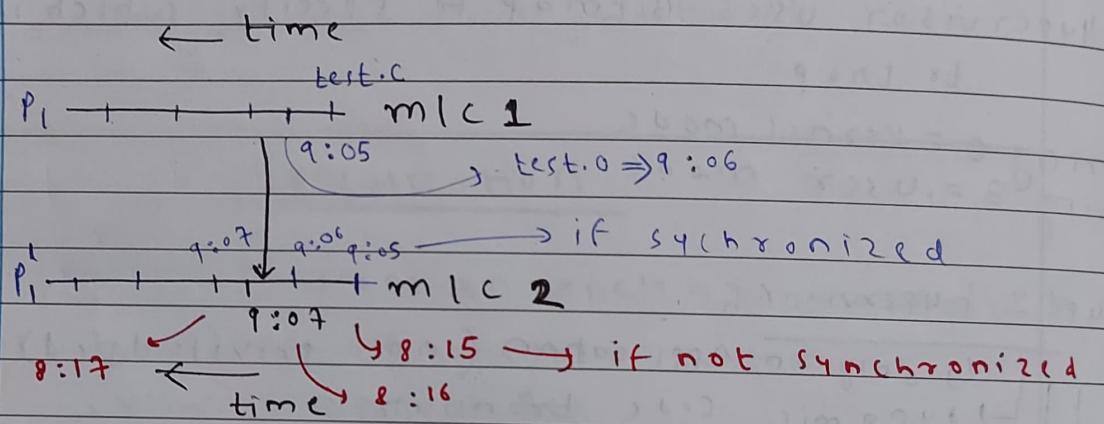


Clock synchronization

- * Solar clock used in ancient times
- * Quartz clock = depend on manufacturer, drift
may change be there
- * If we don't synchronize the clock
 - ↳ synchronizing machines for meeting becomes difficult



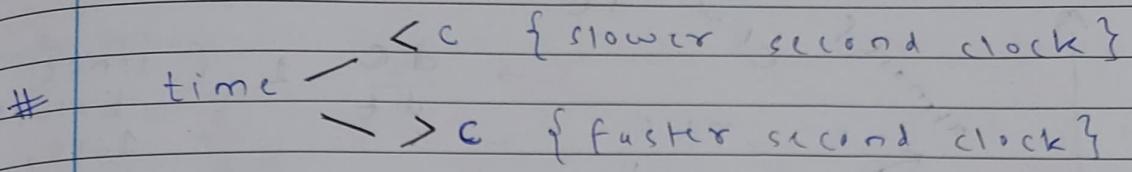
Let test.o send to m1c2, let test.o sleep for 2 sec and after 2 sec, it will run. Let m1c2 start at 8:15, test.o will reach m1c2 at 8:17 (after 2 second sleep).

2 problems:

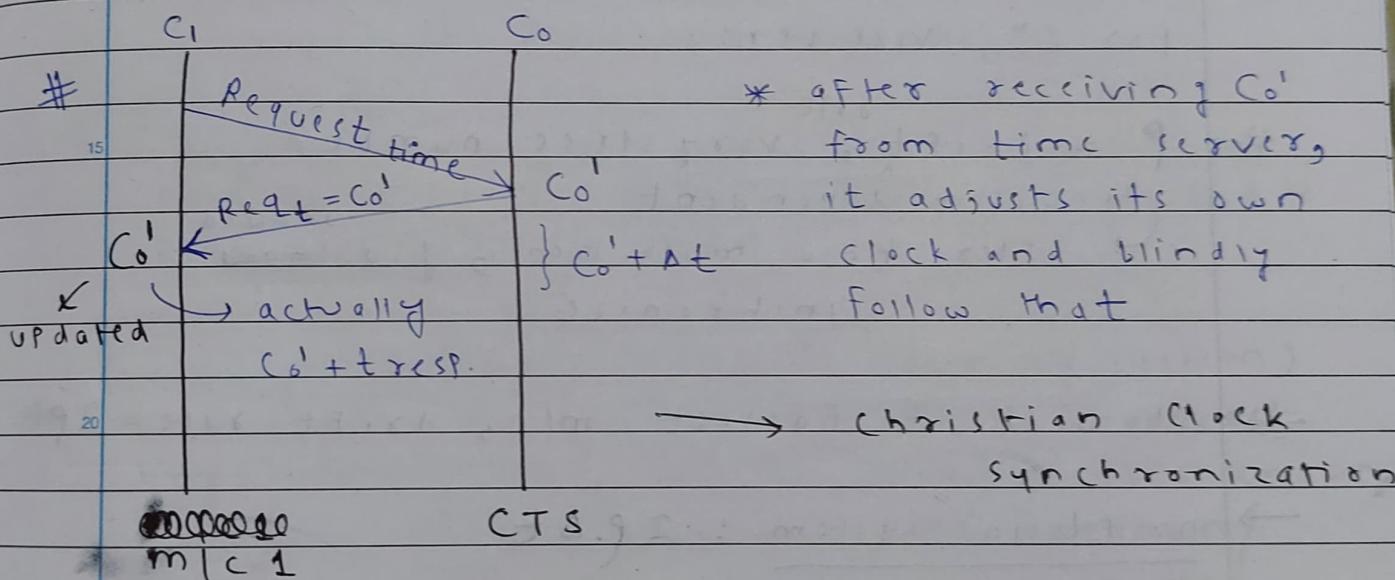
- IF clock of m1c2 behind clock of m1c1, then ~~out~~ code has to wait for some time
- IF clock of m1c2 ahead clock of m1c1, then error will be there and program not be started.

problem
No problem occurs in case of centralized systems

Page:
Date:



- drift
- meeting
- (i) centralized time server (CTS)
- (ii) Berkely Master Clock
- (iii) distributed synchronization
- } synchronization methods



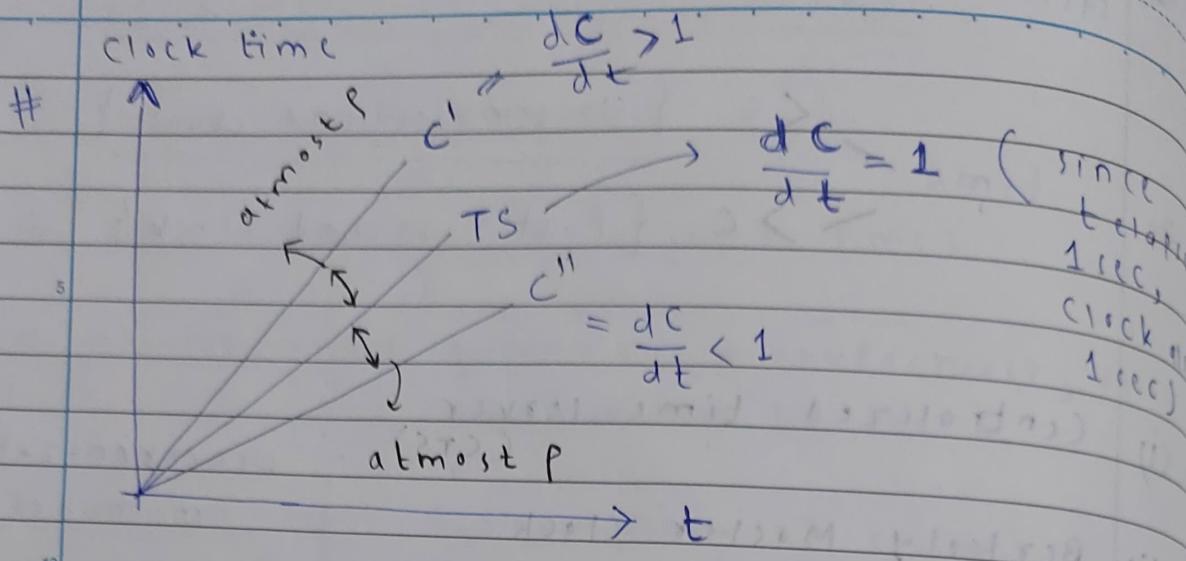
Problem: over the NW, it takes some time
to reach mlc 1 from CTS, so in actual
 C_0' will be equal to $C_0' + t_{\text{response}}$.

$$C_1 = C_0' + t_{\text{resp.}}$$
$$= C_0' + \Delta t \quad \{ \Delta t \text{ time passes at } CTS \}$$

When to synchronize $\Rightarrow ?$

C'' = slower
 C' = faster

Page:
Date:



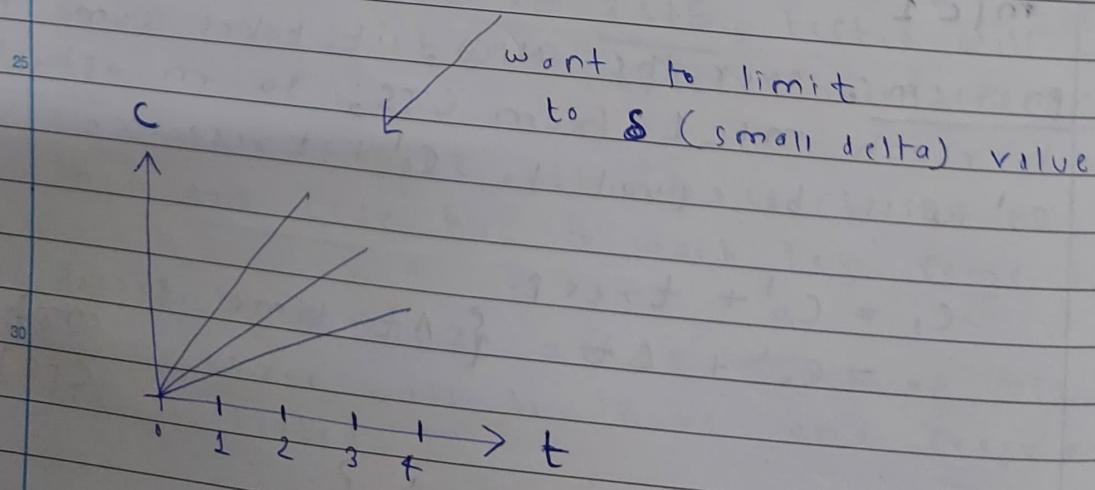
For C' , time elapsed 1 sec, clock elapsed more than 1 sec

For C'' , vice versa

$P = \text{drift rate}$, very less number
 ↴ system dependent
 ↴ after how much time, you should synchronize

