

Distributed Systems

"autonomous" — should be able to perform some tasks independently

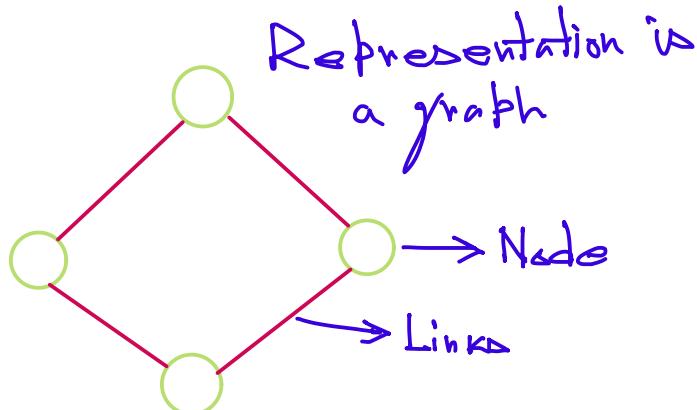
"communicate" — So, your PC Lab is not a dist system

"some task" — achieve something "global"

SIMD

ILM's classification

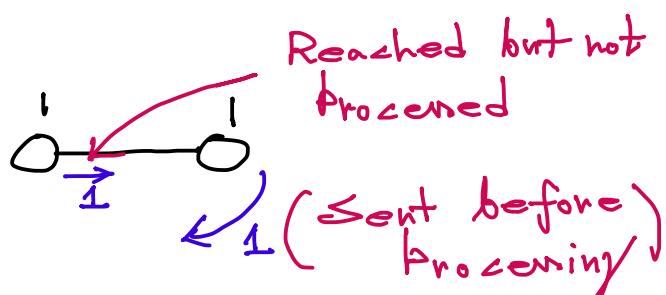
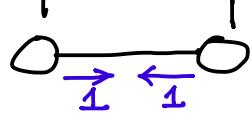
Parallel
Concurrent
Distributed



Fault-tolerant

Replicated

Atomicity losses!



Central!

- ▷ How to elect leader?
- ⇒ How everyone knows who is the leader?
(Each non-leader node)
- ⇒ Leader failure?
- ↳ Structure of the network (How to reach the leader?)

1. Dist Algorithms

2. Fault

3. Replication

4. "Other"



Every algo has some assumptions of the system

Asynchronous:-

- Makes no assumption on timing anywhere in the system.
- No assumption on max msg delay
- No assumption on processing delay

Synchronous:- Makes some assumption on timing

Typically/

- Assumes max msg delay
- Assumes processing delay

$$\geq d \times \alpha + \beta$$

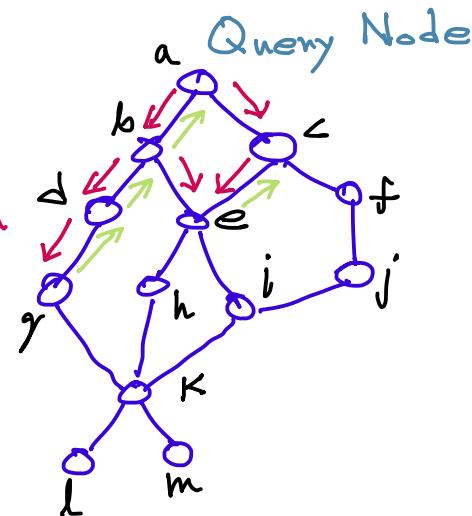
↑
diameter

→ choose as parent the first node you get the query msg from

▷ On receiving a msg

a) First: Send <CHILD> to parent and forward to else

b) Else: Send REJECT



Q) How do we write a distributed algorithm?

Main Idea: Identify events, write event handlers

So ① In initiation (if needed)

② Identify all messages you will use, including the content

For each process:-

- Write the send rule for each msg (When will a msg of that type be sent?)
- Write a receive rule for each msg (What does the process do on receiving a msg)

Add other types as needed! -

- 1) On a timeout
- 2) On a fault
- 3) etc

Dijkstra's guarded commands!

do

guard 1 → action 1
(Condition) (What is done when the condition is true)

guard 2 → action 2

end

Note! Both these ways but no constraint on the order of things done if (i) More than one event occurs at the same time (More than one guard true) and (ii) an event happens in the middle of handling of another event.

Q) Is the real world sync or async?

