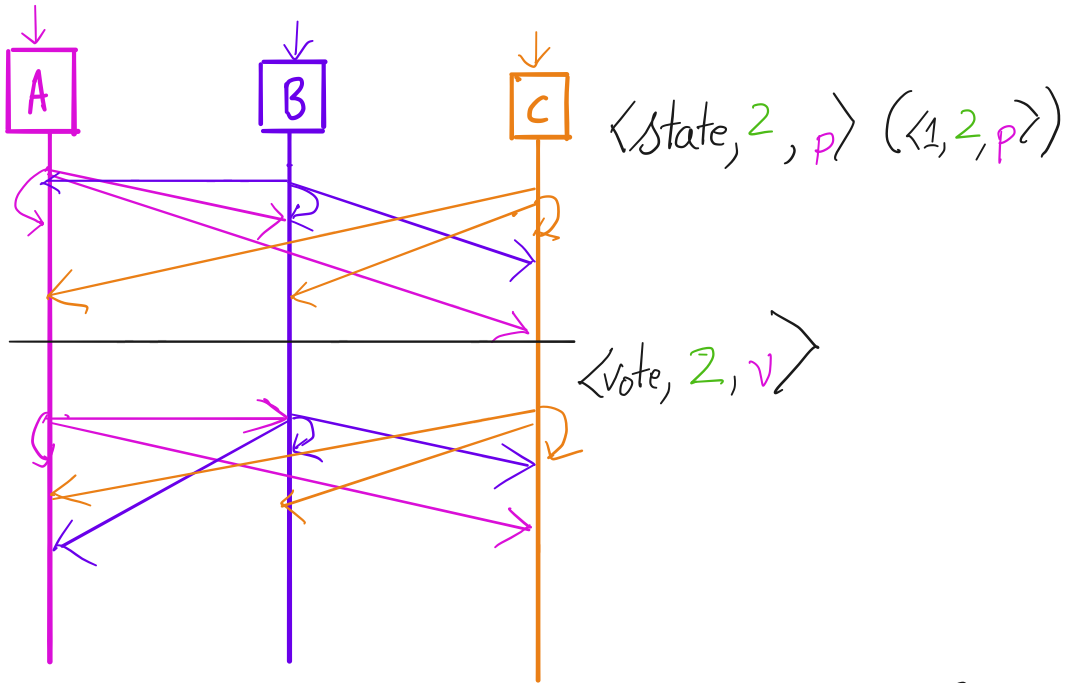
Run Ben-On for
 round 2

Q1. What if all processes
 vote for non-? v ?

Q2. What if $\lfloor \frac{n}{2} \rfloor + 1$ processes
 vote for non-? v ?

Decided the same value

Round 2



Agreement { Claim: If a process decides in round n all live processes decide by round $n+1$. Why?

How validity?

Q. What happens if a node gets only
? votes?

- Common coin?

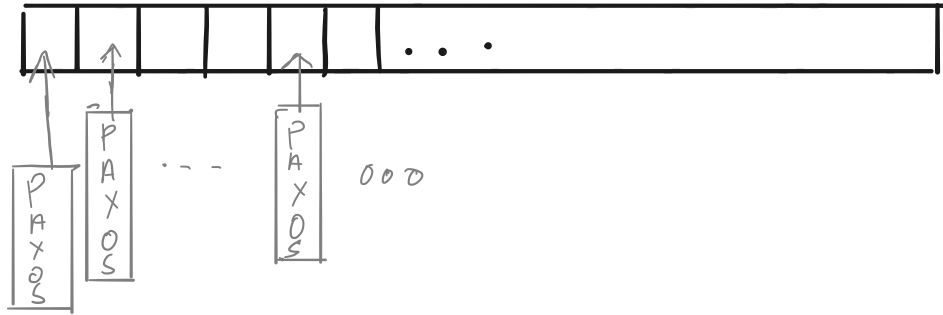
So where we are

- Randomized consensus protocol for $\emptyset, 1$
- Does terminate - see Augilera & Toueg
if interested.

Can we turn this into a RSM?

