

Distributed

Systems o Randomized Consensus

- Final project proposals
 - Sent feedback to some of you about concerns
 - The rest seemed great
 - ↳ Very excited about some of the proposals!
- Poll on guest lecture
 - Currently, we have a class on BFT scheduled in 2 weeks

- But goes over the same material as Daniel
- Two options
 - COVER THE SAME MATERIAL
 - USE THE TIME FOR OTHER TOPICS

- Exam observation

Commit? A persistence property

↳ THIS LOG ENTRY WILL
NOT BE OVERWRITTEN/
DELETED.

↳ Regardless of how
the system evolves/
future transitions
(assuming failure model +
correct impl)

What if it is not

- A node (e.g., the leader) knowing
this.

When are entries committed?

EASY CASE: MULTIPAXOS

Phase 1: Collect any entries
replicated to a majority

Lamport's

P2^b : If a proposal with value v is chosen then

EVERY HIGHER-NUMBERED PROPOSAL HAS VALUE v .

An entry is either committed (replicated to a quorum) or not.

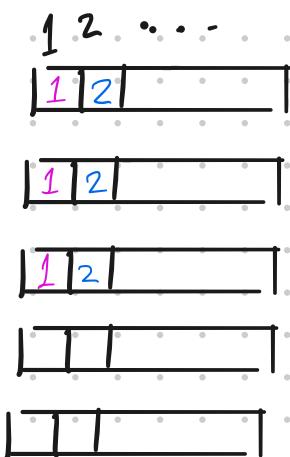
HARD CASE : Raft

- Replicating to quorum is necessary but not sufficient

- But,

T=2 @

Sufficiently replicating entries in the current term allows for commit



For everything else apply the latest log rule.

Randomized Consensus Protocols

Why? Because of FLP [next class]

There is ~~No~~ deterministic, fault tolerant
CONSENSUS PROTOCOL (AGREEMENT, VALIDITY,
TERMINATION) in the asynchronous model.

- But you have already looked at Two CONSENSUS PROTOCOLS
(AND IMPLEMENTED ONE)...

We assumed about more than the asynch. model
↳ Timeouts for Heartbeats & Leader Election

PARTIAL SYNCHRONY [ALSO NEXT CLASS]

There is ~~No~~ deterministic, fault tolerant
CONSENSUS PROTOCOL (AGREEMENT, VALIDITY,
TERMINATION) in the asynchronous model.

Today's goal: Relax a different constraint

There is ^{RANDOMIZED}? ~~No~~ deterministic, fault tolerant
CONSENSUS PROTOCOL (AGREEMENT, VALIDITY,
TERMINATION) in the asynchronous model.

Why?

- Introduces binary consensus \circ Used in FLP.
- May remove some assumptions about deployment — How To SET TIMEOUTS
- Intellectually interesting
- From people here 😊

BEN-OR CONSENSUS

- Binary Consensus

$$l_A = \{0, 1\}$$

A B C

PROTOCOL
RUNS

A B C



DECIDE \emptyset
OR 1

- Agreement: If process p decides v and process q decides v' then $v = v'$

- Termination: Eventually some process p decides.

- Validity: The decided value v must have been some processes input

→ If all inputs are \emptyset then protocol decides \emptyset

$$\emptyset \ \emptyset \ \emptyset \Rightarrow \text{Decide } \emptyset$$

If all inputs are 1 then protocol decides 1

$$1 \ 1 \ 1 \Rightarrow \text{Decide } 1$$

$\emptyset \xrightarrow{1} \emptyset \Rightarrow \text{DECIDE } \emptyset$ Both are OK!
 $\emptyset \xrightarrow{1} \emptyset \Rightarrow \text{Decide } 1$

