

Raft

(Continued)

Reminder: Project proposal — Next week

Reminder from last class (on two)

RSM:— Replicate commands + order in which they are executed

- Execute command once it (and its order) has been sufficiently replicated.

[Not going to be lost due to failures]

Raft - A protocol for this.

Core inv: - Leader based - one leader at a time. Term lockout

- Leader's log is authoritative

↳ Protocol is designed to replicate the leader's log to all other replicas

[Designed to do so efficiently]

AE + Log matching

- Node n becomes leader in term t

\Rightarrow n 's log contains all commands committed before term t

[Leader Completeness]

How? Restrictions on granting votes + becoming leader

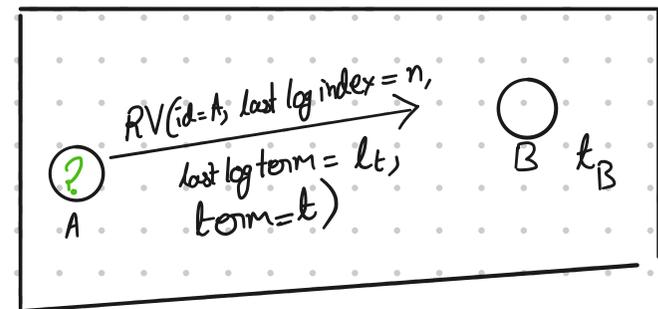
- To become leader in term t , node A must get votes from quorum.

- B grants vote if

- $t_B \leq t$

- Not voted for anyone in term t .

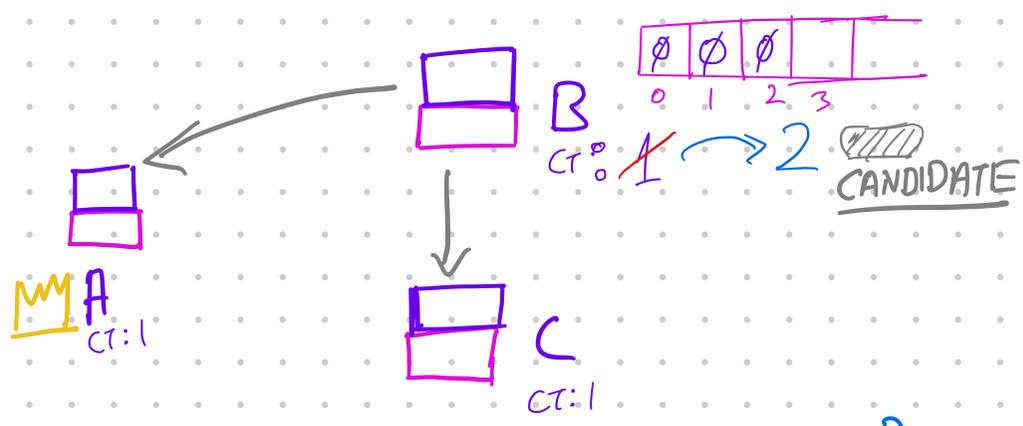
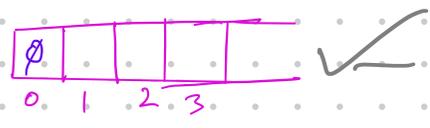
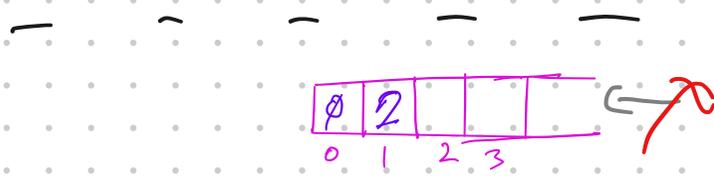
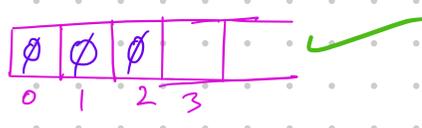
- $t_t >$ term of last log entry at B