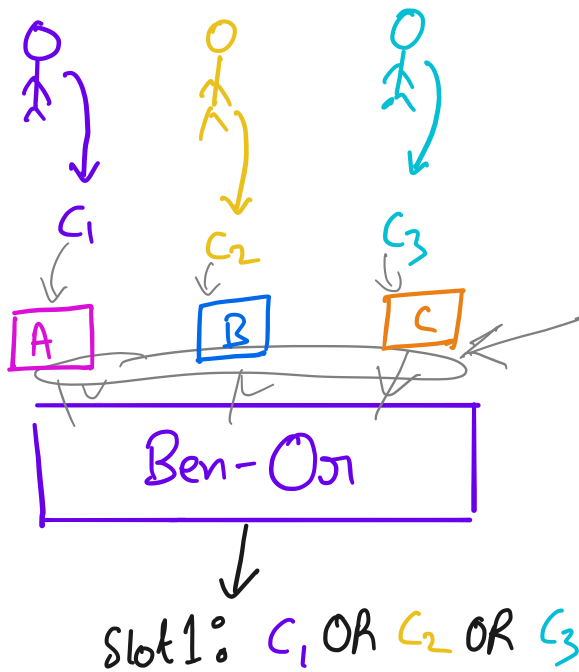
- Borrow an old idea:

Multi-Paxos



Run a Ben-Or round per slot.

New problem: Tyranny of choice



Expects Binary Input

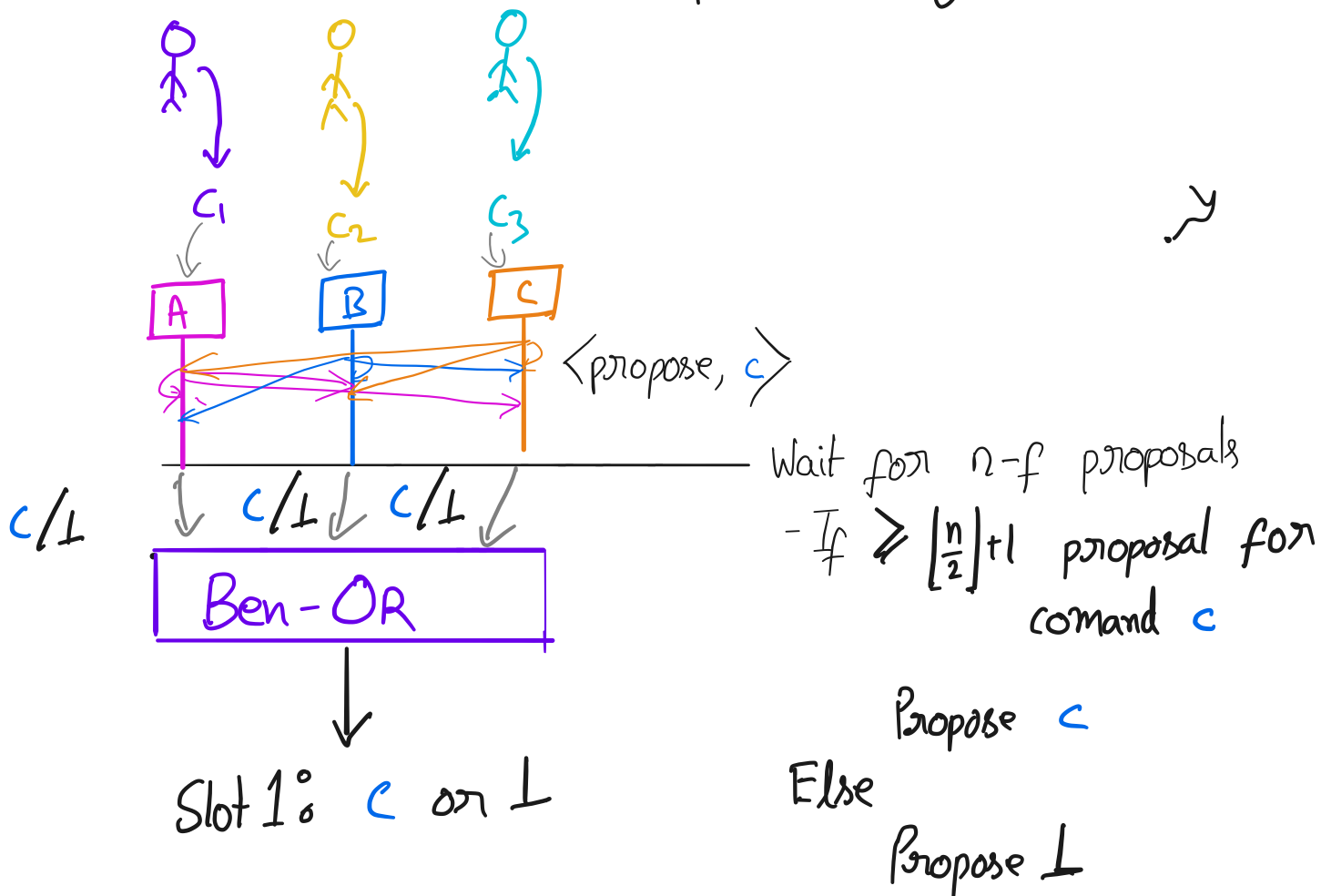
- At most two values proposed (a)