BEN-OR: Binary consensus with

- Agreement
- Validity
- Termination w/ probability 1

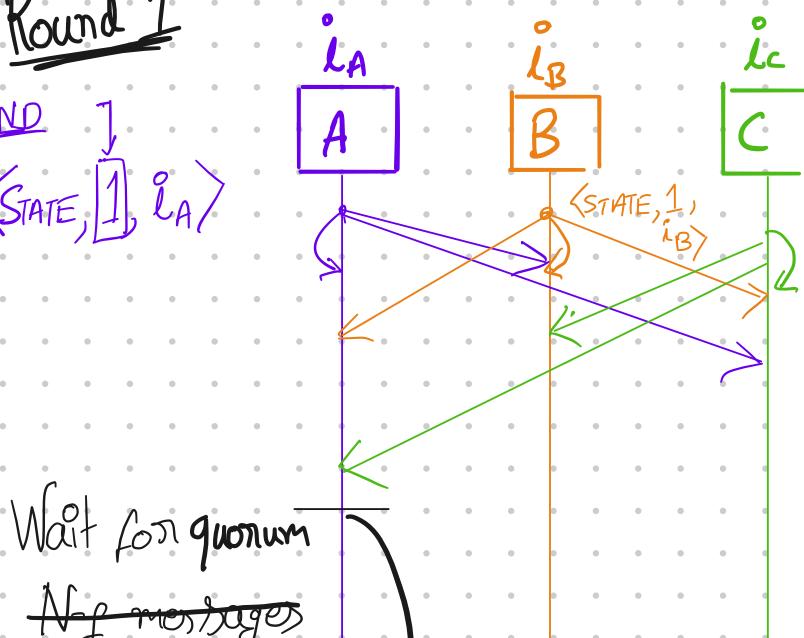
CORE TECHNIQUES

- QUORUM INTERSECTION
- RANDOM COIN FLIPS TO BREAK CONFLICTS

↳ Think back to Paxos discussion

re: termination with competing
proposers Section 2.4

Round 1
ROUND \downarrow
 $\langle \text{STATE}, 1, i_A \rangle$

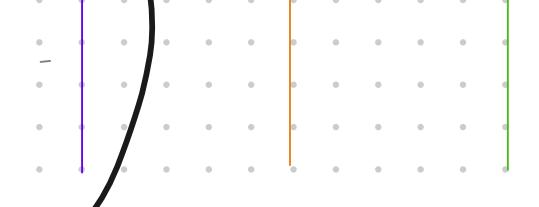


Wait for quorum

~~No messages~~

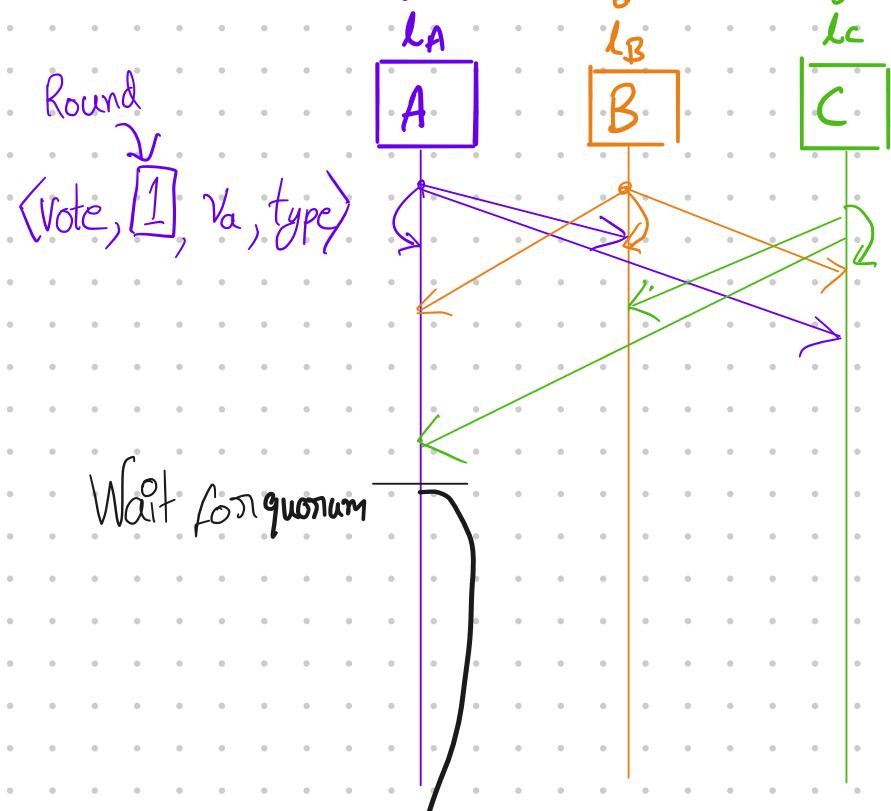
$\langle 1, \dots \rangle$
 $\langle 2, \dots \rangle$
PAPER USES
 MESSAGES OF
 TYPE 1 $(1, \dots)$
 & TYPE 2 $(2, \dots)$