OR $[t_t =$ term of last log entry @ B AND

$n \geq$ index of last log entry @ $B]$

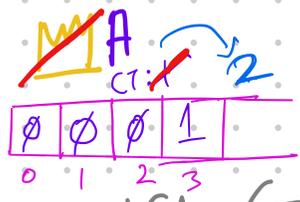
"A's log is at least as up to date as B"



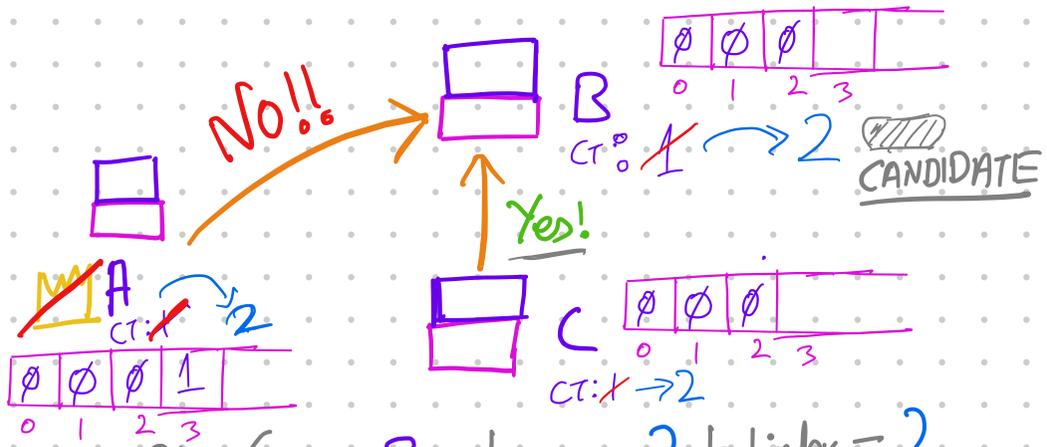
Request Vote (ID=B, term=3, lastindex=2,
lastIndexTerm=∅)



CANDIDATE

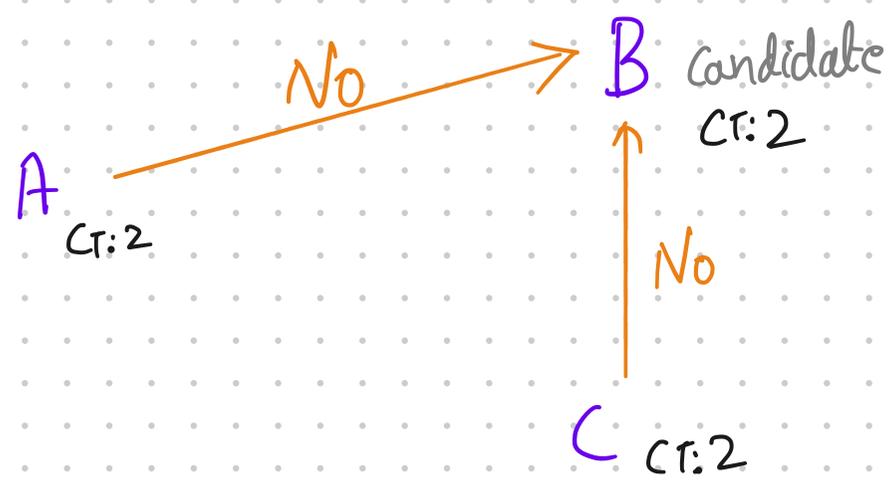


Request Vote (ID=B, term=2, lastindex=2, lastIndexTerm=∅)