- Possible proposals are known ahead of time (b)

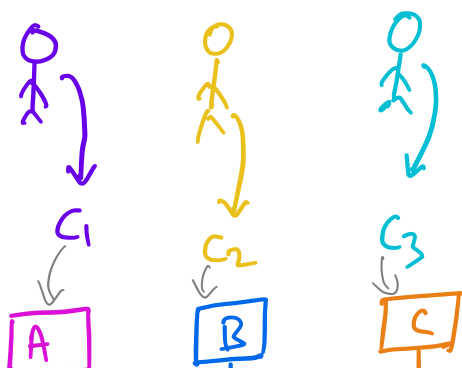
Solving (a): Need to reduce the number of choices

How? Quorum intersection

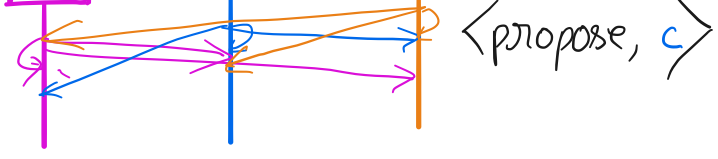
↳ At most one command can be proposed by a quorum



Sub Problem: How to make sure that nodes do not always propose L

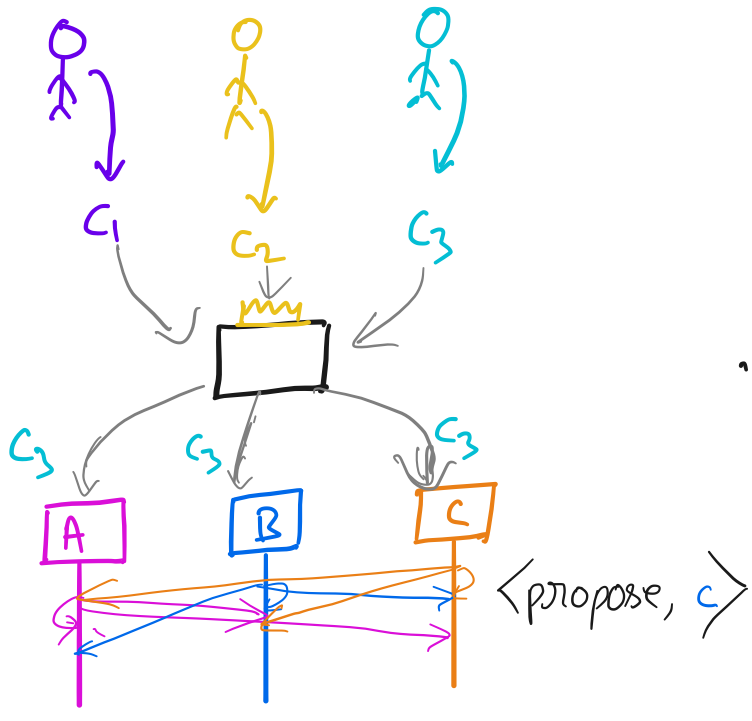


$$c_1 \neq c_2 \neq c_3$$



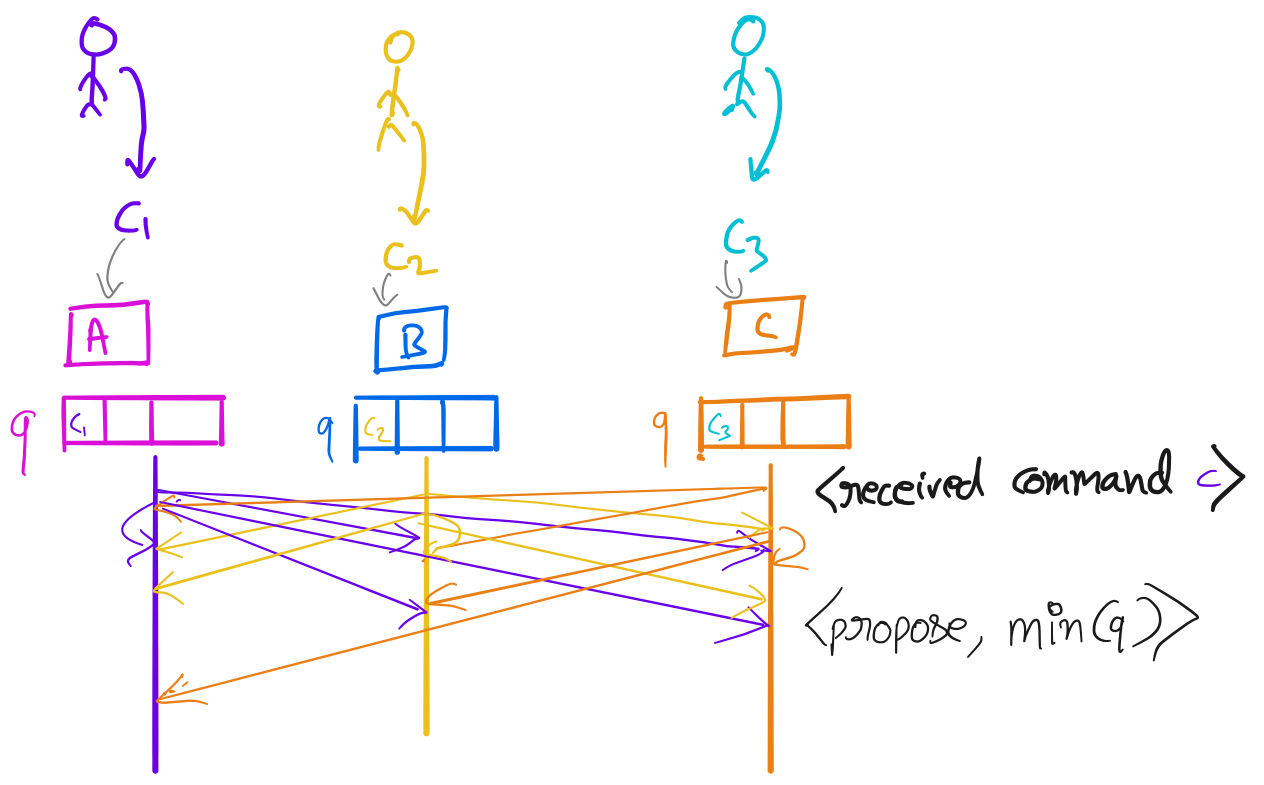
Two solutions

i



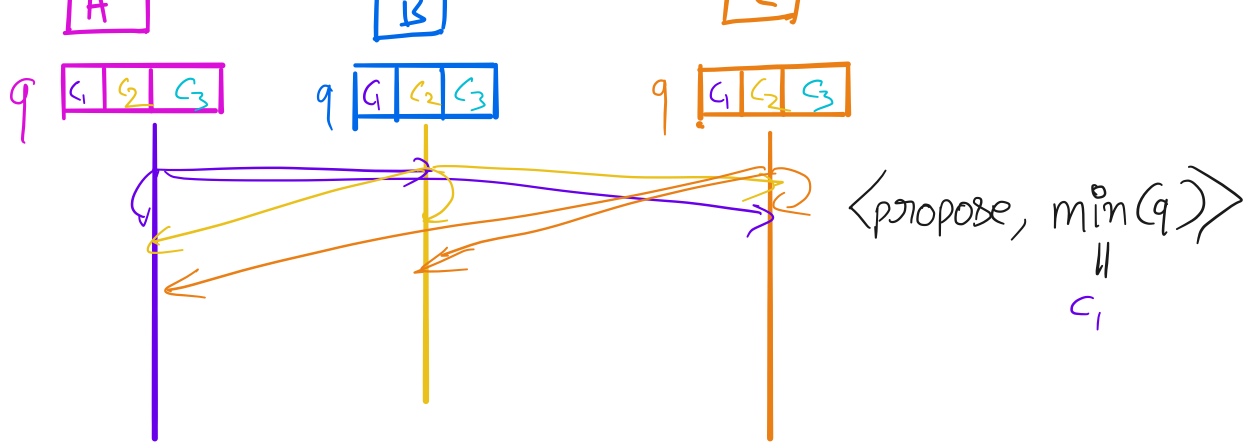
"Proxy" /
"Leader" /
...