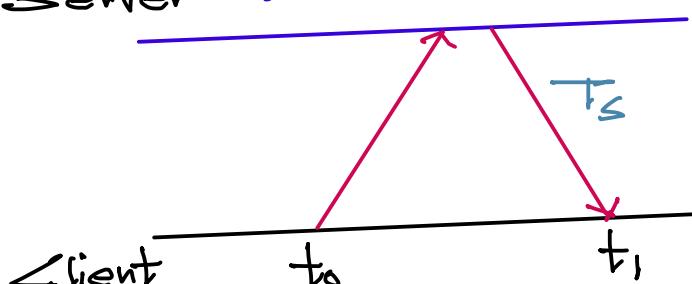
⇒ Async → But in practice, we make it sync with timeouts.

Time in DisSys

Cristian's Args:



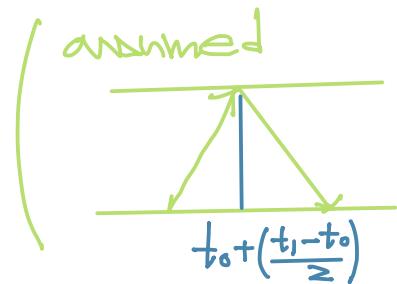
Server time →

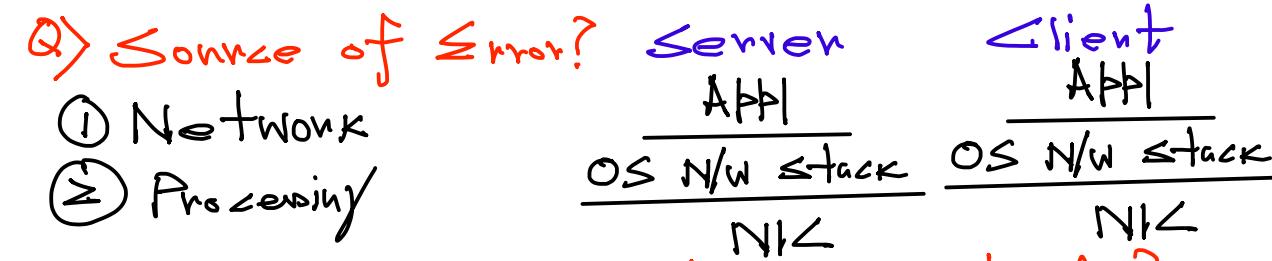


Max Error Possible?

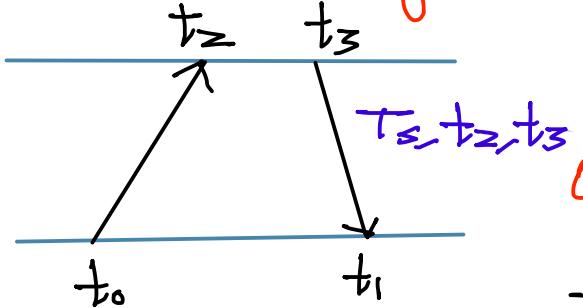
$$\left| \frac{(t_1 - t_0)}{2} \right|$$

$$\text{Set } T_L = T_S + \left(\frac{t_1 - t_0}{2} \right)$$





Q) Would it help if we also know t_2, t_3 ?



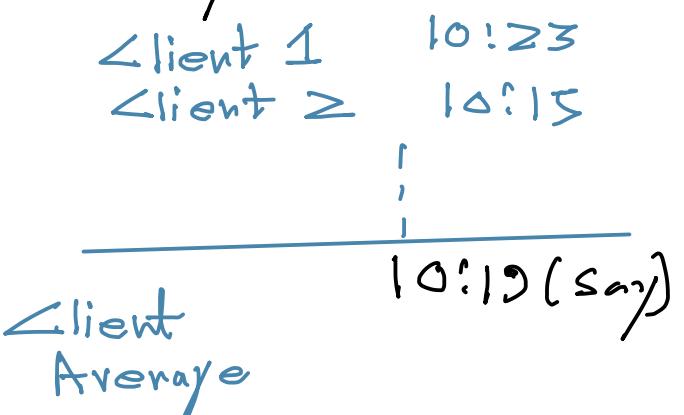
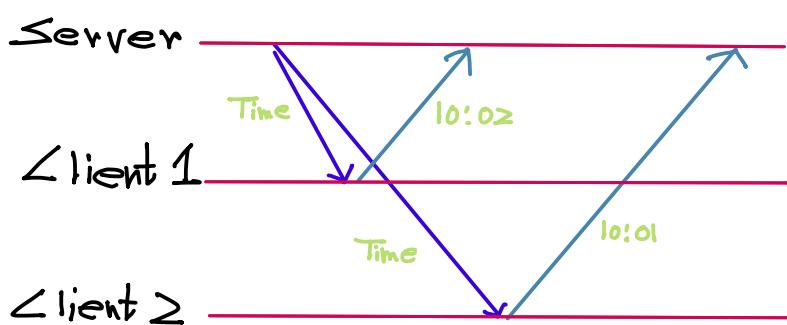
helps!! *(In remove processing time)*