→ slowest to fastest m/c, drift rate = $2P$

→ Drift in system : $2P \cdot \Delta t$



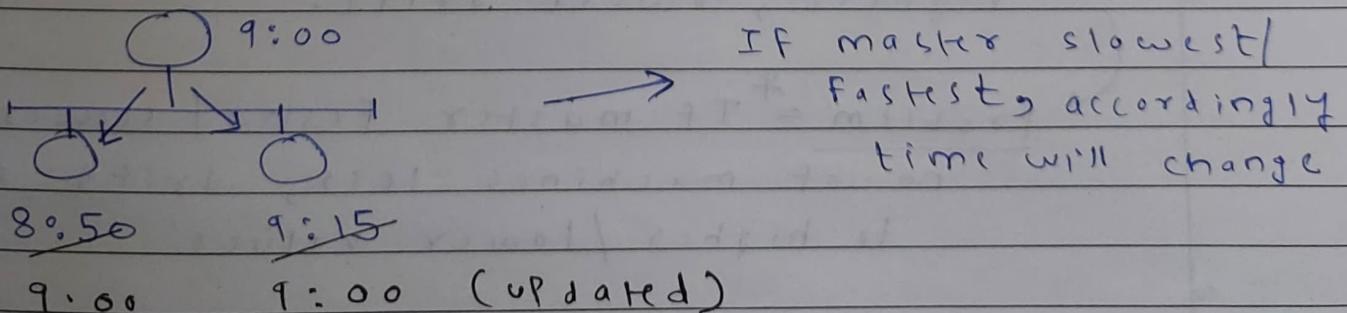
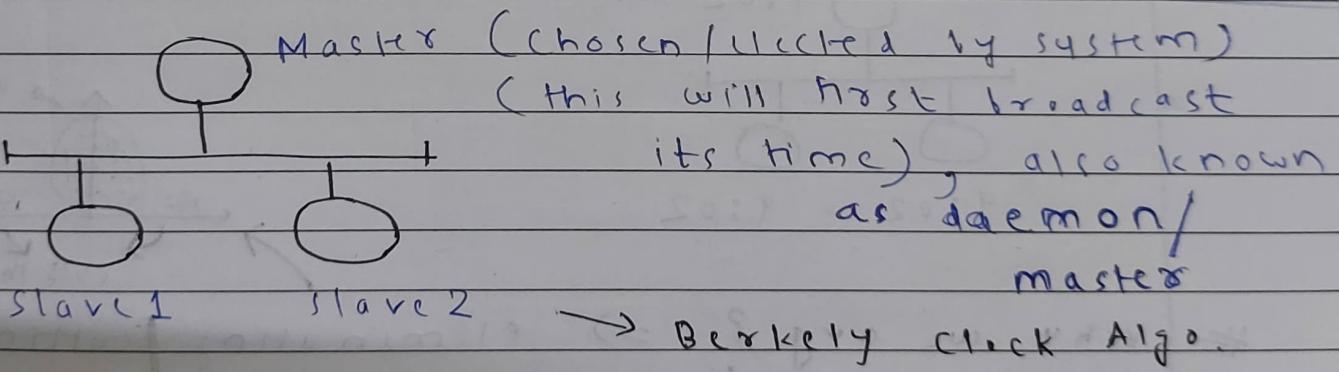
→ When to limit $2P$ to δ , after time unit

$= \frac{\delta}{2P}$, we need to synchronize (send the request to time server and synchronize)

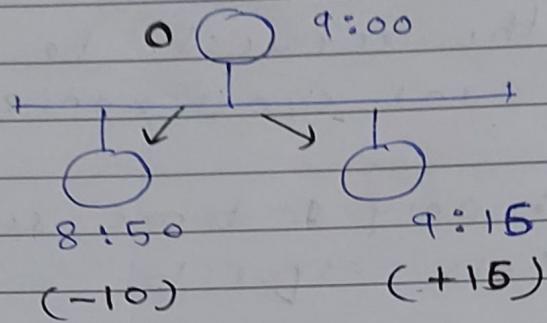
→ Another Problem: only 1 time server, many processes send their request

- ① message loss
- ② link failure since no replica
- ③ Reliability Problems

Therefore had to move from centralized to decentralized system.



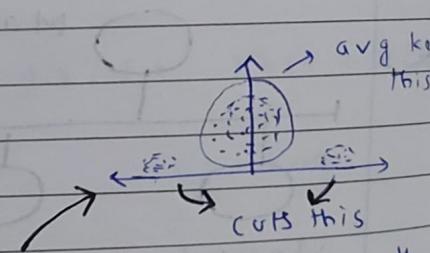
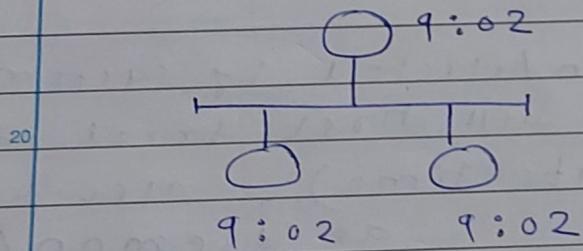
Another solution = Take the average of all drifts and update accordingly



$$\text{average diff} = 0 - 10 + 16 = \frac{6}{3} = 2$$

offset

All m/c updated with Master Clock + offset



Why we avg.? \Rightarrow eliminating / cutting the highest, lowest drift.

problem = If master faster / slower and no. of machines less, drift shifts to higher / lower values

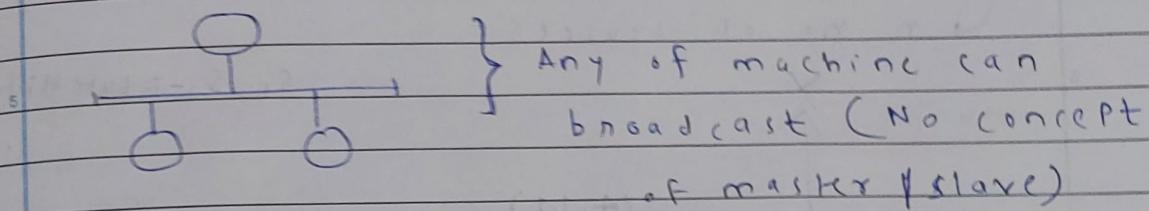
* master broadcasting msg. to slave through 1-way communication
some error may be there (due to single point of failure)

To cluster clock to Avg. clock home
"Reasons for Averaging"

Page :

Date :

Another Approach: no master/slave



→ Benefit = No single point failure, hence reliability is there

→ Problem = Congestion may be there,
Message passing overhead

k + offset

avg keeps this

UB this

Hitting the

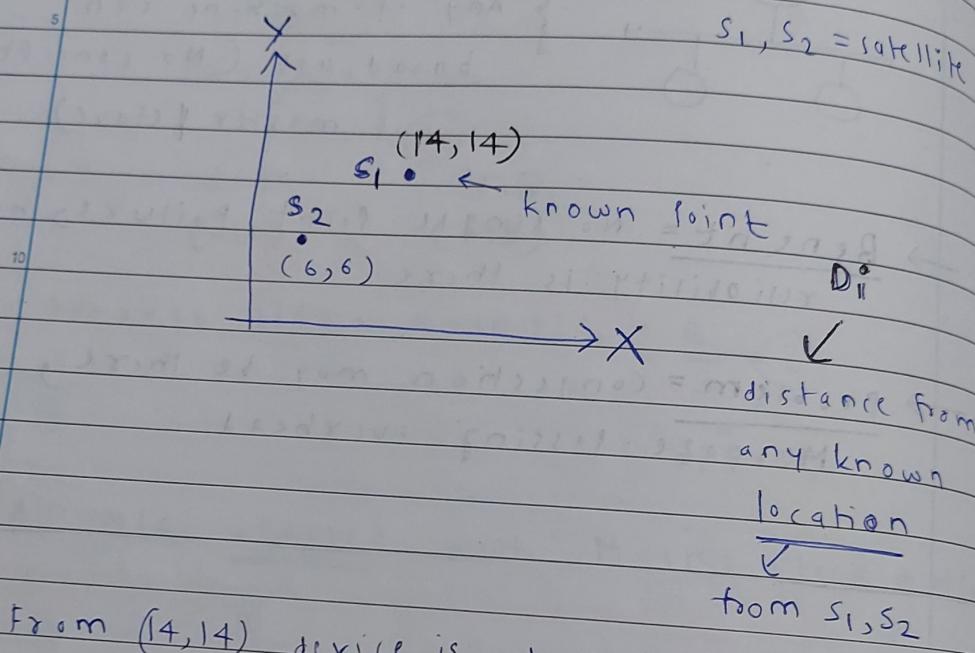
lower and
+ shifted

ng. to
ion,
single
Camlin

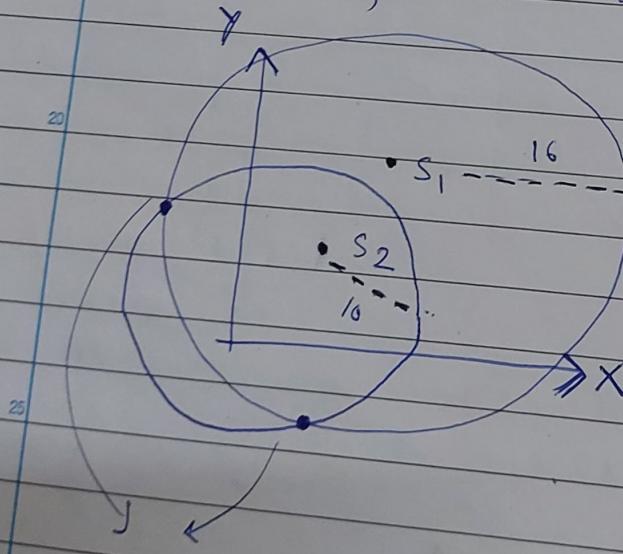
GPS

↳ geostationary satellites

↳ n no. of devices



From $(10, 10)$, device is at radius 16



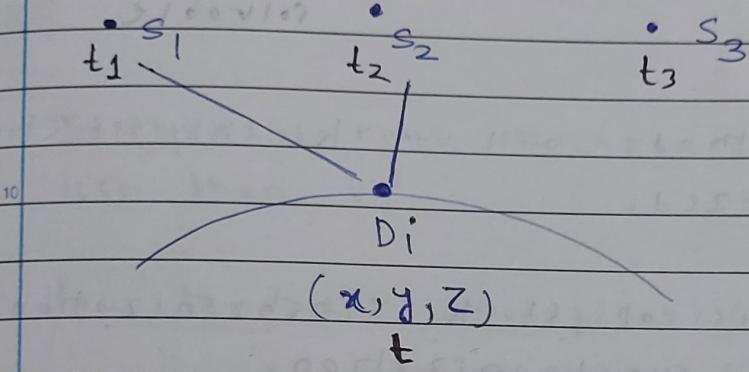
at
device may be either of these points

* Only 2 satellites are not enough for global
location of any device

* Use as many satellites to exactly know the position of device.

How to calc. distance?

$(x_1, y_1, z_1) \rightarrow \text{known}$ $(x_2, y_2, z_2) \rightarrow \text{known}$ $(x_3, y_3, z_3) \rightarrow \text{known}$



unknown = x, y, z

t_1, t_2, t_3 = local
clocks

#₁₅ Beacon \leftarrow its known position
Locat clock

\hookrightarrow sent by satellite to device

$$d = \sqrt{(x - x_1)^2 + (y - y_1)^2 + (z - z_1)^2} \quad \begin{matrix} \text{Beacon} \\ \text{From } S_1 \end{matrix}$$

$$\sqrt{(x - x_2)^2 + (y - y_2)^2 + (z - z_2)^2} = \begin{matrix} \text{Beacon} \\ \text{From } S_2 \end{matrix}$$

$$\sqrt{(x - x_3)^2 + (y - y_3)^2 + (z - z_3)^2} = \begin{matrix} \text{Beacon} \\ \text{From } S_3 \end{matrix}$$



3 unknown, 3 equation but not solvable
atleast 4 equations needed to solve non-linear equations.