ii



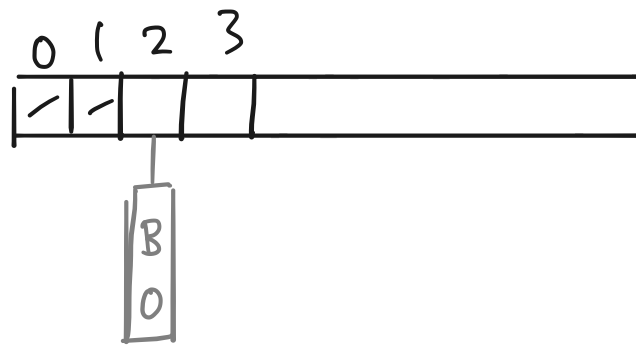
o o o



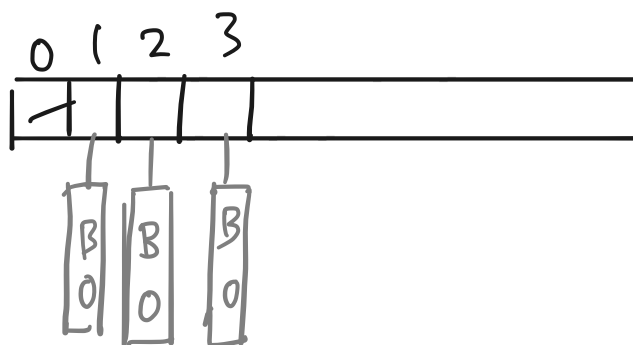


Does it work?

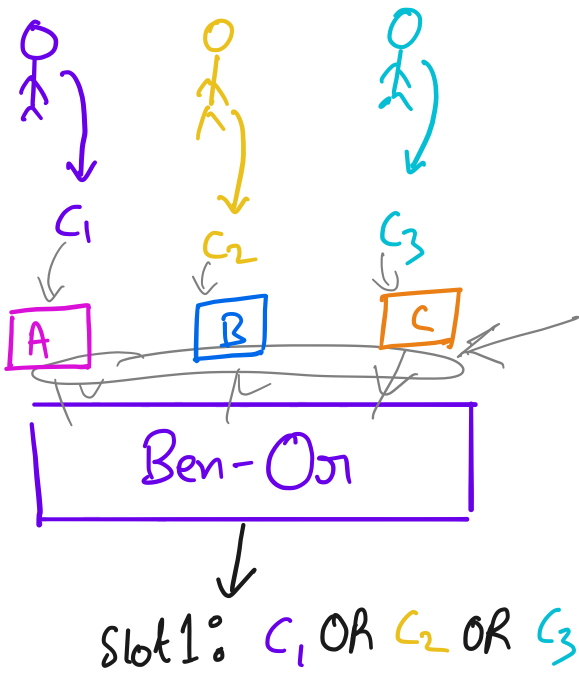
- No pipelining (1 active Ben-On instance at a time)



- Pipelining?



Left as an exercise to the audience.



Expects Binary Input

- At most two values proposed (a)

- Possible proposals are known ahead of time (b)

(b) ← What the errata is about.

Original Idea →

Use Ben-Or to decide 0/1

$0 \Rightarrow 1$ $1 \Rightarrow c$

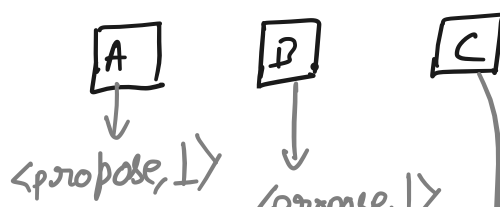
Problem: How to recover c ?