Q) But where does the server timestamp t_2, t_3 ?

→ If at n/w layer, not helpful

→ So, Need to timestamp t_2, t_3 as close to N/W as possible (Use NIC Timestamping)

- Server adjusts for message delay (Same as Christian's)
- Subtles after message delay



Server sends
 Client 1 -0.04
 Client 2 +0.04

Problem in advancing clock immediately (in some case):
 May miss events scheduled to happen in the past jumped.

Pool.ntp.org
 NTP pool

Good estimation of delay

NTP

- ① Use hierarchy to scale
- ② Ask more than one time server.
 → Compute my delay offset
 etc for each, choose a subset (clock filtering algorithm)

— from the subset, select another
subset to eliminate potentially bad clocks
(clock selection algorithm)

good
clocks

— choose one final one based on ??

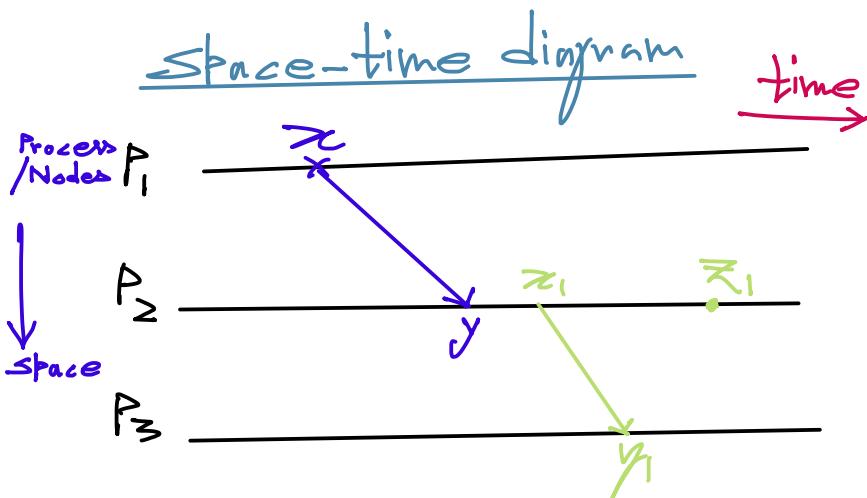
But not a total order

for ex, can't say

$y_1 \rightarrow z_1$

or $z_1 \rightarrow y_1$

y_1 & z_1 are called
"concurrent"
Events



$$\begin{aligned} &x \rightarrow y \\ &x \rightarrow z_1 \\ &x \rightarrow \bar{z}_1 \\ &x_1 \rightarrow y_1 \\ &y_1 \rightarrow z_2 \\ &y \rightarrow z_1 \\ &x \rightarrow z_1 \end{aligned}$$