Going to use
 $\text{type 1} \equiv \text{state}$
 $\text{type 2} \equiv \text{vote}$



→ Cases

- Quorum has the same value ✓
- vote value $v_A = v$ type = D
- Else - vote value $v_A = ?$ type = L



→ Cases

- ① One $\langle \text{vote}, 1, v, D \rangle$
- } Some node saw evidence that a majority had the same value v

$$l_A = v$$

$\textcircled{2} \geq f+1 \langle \text{vote}, 1, v, D \rangle$ } A quorum saw
 message that a majority
 had the same value
 \rightarrow

Decide v

[Note, this also means for all
future rounds $i_A = v$]

$\textcircled{3}$ No message of type D

i_A = Randomly choose \emptyset & 1

Go to Round 2

- State

- Vote

...

Why does it work

- What happens if

$i_A = i_B = i_C = \emptyset$ (or 1)

at round R?

Dec. \emptyset

- What happens if a quorum has the same input, e.g., $i_A = i_B = 1$?

1 or ?

- Claim: If node A decides in round R, then other nodes decide in round $R+1$. Why?

$\gg f+1 \langle \text{vote}, R, v, D \rangle$ } A quorum saw
message } that a majority
had the same value
 v

Decide v

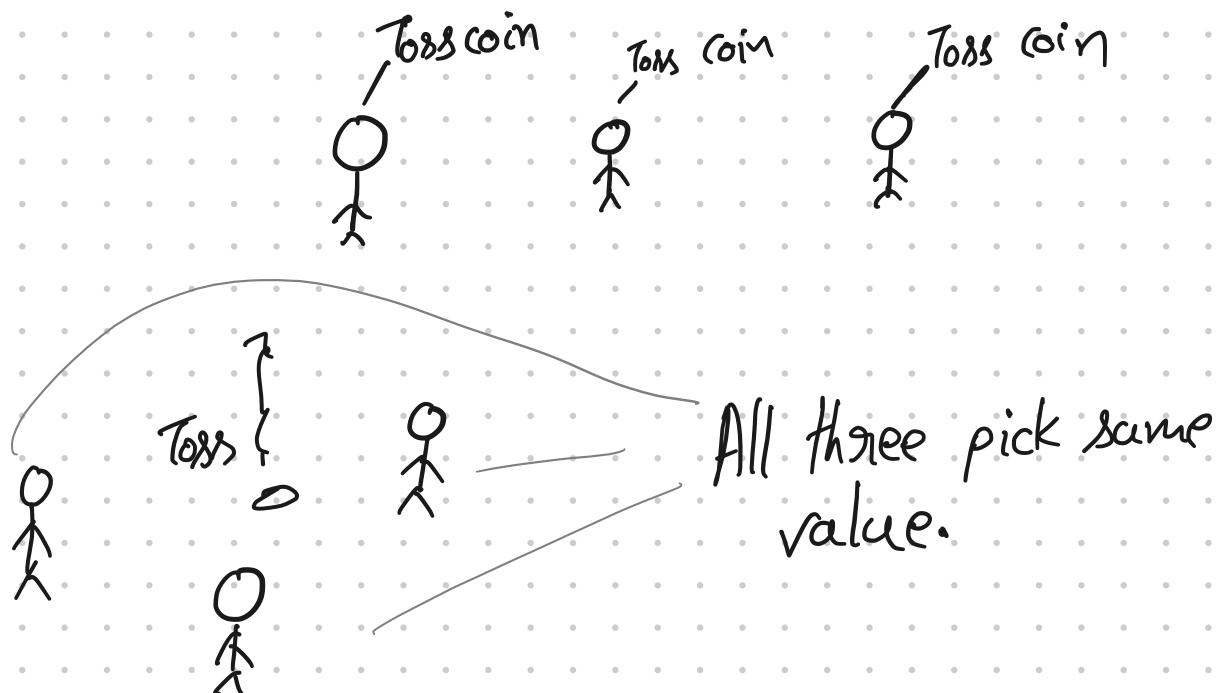
[Note, this also means for all
future rounds $i_A = v$]