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

Observe: Candidate might not become leader

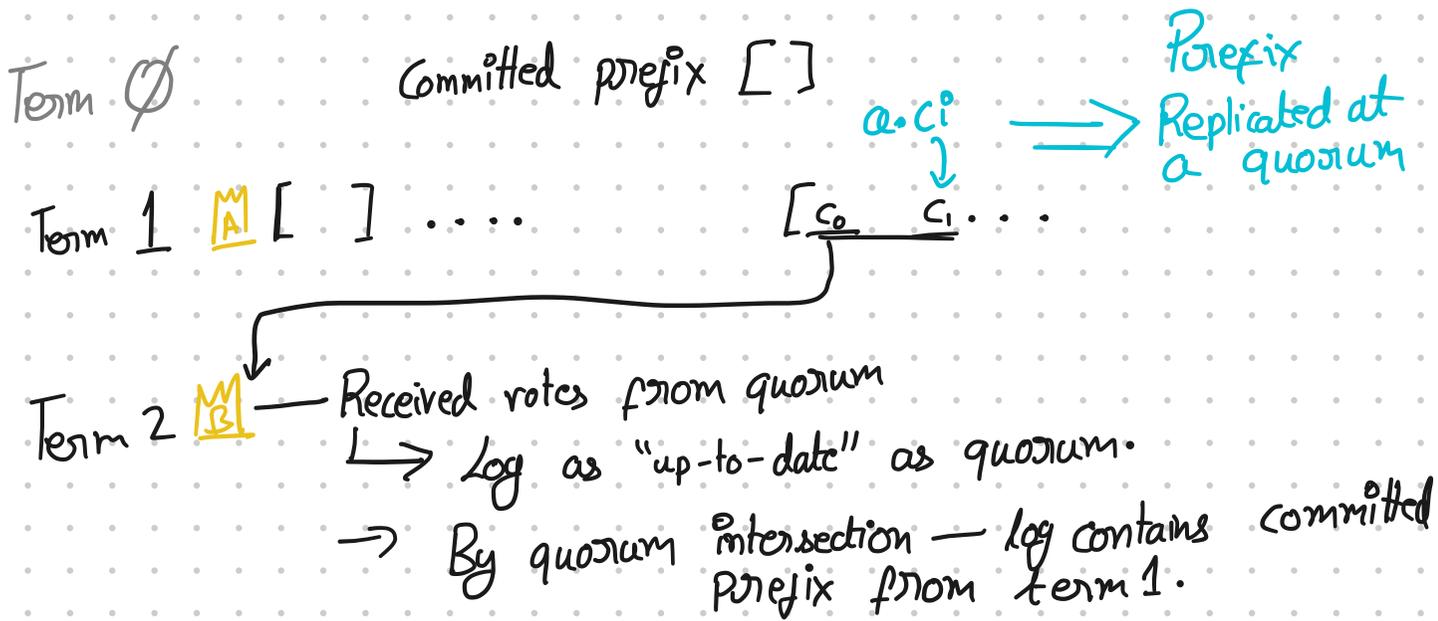


1 1 1 2 1 1 1 1 1 1 1

What now? [Talked about this last class]

- Node n becomes leader in term t
 - \Rightarrow n 's log contains all commands committed before term t
- Observe - Protocol structure \Rightarrow Prefix committed
- [- - - -]

How?



More generally

Easy to show: Term 1 leader's log contains committed prefix
 \hookrightarrow Trivial \rightarrow Committed prefix is empty

Assume term t leader's log contains
committed prefix from term $t-1$,
 $t-1$ leader had $t-2$ prefix, ...

Show any node n that can become
leader in term $t+1$ must contain
committed prefix from term t

↳ Let us say last command committed in term t
is at index p .

Equivalently $[0 \dots i]$ prefix
committed in
term i

⇒ $[0 \dots i]$ replicated
at a quorum

[Protocol requirement to
commit]

Node n can become leader in $t+1$ if it can
get votes from a quorum

⇒ \exists a quorum Q s.t.

$\forall q \in Q$

either

last log entry in q
was committed in a
term **BEFORE** n 's

last last log entry

OR q has fewer entries
than n .