Synchronize after 8/2pm
Time

Page:
Date:

$$s = d/t$$

For satellite 1:

$$\text{Time elapsed} = (t - t_1)$$

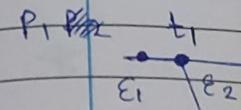
$$d = c * (t - t_1)$$

Now 4 equations are there, therefore solvable

But above methods will work only if clocks are synchronized.

Move from physical clock synchronization to logical clock synchronization.

Lamp



Pipe

In Ph less

How

(i) if ..

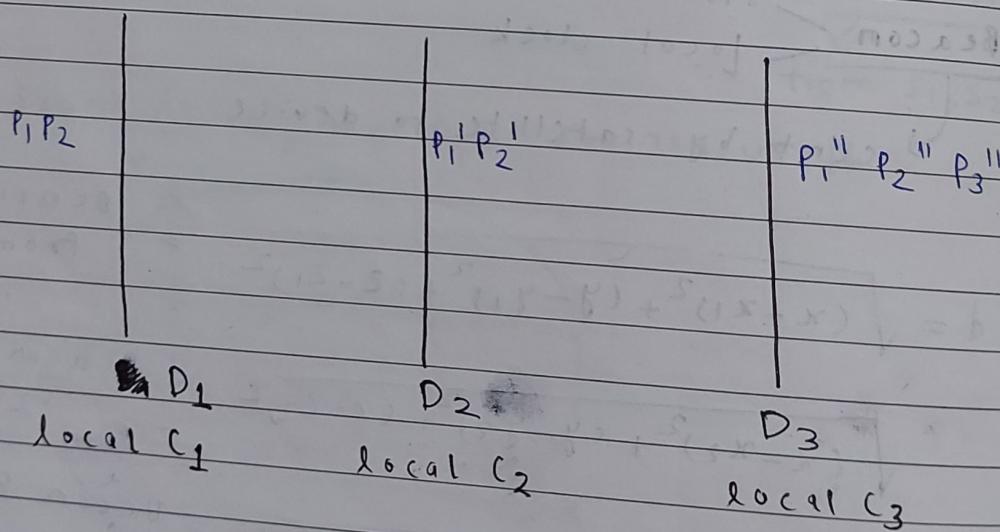
phys

(ii) sen

↳

In a

A →



Among $D_1 \rightarrow D_n$, only subset $\{D_i, D_j\}$ are communicating at particular time

If D_3 not communicating with any other device, no need to synchronize

Ex1:

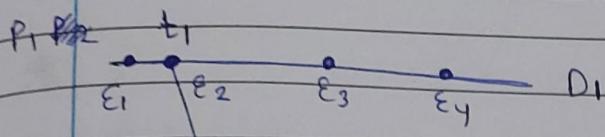
2

E1

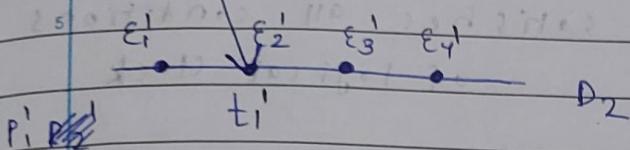
E1

1

Lamport Logical clock (Happened before relationship)



At t_1 , let request sent from D_1 to D_2



$E_1, E_2, \dots, E_4 = \text{events}$
 $E'_1, E'_2, \dots, E'_4 \uparrow$

In physical clock synchronization, t_1 should be less than t'_1

How to order in logical clock synchronization?

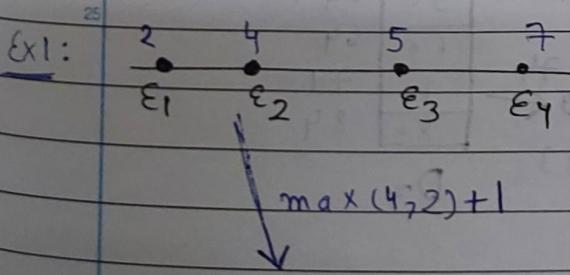
(i) if events belong to same device, ordered by physical clock time of device ($A \rightarrow B$)

(ii) Send-Receive (Request-response)

↳ happen before than the timestamp of Received message

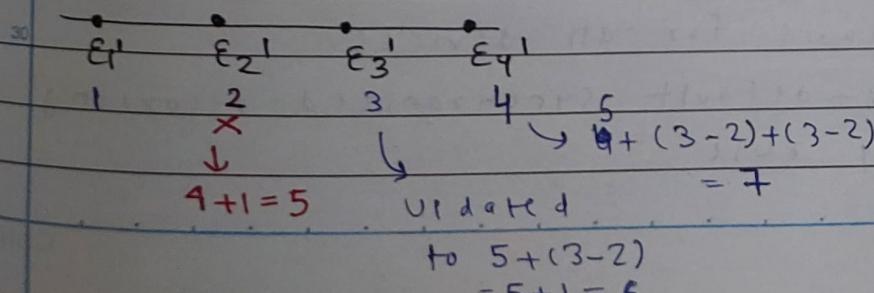
In above example: $t(E_2) < t(E'_2)$

$A \rightarrow B, CT(A) < CT(B)$



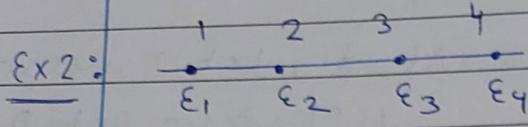
if $t(E_2) > t(E'_2)$

$t(E'_2) = \max(t(E_1), t(E_2)) + 1$



Updated to $5 + (3-2)$

$= 5 + 1 = 6$



\Rightarrow satisfies all conditions
of logical clock

*¹⁰ If $A \rightarrow B$, $B \rightarrow C$, then $A \rightarrow C$

\swarrow \searrow

A happened before B B happened before C A happened before C

Logical clocks follow the transitivity rule.

Ex 3:

Total ordering

0
6
12
18
24
30
36
42
48

$$83 + 1 = 84$$

0
8
16
24
32
40
48
56
64

D₂

0
10
20
30
40
50
60
70
80

D₃

$$\max(43, 40) + 1 = 44$$

54

64

74

84

A \rightarrow B follows for all devices
D₁ & D₃ by default synchronized according to
transitivity rule

event number \leftarrow x.y \rightarrow device number

Page :

Date :

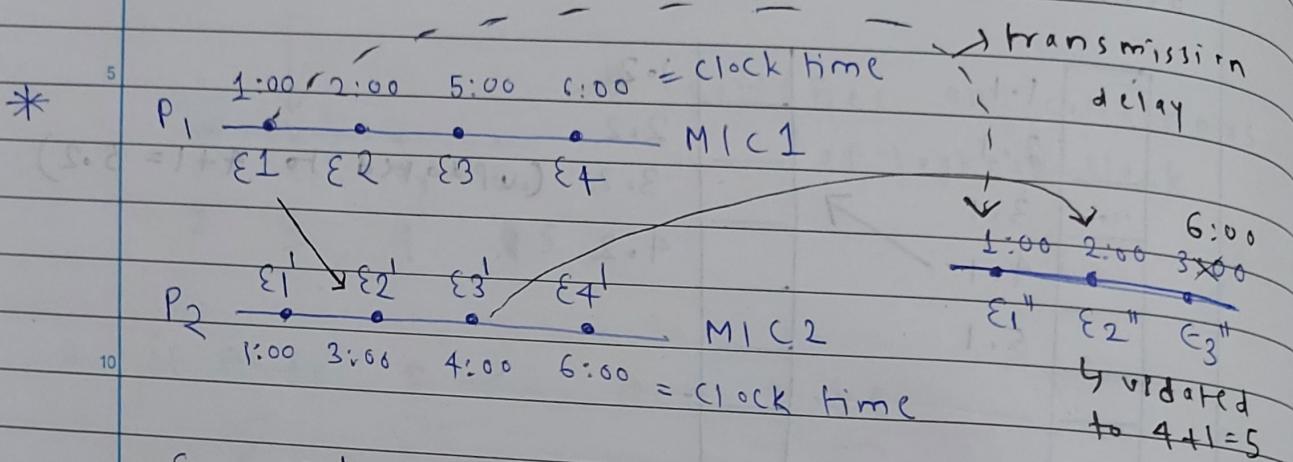
* We can also number the events of different devices.

		1.1		1.2
5		2.1		2.2
		3.1		3.2 (updated to $4+1=5.2$)
		4.1		4.2
		5.1		

Modifications done similarly to previous one.

Lec 3 - AM (14th October, 2022)

if $A \rightarrow B$, then $CT(A) < CT(B)$
 true is also true



$\varepsilon_1 \rightarrow \varepsilon'_1$, ~~$CT(\varepsilon_1) < CT(\varepsilon'_1)$~~ true

logically ε_1 occurs before ε'_1 , due to transmission delay, it does not reach ε'_1 but msg. from ε'_3 reach ε'_2 first.

if $CT(\varepsilon'_1) < CT(\varepsilon'_2)$, we cannot say that ε'_1 occur before ε'_2

↳ limitation of Lamport logical clock
 (knowing clock time cannot conclude which event occur first)

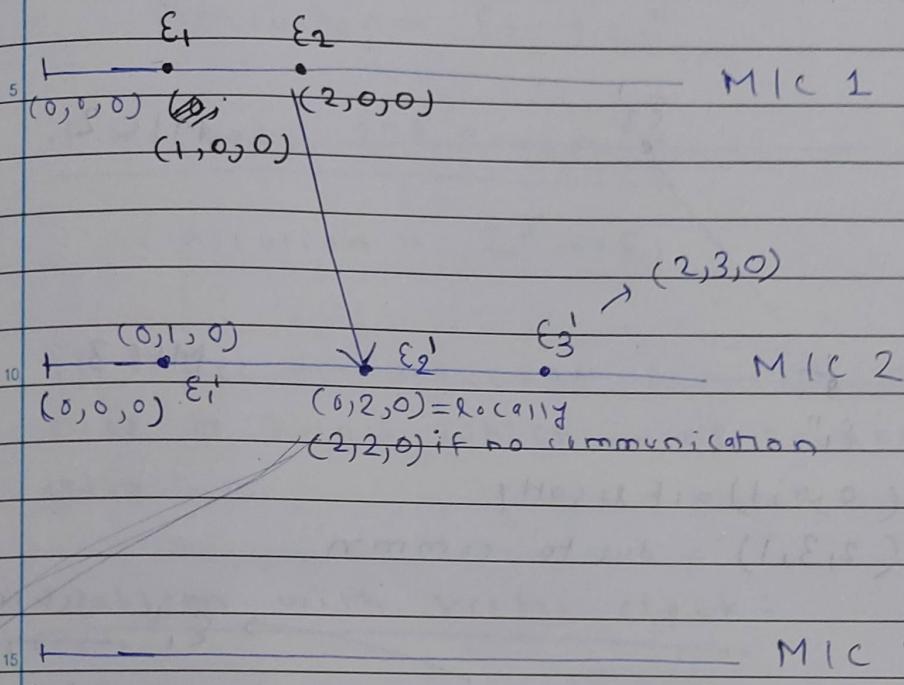


send the clock parameters along
with msg. i.e. Packet

Page:

Date:

vector clock



* Vector Representation:

20 (n_1, n_2, \dots, n_m) = m machines
Timestamp of

Initially vector = $(0, 0, 0)$

25 $(0, 2, 0)$ field
 $\uparrow \uparrow$ = update the ~~field~~ where commⁿ
takes place

*

MIC₁

5

ε_3'

10

↓

ε_1''

MIC₃

(0, 0, 1) = if locally

(2, 3, 1) = due to common.