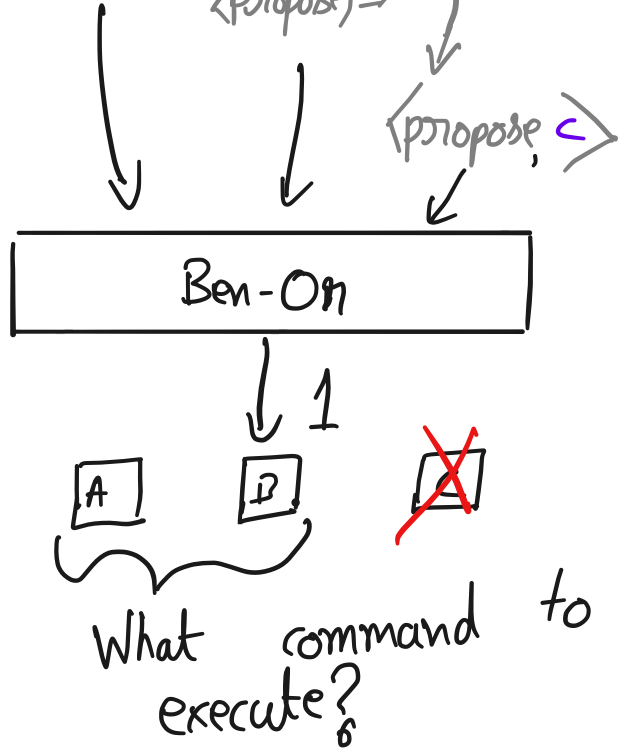
- Quorum proposed c

⇒ At least one live node has c

- But no node might know that

a quorum proposed c





Errata approach.

Ben-On works on c or \perp

↳ During coin toss

Coin must return

\perp or c

How? Builds on the sequence of events requiring a coin toss (no votes other than ?; why nodes vote?).

See errata.

Termination

One thing with Rabin's Weak-MVC (Ben-Or) terminates

```

22: wait until receiving  $\geq n - f$  VOTE messages for round  $r$ 
23: if non-? vote  $v'$  appears  $\geq f + 1$  times in round- $r$  VOTES then
24:   return  $v'$ 
25: else if non-? vote  $v'$  appears at least once in round- $r$  VOTE messages then
26:   state  $\leftarrow v'$ 

```

▷ Termination

```

Algorithm 2 Weak-MVC: Code for Replica  $i$ 
When WEAK-MVC is invoked with input  $q$  and seq:
1: // Exchange Stage: exchange proposals
2: Send (PROPOSAL,  $q$ ) to all           ▷  $q$  is client request
3: wait until receiving  $\geq n - f$  PROPOSAL messages
4: if request  $q$  appears  $\geq \lfloor \frac{n}{2} \rfloor + 1$  times in PROPOSALS then
5:   state  $\leftarrow 1$ 
6: else
7:   state  $\leftarrow 0$ 
8: // Randomized Binary Consensus Stage (Phase  $p \geq 1$ )
9:  $p \leftarrow 1$            ▷ Start with Phase 1
10: while true do
11:   /* Round 1 */
12:   Send (STATE,  $p$ , state) to all       ▷ state can be 0 or 1
13:   wait until receiving  $\geq n - f$  phase- $p$  STATE messages
14:   if value  $v$  appears  $\geq \lfloor \frac{n}{2} \rfloor + 1$  times in STATES then
15:     vote  $\leftarrow v$ 
16:   else
17:     vote  $\leftarrow ?$ 
18:   /* Round 2 */
19:   Send (VOTE,  $p$ , vote) to all         ▷ vote can be 0, 1 or ?
20:   wait until receiving  $\geq n - f$  phase- $p$  VOTE messages
21:   if a non-? value  $v$  appears  $\geq f + 1$  times in VOTES then
22:     Return FINDRETURNVALUE( $v$ )       ▷ Termination
23:   else if a non-? value  $v$  appears at least once in VOTES then
24:     state  $\leftarrow v$ 
25:   else
26:     state  $\leftarrow$  COMMONCOINFLIP( $p$ )   ▷  $p$ -th coin flip
27:      $p \leftarrow p + 1$                  ▷ Proceed to next phase

```

What does termination mean?

ⓐ Stop participating in Ben-Or rounds?

Earlier in class →

Claim: If a process decides in round r all live processes decide by round $r+1$. Why?

ⓑ Mark entry as committed?

- What should a node do if it receives propose, state or vote messages for committed slot?

Takeaways

- Please don't think of this class discussion as saying Raft is complex, broken, or bad.

I really like the protocol,

found that one can implement a version in a day or so.

- Rather just a reflection on the complexity of making these systems practical.

Lab 4

- Original thought:

Implement Raft

Concerns

- Many pitfalls & corner cases, might require more than 2 weeks (Or seems unfair to give it less time than Part)
 - Cuts into final project time.
- For this year:
- Lab 4 gives you a mostly complete Rabia implementation
 - Required task is to integrate it with a state machine (a lin. counter)
 - Extra Credit: Analyze a problem + propose fixes
 - ↳ Probably requires changes to the protocol

While the main lab is simple, the hope is that
a Rabin implementation makes it
easier to connect some of this
to reality.