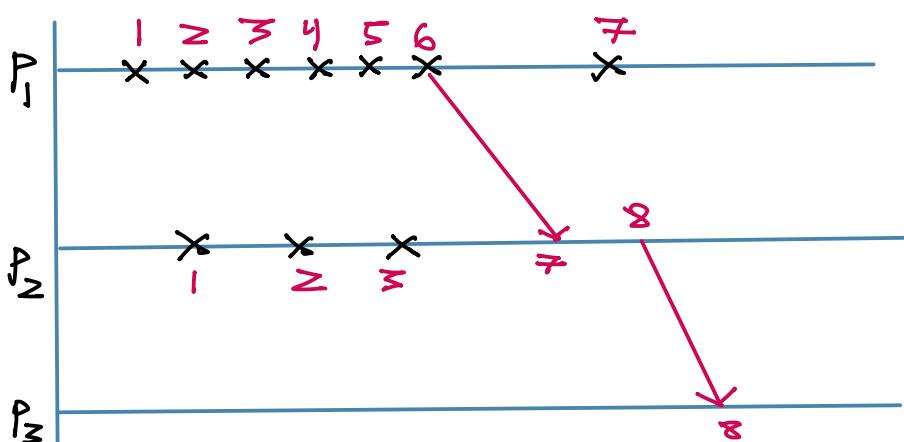
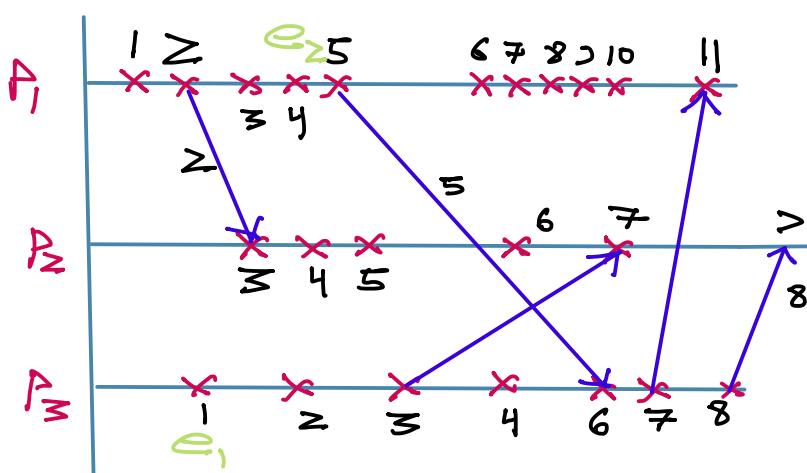
Events

Send (Send Event) of a msg

Receive (Receive Event) of a msg

Anything else (local events)

X-events



$$\begin{array}{l} \angle_1 \nrightarrow \angle_2 \\ \angle_2 \nrightarrow \angle_1 \end{array} \quad \angle_1 \parallel \angle_2$$

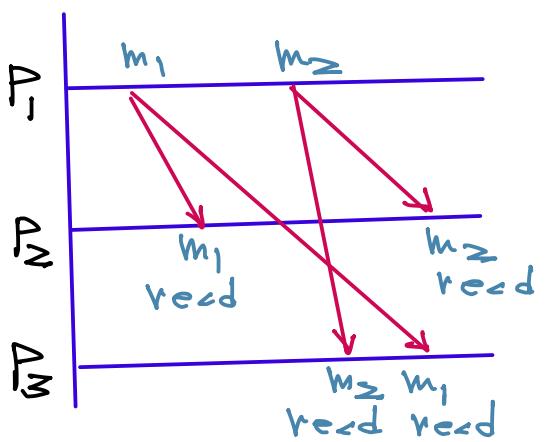
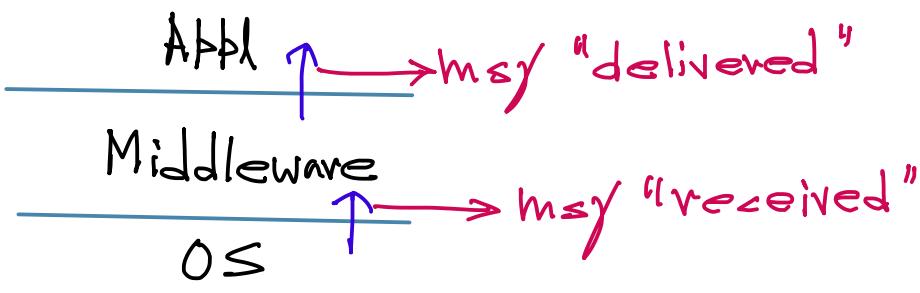
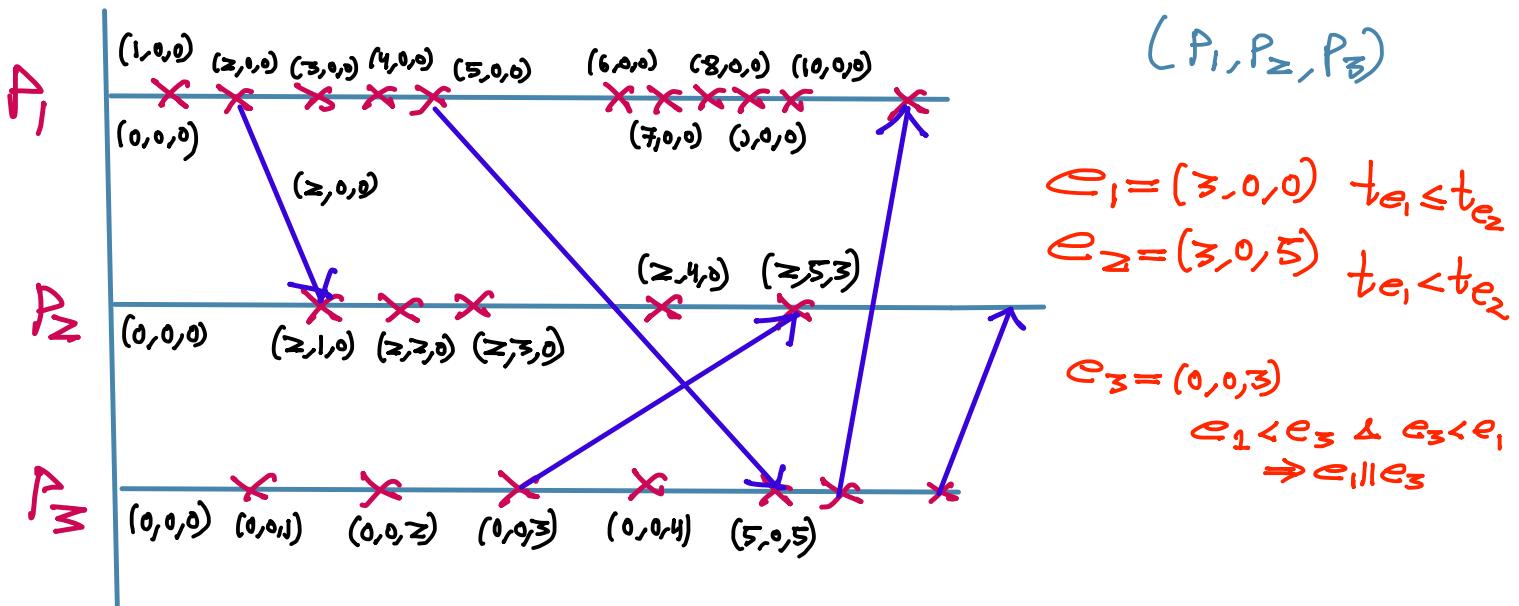
\Rightarrow To make a total order,