Optimization of Common Coin

③ No message of type D

$i_1 = \text{Randomly choose } \emptyset \& 1$

Remember: Want to get all nodes to propose the same value

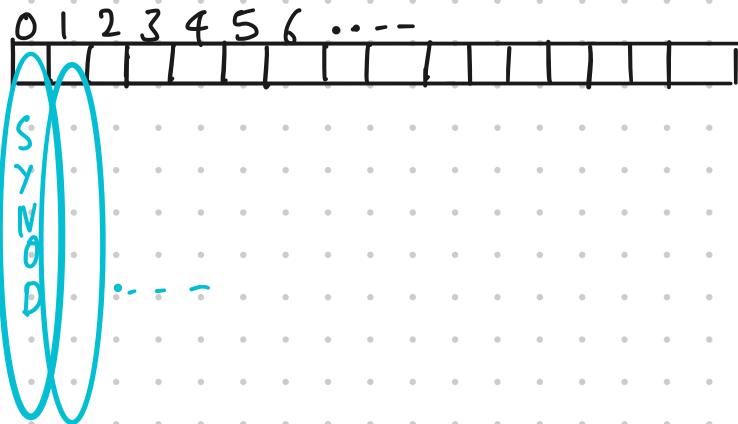


How?

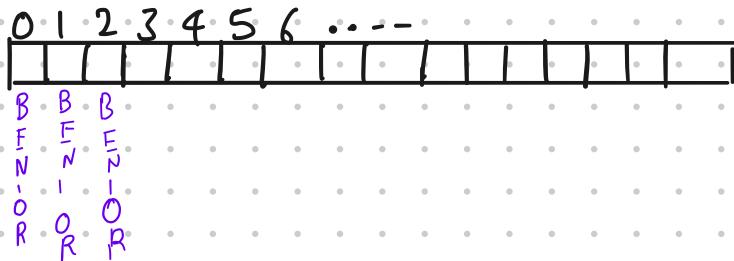
USING BEN-OR IN AN RSM

- Kind of saw this when building multipaxos

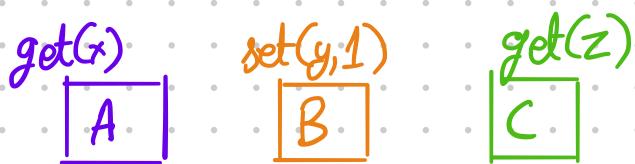
Synod
(One decision) → Multipaxos
(RSM)



- Similar approach



Unique challenge: Binary consensus → Arbitrary commands



WHAT DO A, B & C USE AS
INPUT TO BEN-OR?

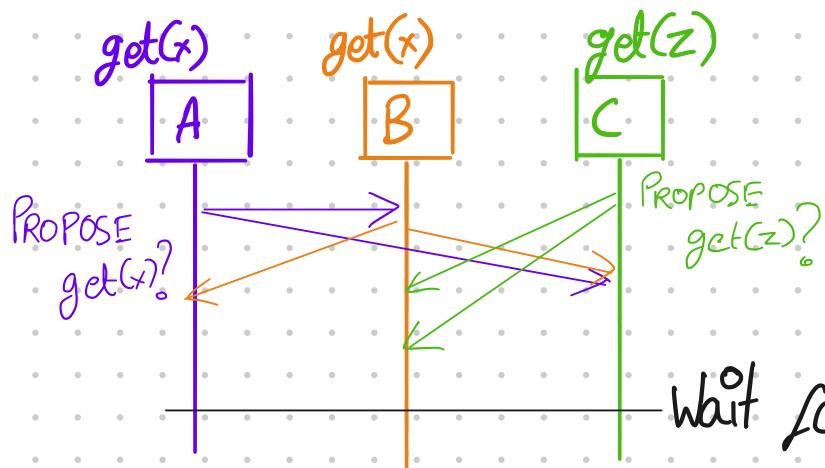
Possibilities — Leader chooses command:

→ That is what we did for multi-paxos.