But $\exists q \in Q$, s.t. q 's log

contains $[0 \dots i]$ log

prefix

[Quorum inter-
section]

Say q 's log is

$[0 \dots i \dots j]$

If command f added at term $> t$

$\rightarrow Q$ cannot exist [why]

If command f added in term t

$\Rightarrow n$'s log contains

[0...i...j]

If command j added in term t' [$t' < t$] & i committed in t'' [$t'' \leq t'$]

$t+1 \leftarrow$ Term @ n leader

$t \leftarrow$

$t' \leftarrow$ Term @ q app.

$t'' \leftarrow$ Term @ i com

$\hookrightarrow t'' = t' \rightarrow$ Trivial

$\hookrightarrow q$ log [0...i...j]

$\hookrightarrow t'' < t' \rightarrow$ Inductive hypothesis

$t''+1$ leader log had [0...i] prefix

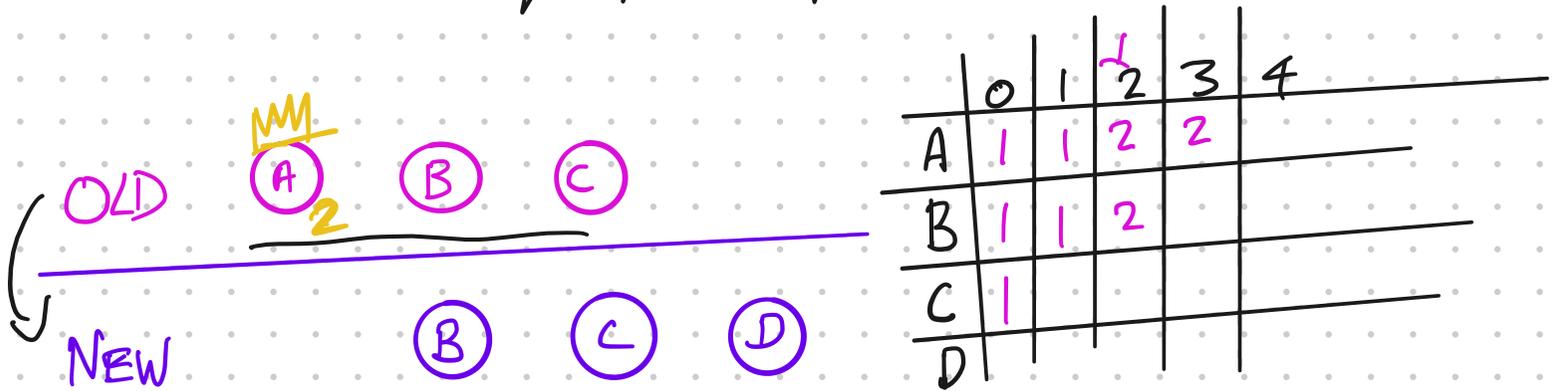
$t''+2$

t'' had [0...i] prefix q log [0...i...j]

Reconfiguration — Joint quorums

— Disclaimer :

- Want to add and/or remove nodes
 - ↳ While continuing to process requests
 - ↳ handle failures

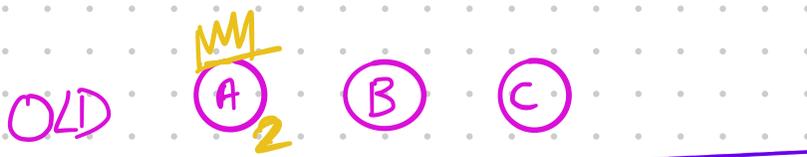


Leader completeness: Who in the new configuration (B, C, D) can be leader?

Leader election protocol: Who in the new configuration (B, C, D) can be leader?

	LI	LI
B	3	2
C	1	1
D	-	-

Can be worse.



NEW

(B)

(C)

(D)

(E)

(F)

Observation: Need some way to maintain leader completeness during reconfig.

- JOINT QUORUM

Building Blocks

(1) Configuration stored as log entry

(a) Usual rules for appending.

⇒ Node A: config C at index 5 in term 2

Node B: config C at index 5 in term 2

What can you say about A's log & B's log?