15

#

Given: $TS(\varepsilon_2) = (2, 0, 0)$

$2 \leq 2$

$0 \leq 3$ before

$0 \leq 1$

ε_2

$TS(\varepsilon_1'') = (2, 3, 1)$

ε_1'' occurred
~~before~~

20

if $CT_i < CT_j$, compare all the values of vectors

NOTE:

if $CT_i[1] \leq CT_j[1]$, then event j occurred before event i

25

1 3rd

if we take 1st parameters: (only)

(inclusion = $\Sigma_3 \rightarrow \Sigma_1''$)

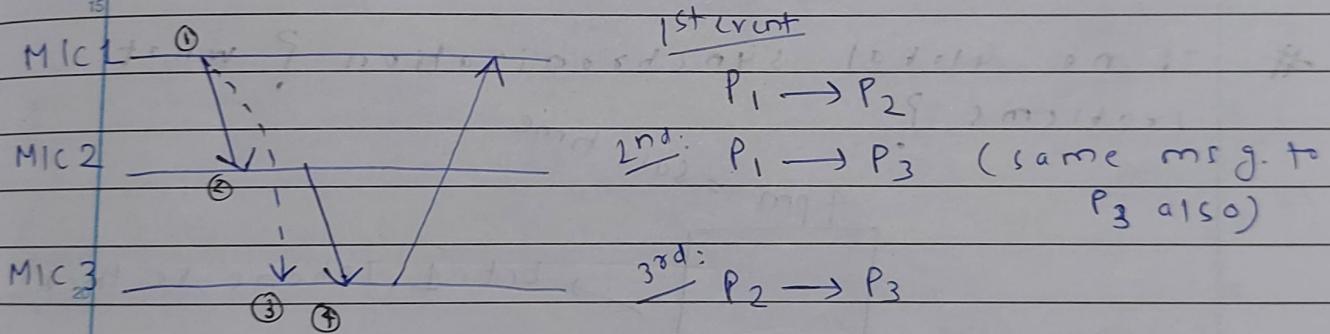
5

if we take 2nd parameters:

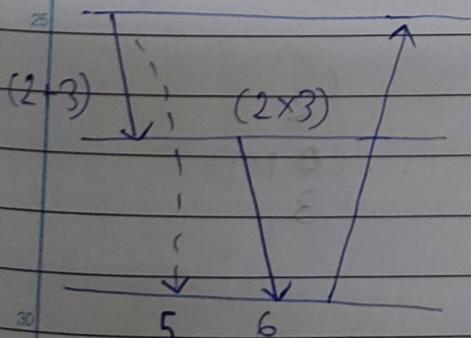
Conclusion = $\Sigma_1'' \rightarrow \Sigma_3$

if there is message passing, then only we need to know which event occurred before/after.

Problem with vector clock:



we expect in this order $① \rightarrow ② \rightarrow ③ \rightarrow ④$



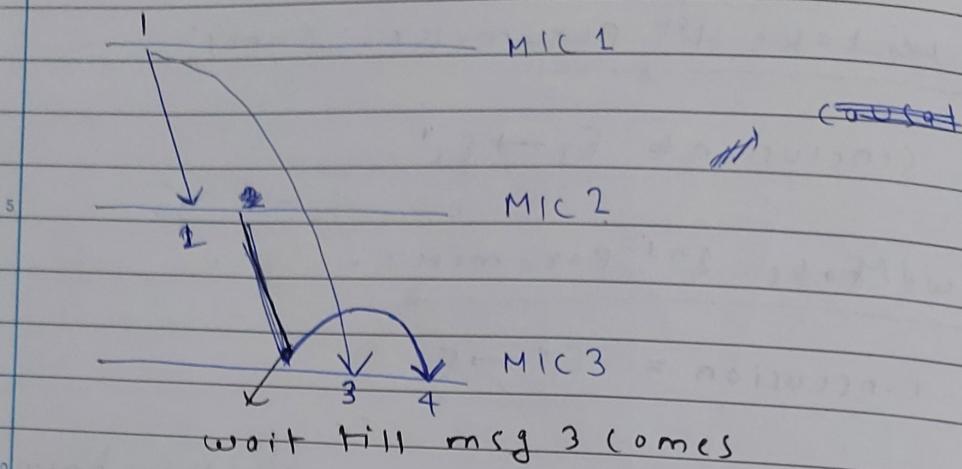
due to New delay, if msg 4 arrives before msg 3, there will be problems.

NOTE:

* Along with syncing clocks, the order of messages must also be ~~reserved~~ preserved.

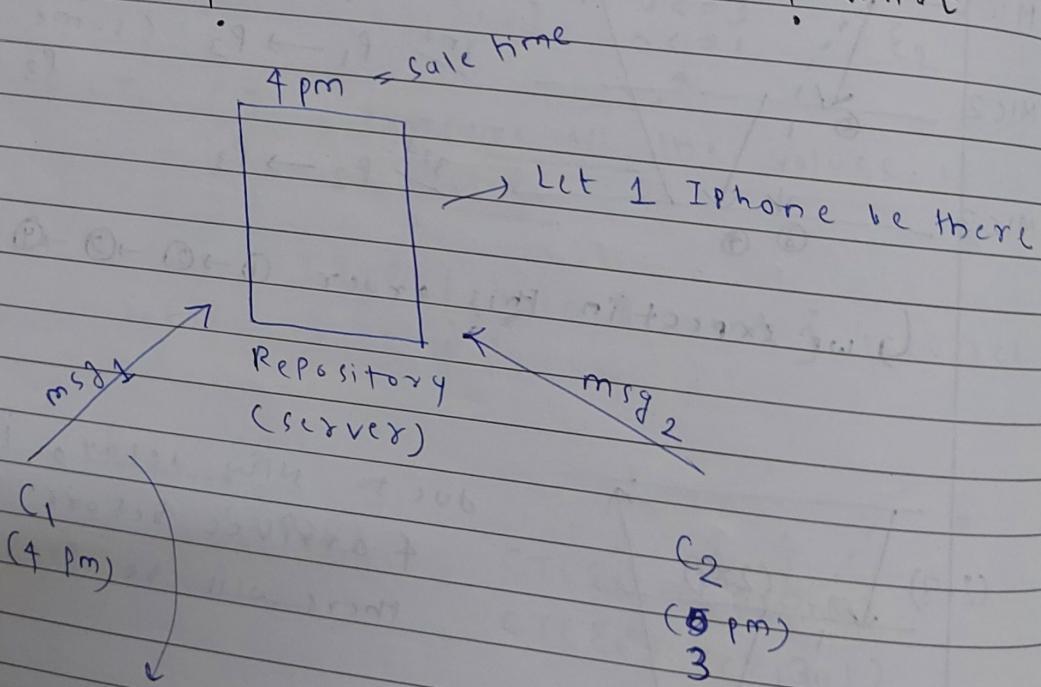
Vector clock = global clock.

Page:
Date:



Causal ordering = order of received msg.
is same as order of sending msg.

If no global synchronization? what
problems?



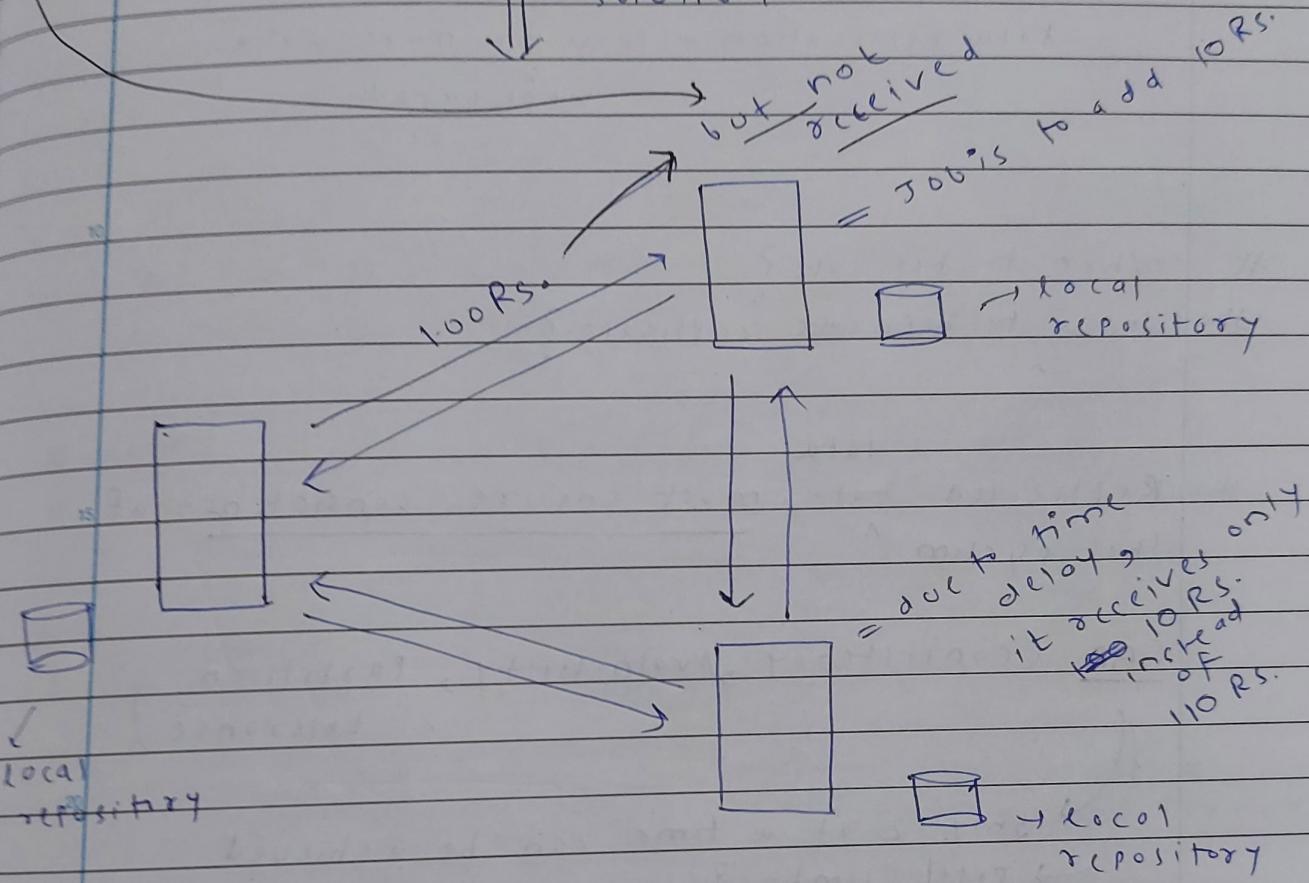
Let transmission time = 1.5 hr
(So C₂ get phone instead of C₁ due
to delay.)

- ① data loss = IF msg. lost
 ② data duplicacy = if ACK. lost, 3rd MIC receives 210 RS.