PAXOS

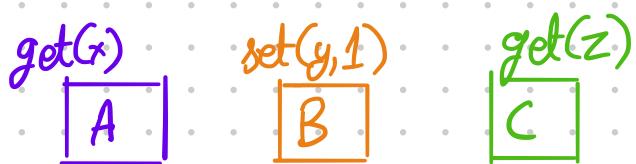
{ - USE A PROTOCOL TO CHOOSE ONE (OR TWO) PROPOSALS

→ Be wary: Why isn't this just consensus?



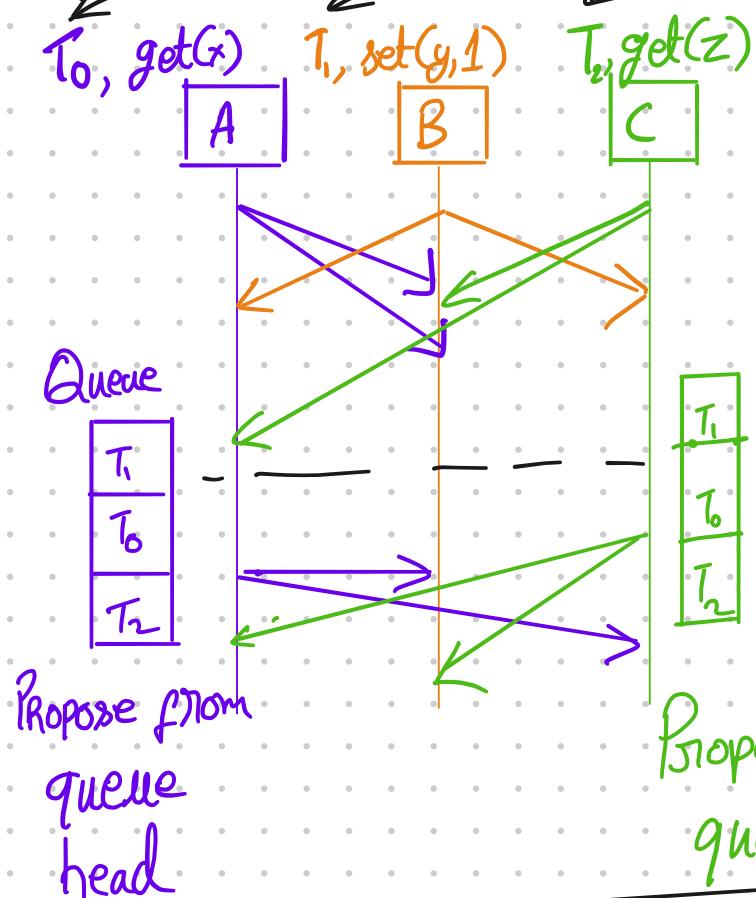
wait for n-f (quorum)

- If quorum has some command
 - Ben On input [for round 0]
 - = C [or COMMON COMMAND]
- Else
 - Ben On input
 - = ⊥



How to deal with this situation?

Purpose of timestamps?



Propose from queue head

- Any node (A) that sees the same command C from a quorum enters Ben-Or with C.
- Other nodes use L

What to do with coin toss?

Remember:

③ No message of type D

i_A = Randomly choose $\emptyset & 1$

Want coin to randomly choose b/w
 $c \wedge \perp$

But, getting to this stage \Rightarrow

- ① Some command c was agreed as a proposal by a quorum of nodes.

Why?

- ② ≥ 1 node did not see a quorum support the command c

Why?

→ Might not know c ?

Errata: Any node arriving at coin toss can recover c

Lab 4

An implementation of Rabia

↳ Mostly done.

→ Aim is to give you a chance to examine
a running version
Without taking significant time from
final project.

A B: vote, 1, \emptyset , D
 ↓ Rd ↓ Val ↙ Tapp
 c. vote, 1, 1, D

① Is it possible

A	B	C
state, R, \emptyset	state, R, \emptyset	state, R, 1
A, B	B, A	B; C
<u>vote, R, \emptyset, D</u>	<u>vote, R, \emptyset, D</u>	vote, R, ?, L
A, B,	B, A	B; C
decide(\emptyset) t.	decide(\emptyset) $L_D = \emptyset$	$L_c = \emptyset$