(b) Unusual Rules for applying

↳ Node uses last configuration in its log: does not wait for commit.

→ But configuration change COMPLETE
entry is committed

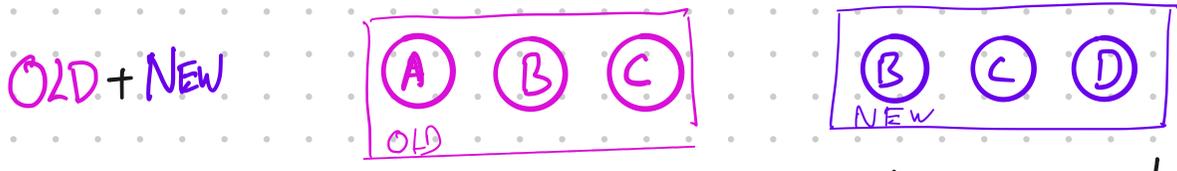
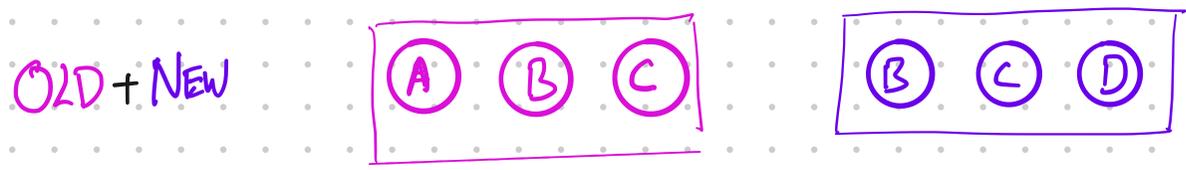
only when entry is

② Joint quorum

Old \rightarrow New implemented

as old \rightarrow old + new $\xrightarrow{\text{Commit}}$ new

Configurations



\rightarrow Same node needs to be leader in both
old & new

Leader Election

	OLD		NEW	
	LI	LT	LI	LT
A	5	2	5	2
B	5	2	1	1
C	1	1	-	-

	1	2	3	4	5
A	1	1	2	2	C_{otN}^2
B	1	1	2	2	C_{otN}^2
C	1				
D					

offer to

— Any entry must be committed to both old & new

OLD			NEW		
	LI	LT		LI	LT
A	5	2	<u>M₃</u> B	5	2
<u>M₃</u> B	5	2	C	1	1
C	1	1	D	-	-

	1	2	3	4	5	6
A	1	1	2	2	C_{otN}^2	
B	1	1	2	2	C_{otN}^2	
C	1					
D						

TRICKINESS: When to go from

Old + new \rightarrow new?

After old + new committed. Why?

OLD			NEW		
	LI	LT		LI	LT
A	5	2	<u>M₃</u> B	5	2
<u>M₃</u> B	5	2	C	5	2
C	1	1	D	-	-

	1	2	3	4	5	6
A	1	1	2	2	C_{otN}^2	
B	1	1	2	2	C_{otN}^2	C_N^3
C	1	1	2	2	C_{otN}^2	
D						

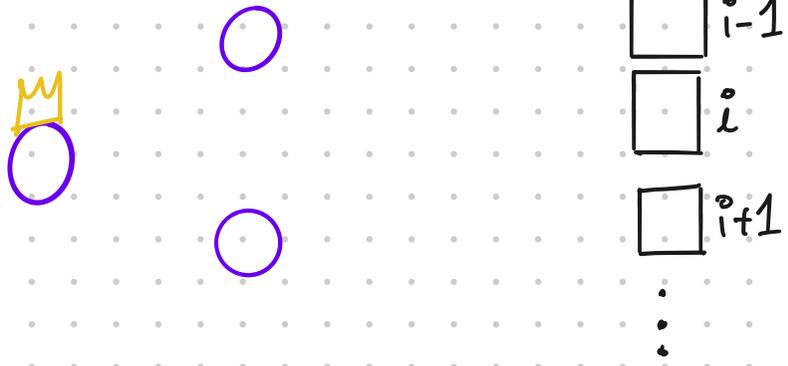
— Where does B replicate C_N^3 ?