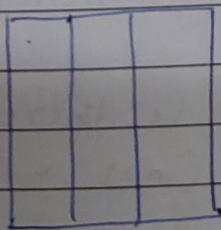
Page: _____
Date: _____

Inconsistency (if first come first serve), if no proper ordering of messages (no global clock synchronization)

↓
solution



* Logical meaning of snapshot:



combine to form whole Image
(no data loss, no redundant picture)

if snapshot is consistent

Take Global snapshot of MIC to avoid Camlin data loss & duplicate

Dec (31st October, 2022)

Page:
Date:

Replication & consistency

Replication → data replication

→ computation replication

✓ replicate computing power

data replication = few or more data replicated

When to replicate?

How to ~~replicate~~ replicate?

* Replicated data must ensure consistency of the system ^

CAP (Consistency, Availability, Partition tolerance)

only 2 at a time can be achieved

Total umt = 3 (3)

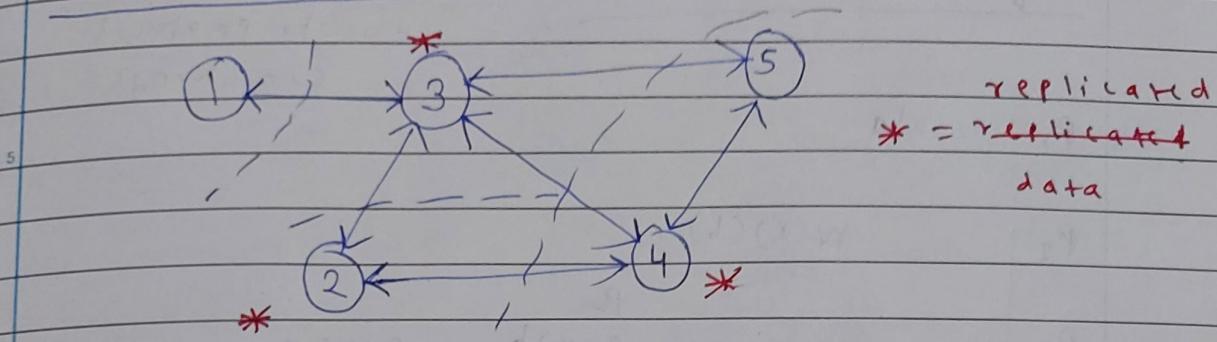
Consistency = machines must have the most updated data. everyone reading same data

Availability = update of data should be there at all machines at that same time (not wait for any other request)

Fault tolerance > scalability > replication

Page: _____
Date: _____

partition tolerance \Rightarrow partition fault tolerance



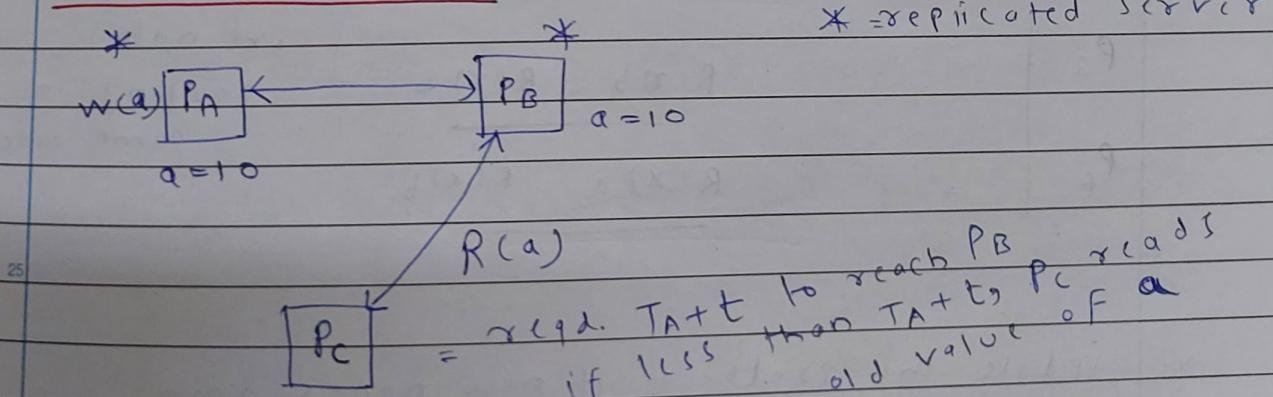
* = replicated data

10. if 3 logs off, graph divides into 2 parts
whether the system can continue the work
or not after 1 device logs off?

- * sync. with help of logical clocks
* computation stuck in 1-to-1 mapping

- * response of the server \Rightarrow depends on frequency
of server.

Strict consistency



- * let update at T_A & t time to reach replicated system

$T_A < T_B < T_A + t \Rightarrow$ this does not happen
in strict consistency

- * very hard to achieve in run time.

easy to achieve.

sequential consistency → weaker than strict consistency
due to protocol can change

P ₁	w(x)a
P ₂	w(x)(b)
P ₃	R(x) b R(x) a
P ₄	R(x) b R(x) b R(x) a

communicating with each other, and agreed upon a protocol in which order they should ~~read~~ execute (All of them would have same data)

P ₁	w(x)a
P ₂	w(x)b
P ₃	R(x)b R(x)a
P ₄	R(x)a R(x)b

} not allowed since P₃ & P₄ are not reading in same order

* Update the protocol if new nodes join

→ valid in same order
as write

Causal Consistency

* maintain same "happen-before" relationship

Linearization

Linearizability → between strict & sequential

P₁ P₂ P₃
x = 1 y = 1 z = 1
print(y, z) print(x, z) print(x, y)

x, y, z = 0 (Initially)

if this sequence:

* x = 1
 print(y, z) o/p: 0 0 1 0 1 1
 { y = 1 → —
 print(x, z)
 z = 1
 print(x, y)

Valid since we are reading either
not updated value / immediate updated value

* x = 1
 print(y, z) } Consistent
 y = 1
 z = 1
 print(x, z)
 print(x, y)

$TS(W) > TS(R)$
↳ not allowed

* $x = 1$

print(x, 2)

print(y, z)

5 $y = 1$

$z = 1$

print(x, y)

} not valid

(not consistent)

10

15

20

25

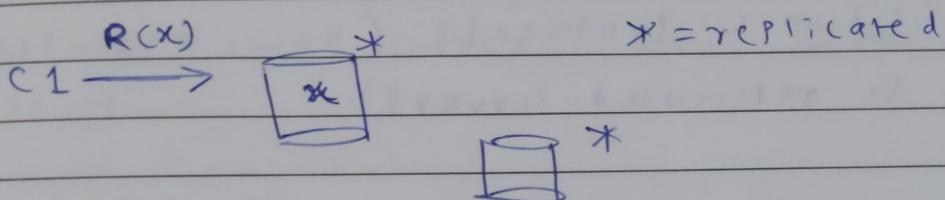
30

Dec (4th November, 2022)

Client-Centric Consistency Models

From
perspective of a single client

- * Monotonic reads \Rightarrow returns same/more recent data



a) $x = 10$

b) $x = 11$

older

c) if $x = 10$ again, older & newer to cannot be distinguished

- * Monotonic writes \Rightarrow replication of data moved in order to ~~serves~~ servers.

- * Read your writes \Rightarrow whenever we write / update data, read the new data always.

- * writes follows reads \Rightarrow

abc@ - time 0

xyz# - time 1

↑ monos

update done on most recent update (on time 1
not time 0)

assumption
~~assumption~~ that updates
will be done to replicated
servers eventually

triggered updates = updates only when one of replicated servers update
Page: _____ Date: _____

Periodic updates = updates after a certain period of time.

epidemic protocols

How to send update copy to other servers?



through epidemic protocols (How msg. spread out to replicated servers)

mostly all systems go for triggered based updates

Spreading an epidemic

Anti-entropy (since we are randomly picking any server)

Rumor mongering

#

Updation - pull = P updates, other server pulling
push = new update, sending update

pull-push

pull the new update and update to other server

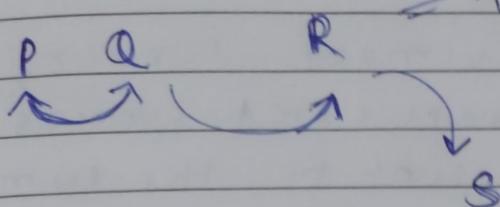
also push the new update

if only pull => one way communication, longer time to spread the epidemic
& push

push-pull (Pairwise update) = takes less time

(Gossiping)

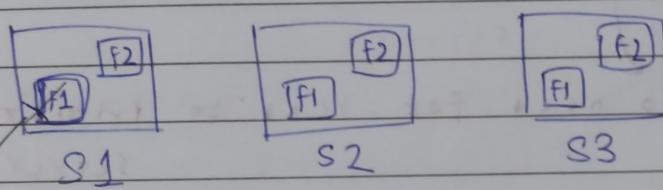
Rumour mongering $\Rightarrow N$ replicated servers,
 N^2 pairs, if any update among pairs, updates spread



has
'if already update, not
send update to any
other server
(spreading)

5 (Anti-entropy)
Faster than previous, but no guarantee that
all replicated servers have same/new update
update.

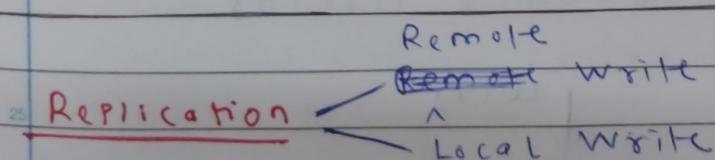
→ Removing Data



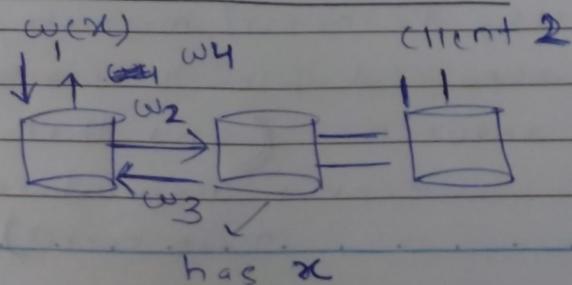
15 if deleted & at time takes to move to other server

SOLUTION = Death certificate issued.

20 each time, validate death certificate, therefore maintaining a log incurs overhead.



Remote write protocols



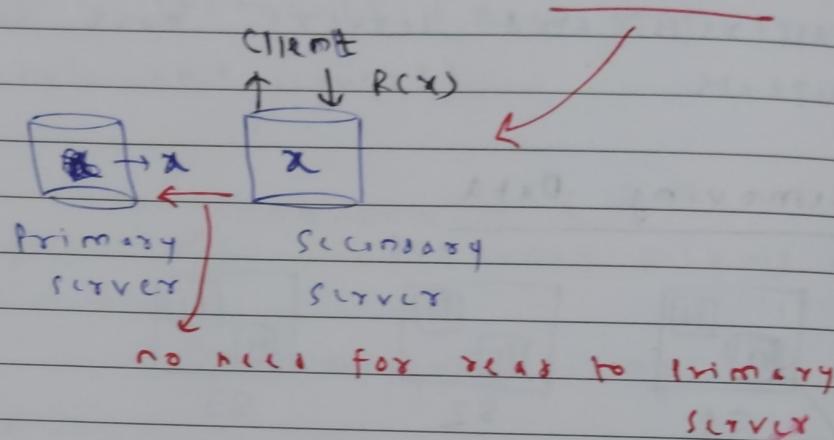
25 1 SERVER
= PRIMARY SERVER
other servers
= SECONDARY SERVER

Remote Write = update to every other servers
if update on Primary server

Page:

Date:

- 1st update on the primary server, then only updates to other replicated servers
- Acknowledgement is sent to the primary server
- Forward the request to the particular server since ~~only~~ all servers are replicated



Backup server = if primary server fails, it takes over charge.