timestamp = (t, P_S)
of e_1

timestamp = $(4, P_1)$
of e_2

Vector Clock!

$\angle_i[k]$ = The no of events at P_k (Processor k) that i knows of.



Atomic Order
= Total Order

Will not be delivered to until m_1 is rec'd

FIFO, Causal, Atomic — Base

Can have Atomic FIFO, Causal Atomic

Q) What "event" are important for the "clock"?

→ Depends on the application

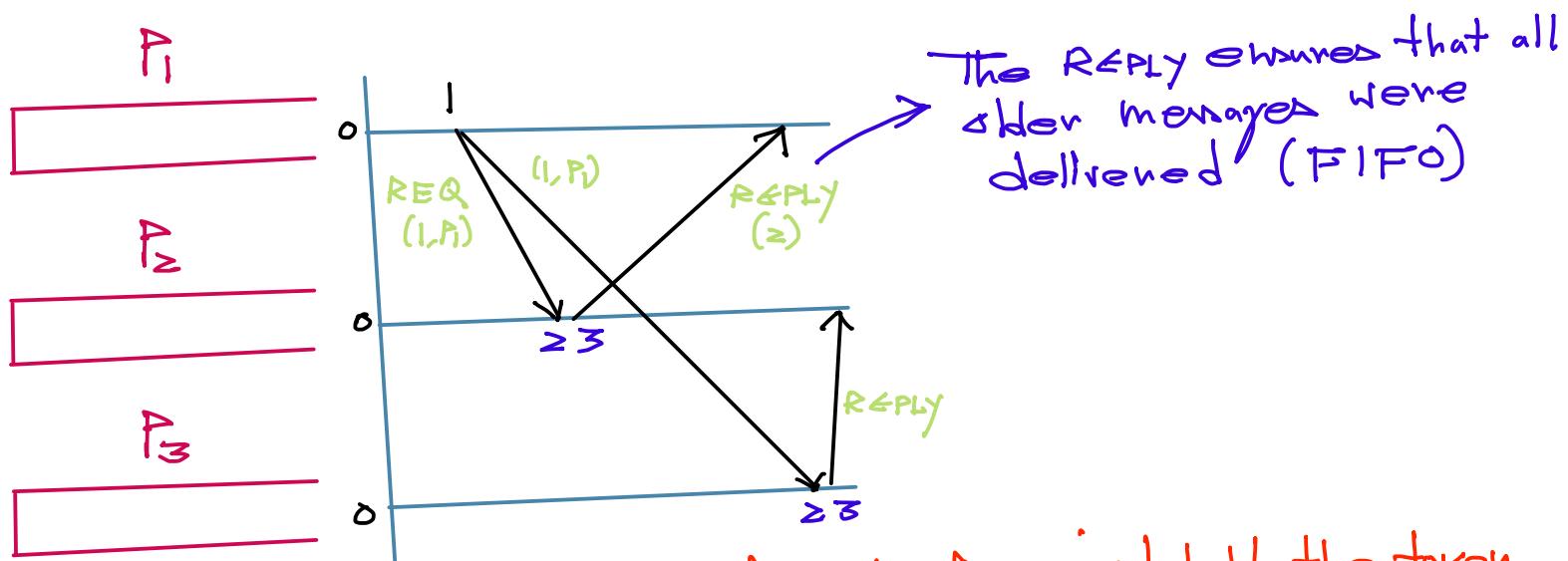
Q) Show Lamport clock does not work for this broadcast