- Is $\{A, B\}$ a quorum for $N=3$?

Paxos

Above Ray: Observe - Protocol structure \Rightarrow Prefix committed
 [- - - - [c]]

What if we relaxed this.

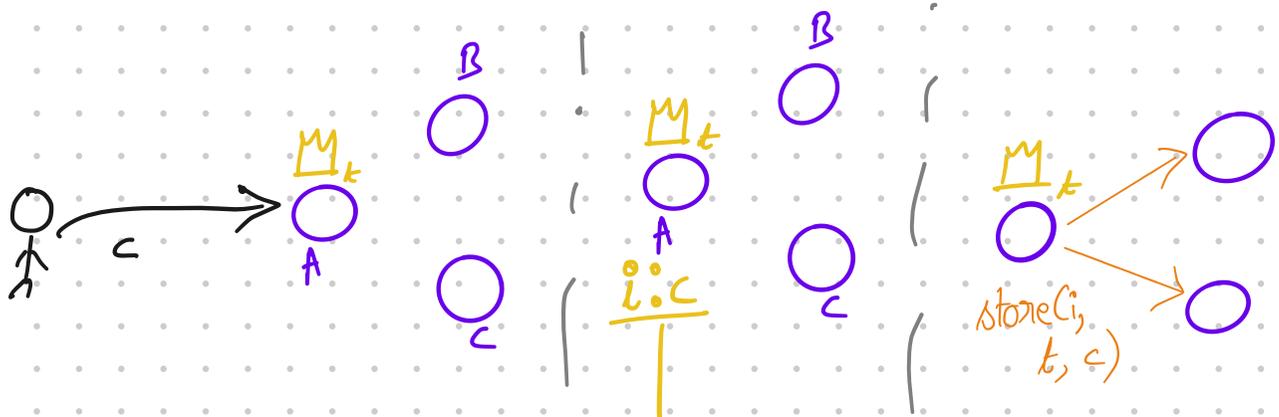


Big Requirements

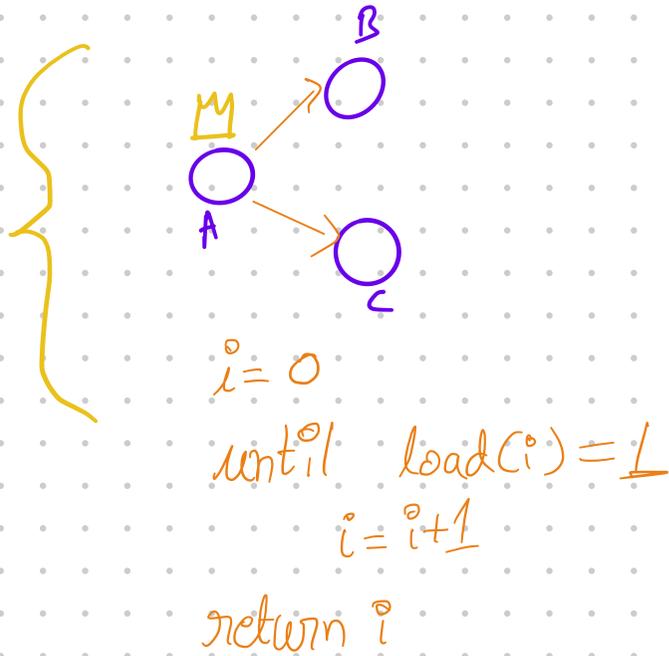
(a) Make sure leader does not overwrite a committed command.