Local Write protocols / no concept of primary server

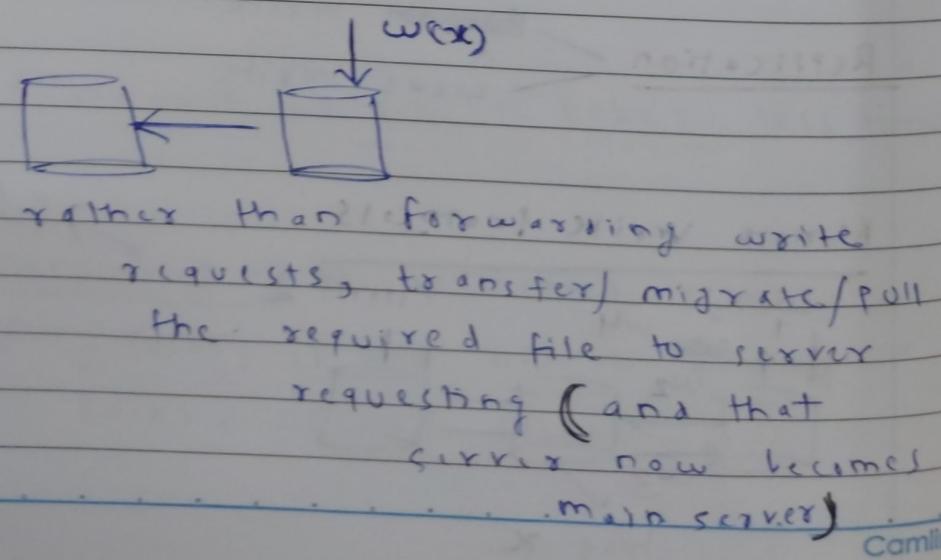
→ current/main server is

Migrating the file to other servers. There

Migrating local data

to local place where access rate is

high



Camilin

Problem = replication transparency must be maintained, much more write frequency by different clients, local write protocol not feasible. (go for remote write in this case)

10 Synchronous replication = wait for ACK from replicated servers whether sync done or not, coordinator or leader will wait for ACK

15 Advantage = if ACK have come, full surety that all replicated servers have updated

Disadvantage = waiting time more for one single operation.

20 Async replication = like ~~eventual~~ consistency, not wait for all ACK, operation successful if ~~no~~ ACK from majority of servers.

25 Advantage = must faster

Disadvantage = ① no guarantee all servers in consistent state.

② error when access to server from which no ACK

lec(5th Nov, 2022)

→ Fault tolerance

Fault less with help of replication

If 100 m/c instead of 10, prov. of Fault more

Probability of failures (depends on) no. of m/c attached to system

Single m/c → Failures all or nothing

Distributed systems → Partial failures

→ Dependable systems : → Maintainability
(can recover from failure without any loss)

availability (measured through 9's concept)
(System should be live at time of access)

4 9's concept = 99.99%
available, 5 9's
= 99.999%.

System available

High available system

↳ 5 9's = 99.999%
↳ For critical request-response
↳ like rocket launcher

Reliability

Safety (partially affected client should not harm the whole system)

→ Types of Faults:

- 1) Transient faults = failure occurred only once and it never comes back (even if comes back, freq. is less)
- 2) Intermittent faults = getting error in time intervals.
- 3) Permanent failures = one major fault, entire system crashes (not able to recover). → much more easy to detect & maintain
 → major fault (since 0/1 = either working or not working)

① & ② = hard to detect & maintain

→ Failure Models:

→ Rash Failure ⇒ A server halts without any notice & no chance of recovery, easy to detect and maintain.

→ Omission Failure ⇒
 receive omission
 send omission
 handles response request from client

send = not able to send msg.

receive = not able to receive

→ Response Failure → Value failure (expected result not valid, like garbage value)

5 State transition failure,
(order of responses
not right)

→ Arbitrary failure ⇒ not able to get response
10 i.e. let time sometimes it timeout

Failure, response failure, system is confused
which fault tolerance is required.

Byzantine fault

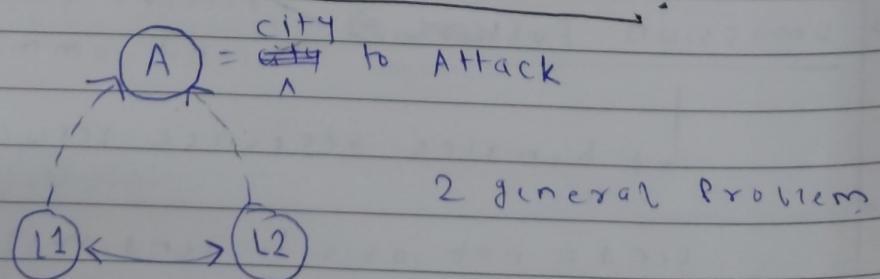
15 → Agreement in faulty systems:

* At least $k+1$ replicated systems for fault tolerant system

20 * ~~k-fault~~ tolerant systems

* How to detect erroneous machine?

25 → How to detect Byzantine failure?

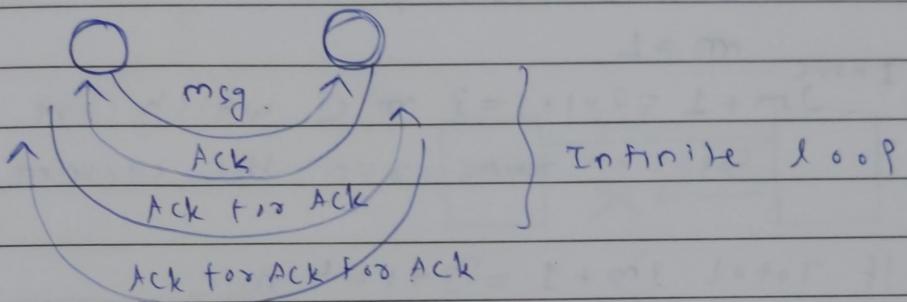


30 L1, L2 = generals... Let's attack at 9...
ACK for every request
or at 10

→ communication channel = eavesdropping may

be there, security issues, due to which expected results not there.

→ infinite loop of final acknowledgement of the acknowledgement



Byzantine general problem

→ Iterative ~~algo~~ algorithm

→ Faulty machine sends different msg. (garbage values)

→ may give correct result / wrong result

a, b, c, d = arbitrary msg. from m/c 3

→ comparison of results (rows) of all m/c, if they are different, that m/c is faulty

(\equiv) (\equiv) (\equiv)

→ n machines, k machines may be failure
 k subset of n , $k \in [0, n]$

→ If mIC come more, detection of fault is difficult due to large no. of comparisons.

5 → $3m+1$ process & $2m+1$ agreement

Ex: 4 mIC

Let 1 mIC = faulty

10 if I have $m=1$
 $2m+1 = 2 \times 1 + 1 = 3$ mIC which are perfect,
we can ~~run~~ run the system

If Total $3m+1 = 3+1 = 4$ mIC

15 Out of 4 mIC, 3 are working fine, then
system works well

20

25

Dec (11th Nov, 2022)

Page :

Date :

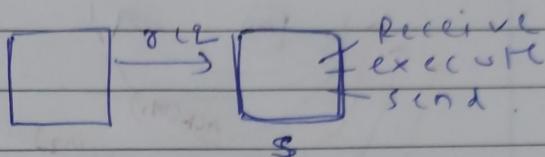
Agreement

2PC

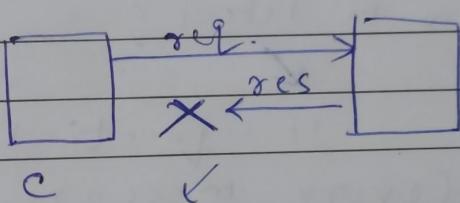
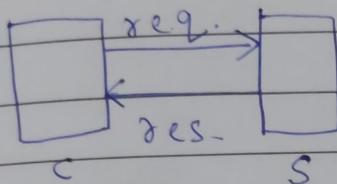
Reliable

3PC

Comm.



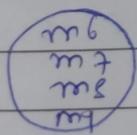
(if req lost, timeout
can apply)



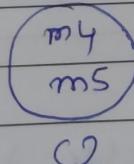
response not reach
client

Faulty unicast
system

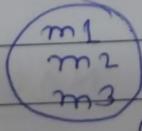
(concurrency)



C1



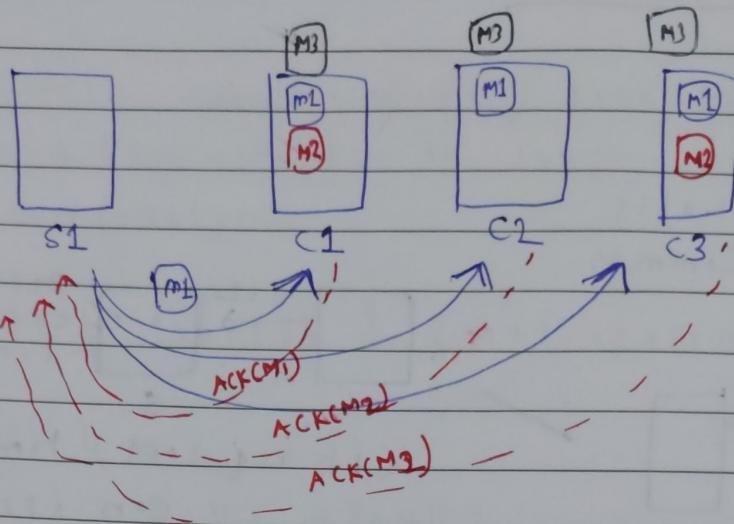
C2



C3

ACK

NACK



How to detect?



if sequential number

(every message has a unique sequence number, and if it does not find seq.no, then it sends NACK)

Acknowledgment explosion = Scalability Problem



Problem of ACK based Protocol

No. of messages is less in NACK based Protocol



only m/c which is in Failure, sends NACK

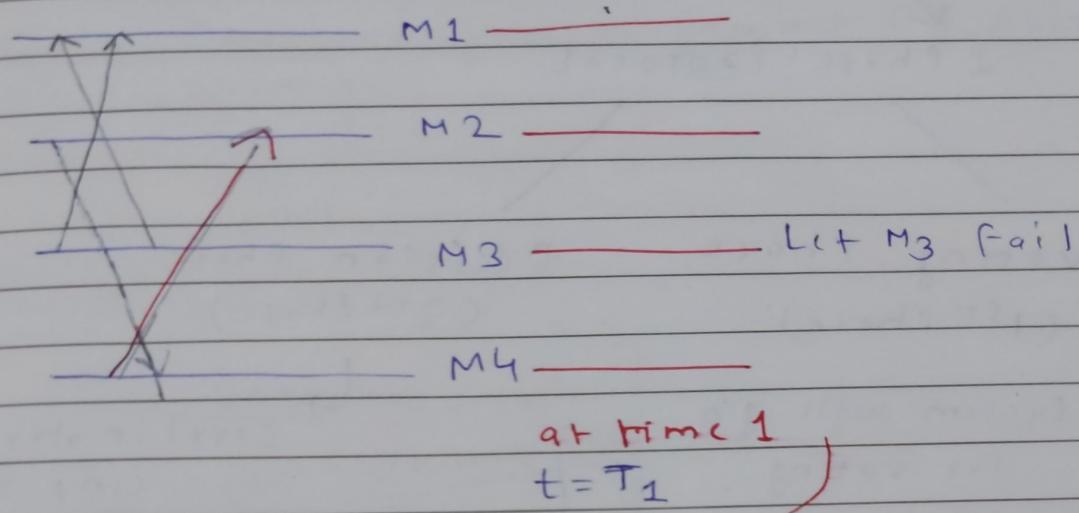
Atomicity Property = either whole transaction occurs or nothing occurs