Metrics for evaluating a dist' mutex algo

1. # of msgs per CS entry req served

≥. # Synchronization delay (Gap in time between the process leaving CS & the next waiting process entering CS) → Sync System (Otherwise no Guarantee)

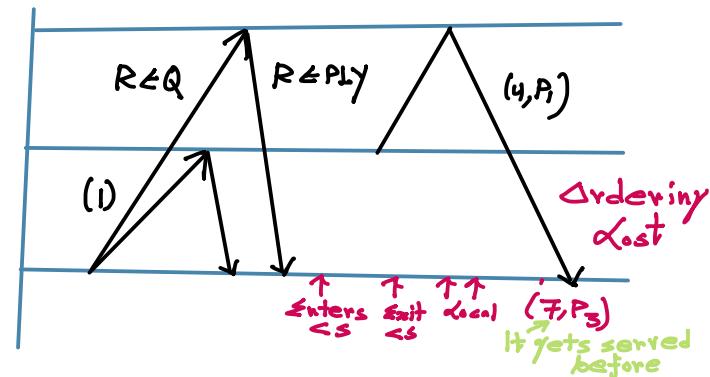
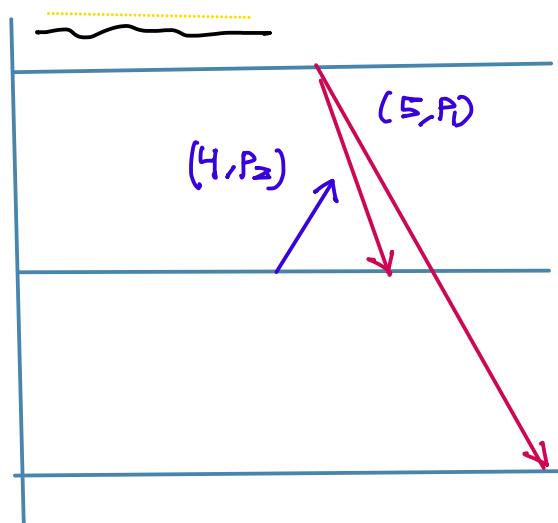
3. Response Time