If command c committed at index i then index i should never contain command c' $c' \neq c$

Last class :- Load & store



Must find index i so that i does not contain a committed command



Synod Protocol

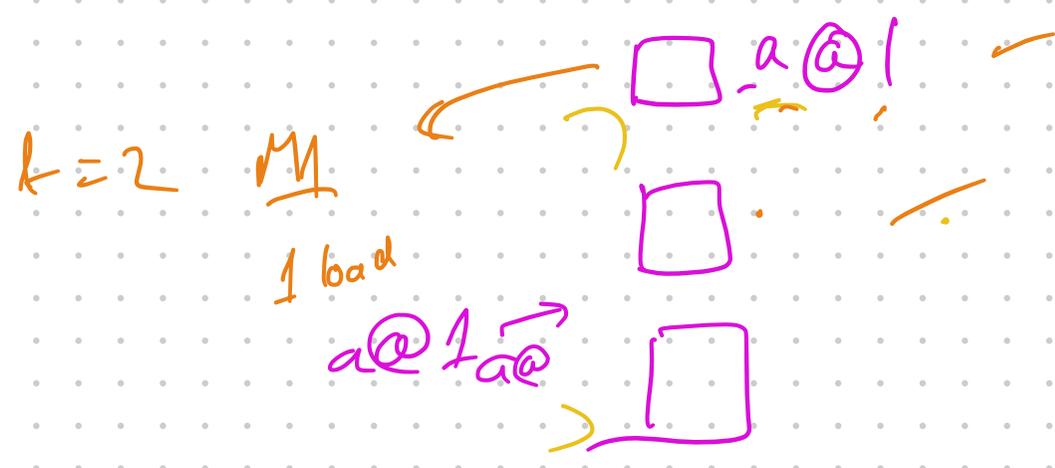
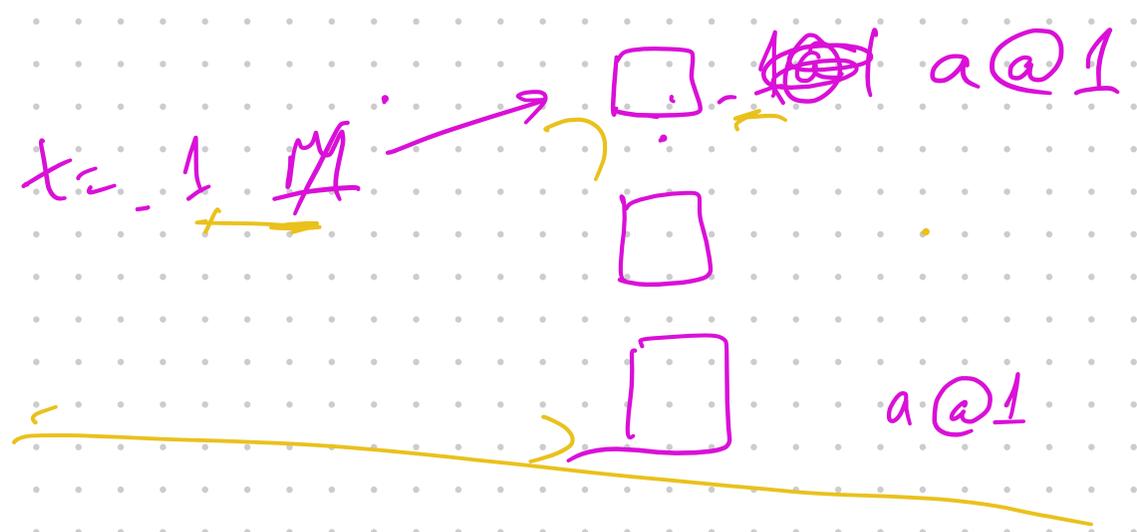
Round / term t

Phase 1 :- Load to find value at slot/index

Phase 2 :- Term 1 Round

(base 2. If phase 1 found
 values $\{v_1@t_1, v_2@t_2, \dots\}$
 pick value v_i at max t_i | Why
 store (v_i, t)

else
 store (c, t)



A 2

$k=7$ B 1

$k=3$ C 1.3

$k=3$ M D 1.3

$k=3$ E 1.3