Page :

Date :

analogously same as Global Snapshot

* Group view Semantics



Initially machines form a group, Let
group = {1, 2, 3, 4}

- regrouping {let regroup = {1, 2, 4}}

In case
of
fault
abort

at time
at ~~t₁~~ = t₂

M₁

M₂

M₃

M₄

(either all are

aborted or again a

regrouping (through agreement protocol
with other live m/c))

Distributed Commit

↙
2 Phase Protocol

voting Phase
(1st phase)

Decision Phase
(2nd Phase)

{ System will go
for voting
whether all
agree to
commit, regroup
or abort

↳ Coordinator will
send the
global decision
to
subordinates
whether to
commit/
abort

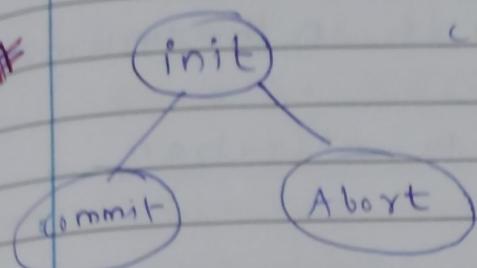
even if 1 machine
votes for abort, whole
system aborts.

COORDINATOR = 1 m/c

Subordinates = $(n-1)$ m/c ↙ if n m/c

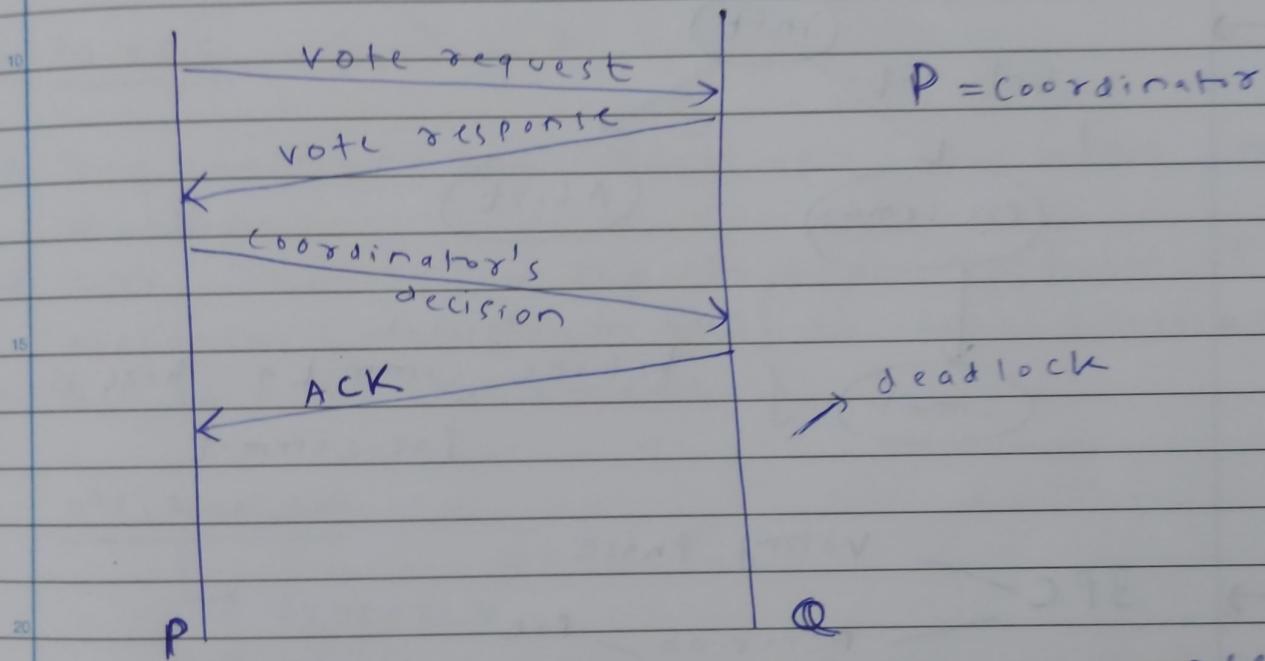
25) ↗ These selected to leader selection algorithm
according

All or none (Atomic)



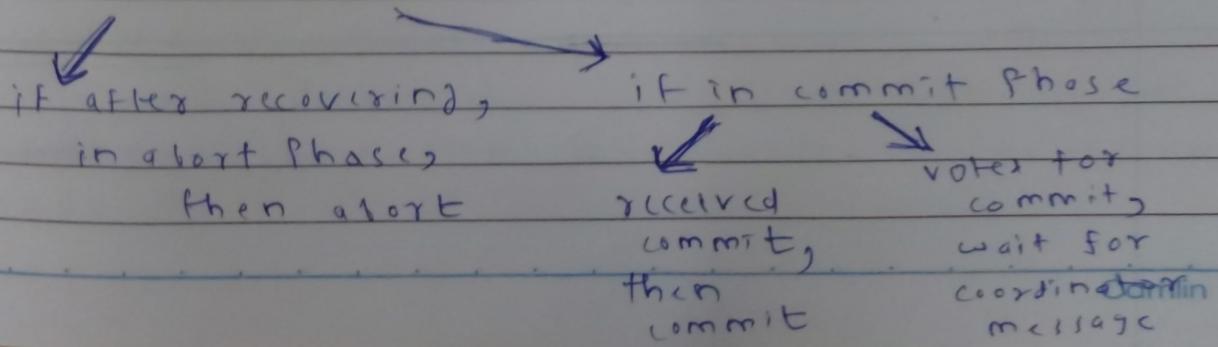
coordinator asks for vote req. from subordinates, then it goes to Commit OR Abort.

Problem with 2PC:



- * After sending the ~~receiving~~ vote ~~req~~, coordinator fails (single point of failure)
- * If Q sends abort, then whole mc abort
if Q sends commit, whole system deadlock

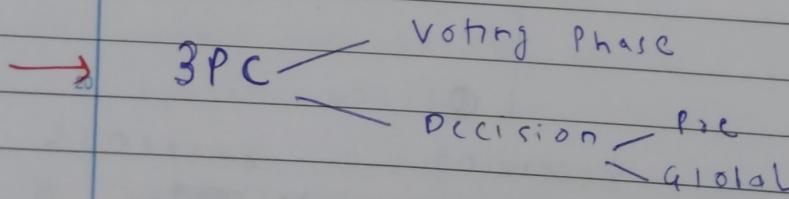
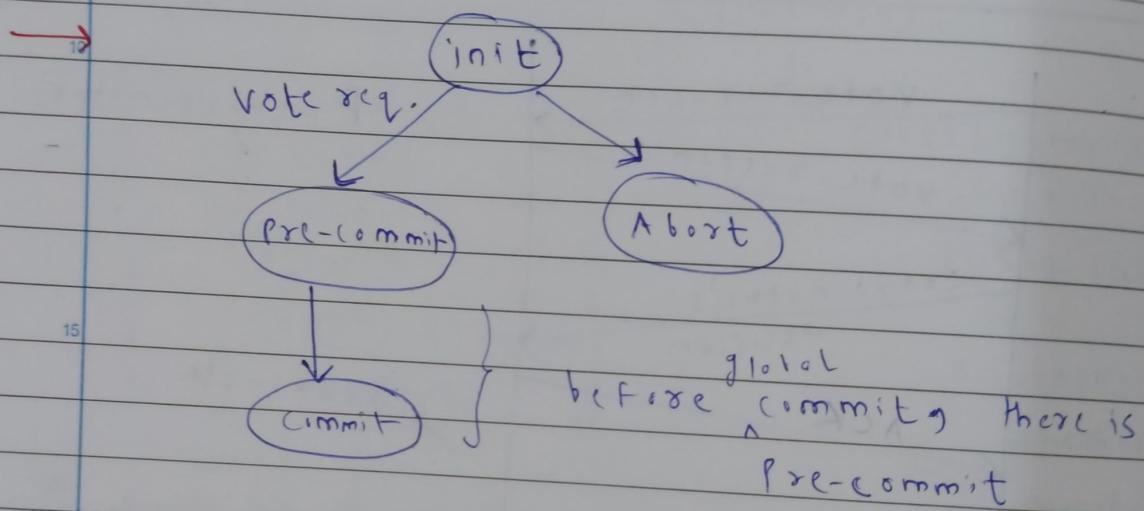
* if Q crashed



Safety = ensures no harm is in the system

Liveness = no progress due to operation

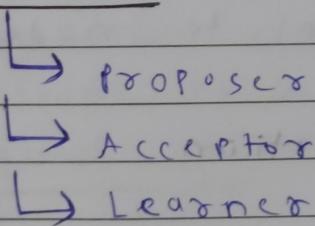
3-Phase Commit (3PC)



- if everyone sends ACK in Pre-commit,
then system goes for ~~commit~~ global commit.
- if 1 commit & N-1 precommits, then
broadcast message of commit

- In Byzantine fault, agreement cannot be done by 2PC/3PC.
- Byzantine Fault handling harder to implement.
- Crash Failure easy to implement.
- Paxos → Iterative Algorithm
- They do not take decisions from all the machines
- Here, the leader try to take majority of response (not from all) to take decisions.
- Not easy to implement

Mechanism



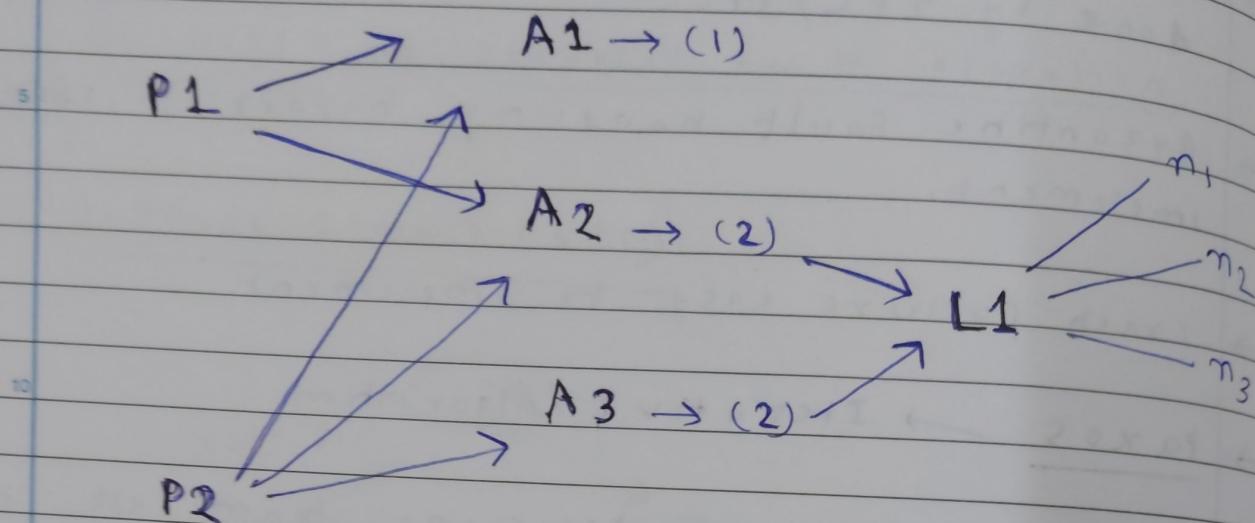
* Proposer \Rightarrow will propose it wants to update / write some value that I am proposing that these will be the new value. Can be single / multiple Proposers.