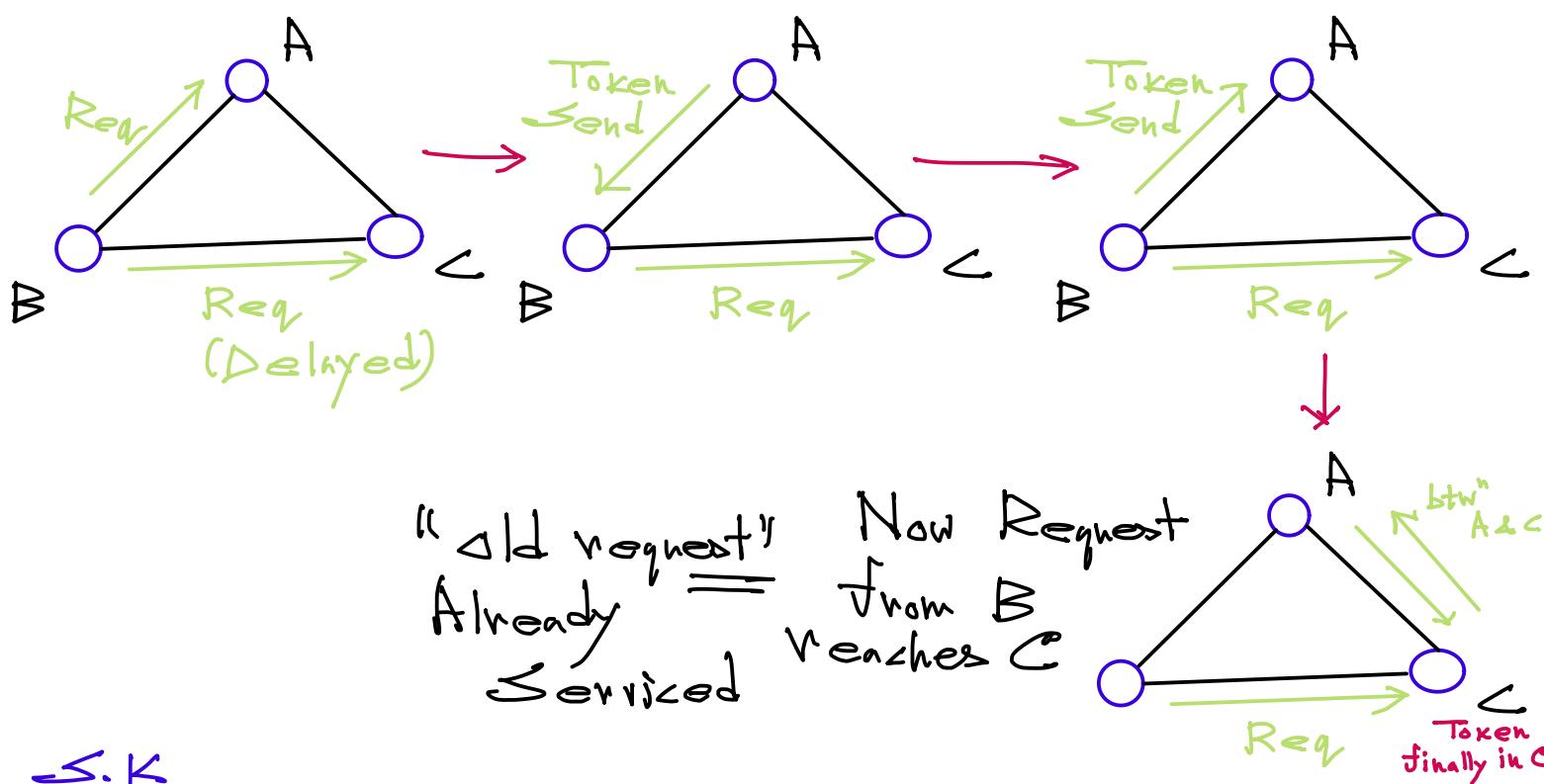


Q) What if I just hold the token if I will be making consecutive requests?

→ Adaptive

Raizinal
Carvalho





S.K

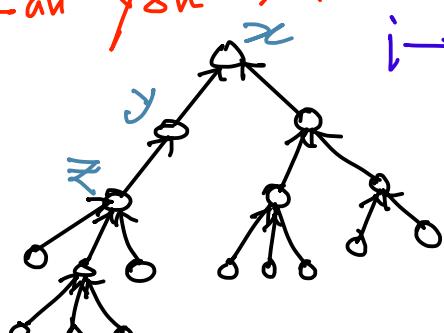
2 types of messages

- ① Request
- ② Token

Token Idle! Node not having token
and no pending request in queue

Assume, the node can make one request only at a time.

Q) Can you build a token-based mutex algo along this line?

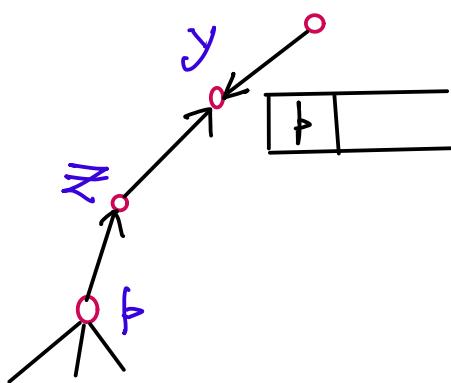


Naive algo!