* Acceptor \Rightarrow ~~AFTER~~ accepting can accept / reject proposal. can send ACK / NACK. ~~AFTER~~ send the response to Learner node.

Not necessarily, proposals will be sent to all Acceptors.

Page:
Date:

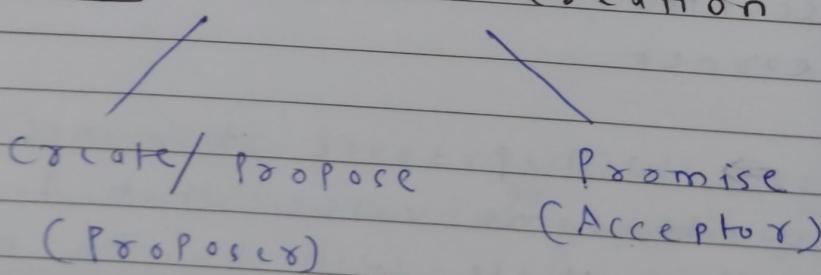
Let $N = \text{nodes}$.



(1) = Accepted from P1
(2) = Accepted from P2

System looks for majority of proposals.
Hence (2) gets accepted and sent to
leader and further sent to nodes.

Phase 1: proposal creation



- * sends a delimiter kind of msg. (sent number to other node that it has a proposal)

- * sends ' n_i ' from i^{th} Proposer (n_i)

* Accept / Reject the
 n_i

* Check if n_i is greater
than ~~previo~~ number
of previous proposal

Ex: $P_1 \rightarrow n_1$

$A_1 \rightarrow$ check if $n_1 > n_2$

↙ $\underbrace{\quad}_{\text{Previous Proposal}}$

if true, then proposal accepted
else rejected

* if accepted, then it sends a ACKNOWLEDGEMENT (Promise message). After receiving promise message, it sends the original message.

↳ (n, v)
↙ \downarrow value

↳ proposal
number

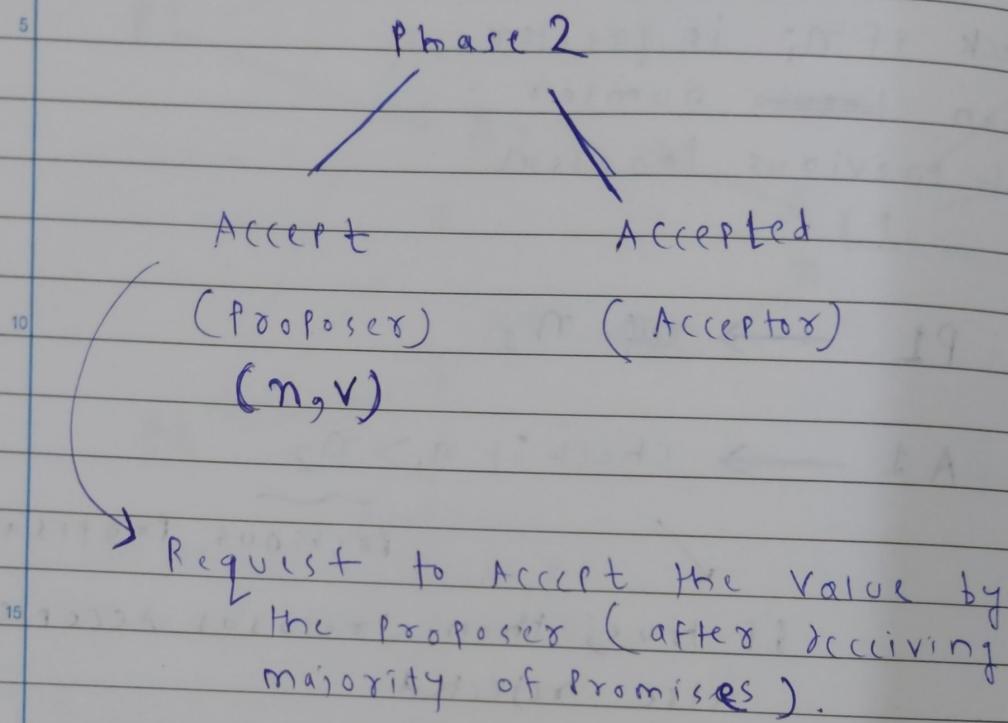
* if rejected, it sends NACK along with reason of not accepting.

Validity
Reliability } By sending ACK/NACK
Safety }

Check Bounce (when not enough money
in bank account) \Rightarrow Paxos

Page: _____
Date: _____

Phase 2 : After phase 1, how database
updated



- * Accepted = final / termination phase,
receives the value and updates
the database with the value accepted.
- * if proposer fails, Acceptor waits for
some timeout period since it is iterative
- * if "no majority", all the proposers restart
with new sequence numbers and according
to ~~Paxos~~ Probability theory, 2 or 3 iterations
the majority of Acceptors are there.
Proposals are accepted.

can also be used in "leader selection" problem not only as a "Fault tolerant" mechanism.

We count the no. of Acceptors where most of promises are sent.

$$\underline{t = T_1}$$

$P_1 \quad n_1$

$P_2 \quad n_2$

If P_1 has most no. of promises, then it is the leader, in case of Tie-breaker, start from the next iteration.

every replica has most recent update from the same proposer (not necessarily, most recent update from system).

No conflict of instruction

Majority can learn updatations from proposer

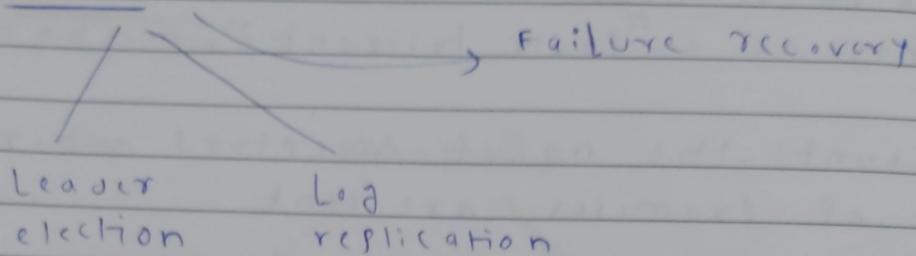
Validity, Liveness ~~is~~ is there
failure

Crash Failure / Fail stop ~~failure~~

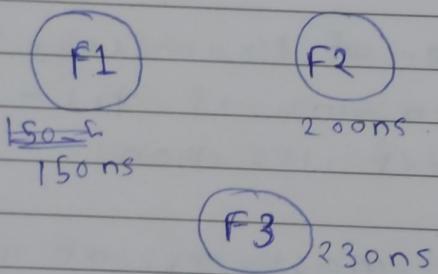
lec (16th Nov, 2022)

Page :
Date :

RAFT = 2 Phases

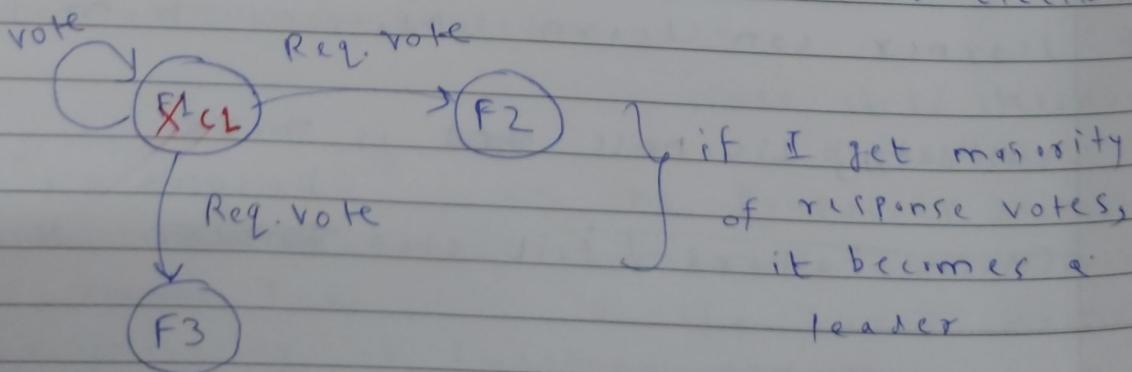


Leader candidate Initially all m/c are
follower in follower state



150-500ns.

→ After a timeout expires for 1st candidate, it becomes a candidate for leader election



→ if 2 m/c have same timeout, both will become candidate for leader selection.

- whichever mc have high votes, that will be leader
- if tie, start the iteration again as all as follower and timeout may be different.
- split vote
 - leader election gets activated when
 - Split vote
 - Leader selected in 1st iteration gets crashed

- L
- F2
 - * Followers neither execute any request
 - * not ask for any vote
- F3
 - * Just follow the leader as per the instruction sent by the leader.

- RAFT use RPC
 - to request vote
 - to update log table
- Each server maintains a file, containing sequential requests that come up to that server.
This is called log file.

Adv. of RAFT = decompose failure recovery
into subproblem

Page:
Date:

Index	1	2	3	4	5	6	7	8	Term No.
Term NO.	x	y	z	x	y	z	x		
Value/ data	4	1	3	0	2	0	6		

Let this be
leader

missing value or fill it
→ Use log Table ← New update need to be
made

	1	2	3	4	5	
Case 1:	F1				y ← 2	→ revert back
Case 2:	F2	6	7	8	9	10

	6	7	8	9	10	→ revert back
	*	2 ←	x ←	x ←	y ← 3	hit match append rule

if term value greater than term no.
follower has, then accept log file.

→ In case 2, revert back to 6, append the

6	+	8	9	10
			↑↑↑	

list of log file from Leader

"Strict write only" — If new leader,
new updated values need to be made

Quiz 2 - 3.2, 3.3 till RAFT (6.1) (3.3 - 6.1)

Page:

Date:

end term = whole syllabus

TILL RAFT = end term.

- Log Replication uses 1 RPC call = "Append Entry"
- When updation, each follower will send an ACK
- Leader will send COMMIT msg. of the last (will wait for majority of ACK from followers, send the termination msg.)
→ Check the liveness of leader
- Heartbeat msg. — Used to check whether still any leader present, or need to elect a leader, Heartbeat msg. sent periodically.
Leader will select a time period after which it will send Heartbeat msg.
- RAFT ensures $\begin{cases} \text{Safety} & (\text{always progress even when leader crashes}) \\ \text{Liveness} & \end{cases}$
- Safety = ~~No other leader can change the index in log file, only 1 leader each time.~~
- Implementation also easy
faster Algorithm than Paxos.
- Recovery is more easier
(crash) fault detection easier
Implementation also easy
faster Algorithm than Paxos.