- Keep token always at root
- Ask root to send token if needed
- Return to root once down

Better!

- ① Make actual root floating (but keep the defⁿ root = node with token)
- ② Let each node "take charge" for all nodes in the subtree rooted at it (Keep queue Get token etc.)



At the same time token reaches
y, y makes a request R

Possibility 1 :-

- Token is handled before R
- token sent towards f
- To token(y) = z
- put R in Q
- Send a Req for R to z

Possibility 2 :-

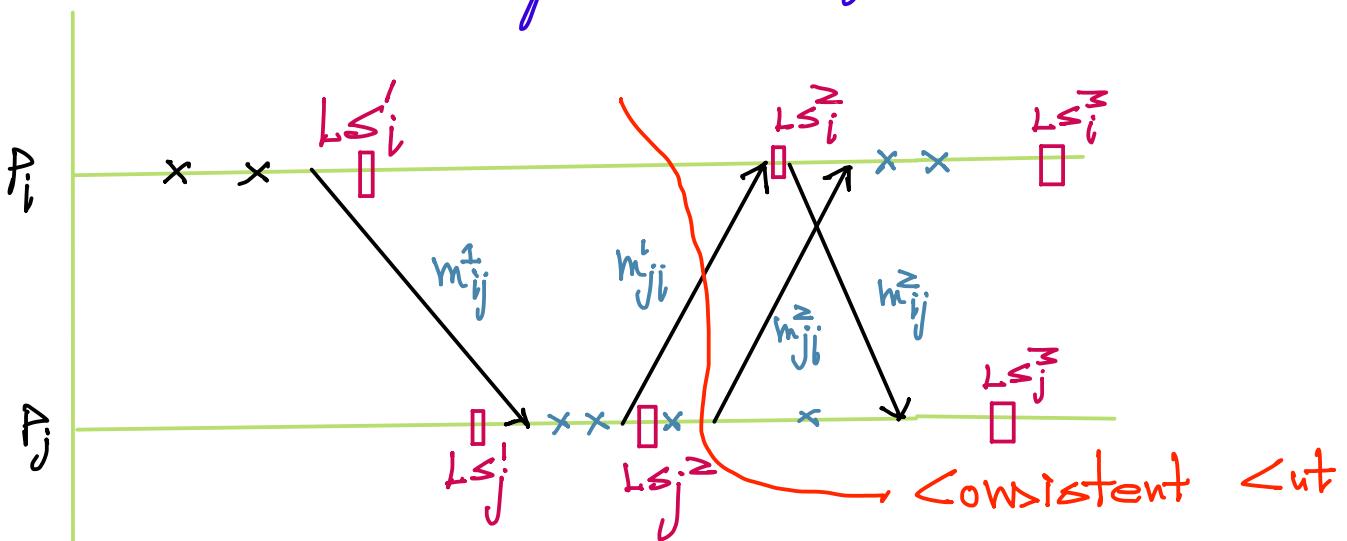
R is handled before token

Q) Ensure all sort of atomicity

Stable property = a property (predicate)

that once true, stay true until some external intervention happens.

For token ring, stable property: Exactly one token
But may keep on going



$$\text{transit}(Ls_i^1, Ls_j^2) = \emptyset$$

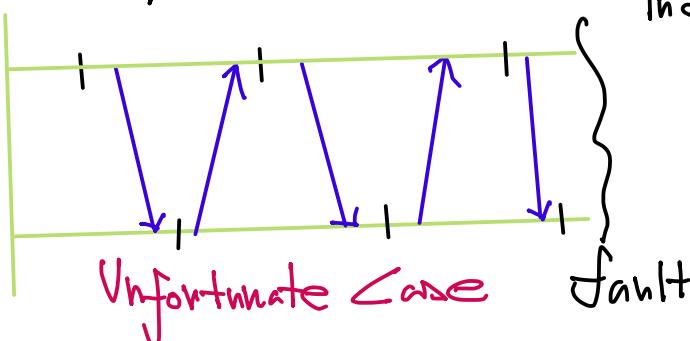
$$\text{transit}(Ls_i^1, Ls_j^1) = \{m_{ij}^1\}$$

$$\text{inconsistent}(Ls_i^2, Ls_j^2) = \emptyset$$

$$\text{inconsistent}(Ls_i^2, Ls_j^3) = \{m_{ij}^3\}$$

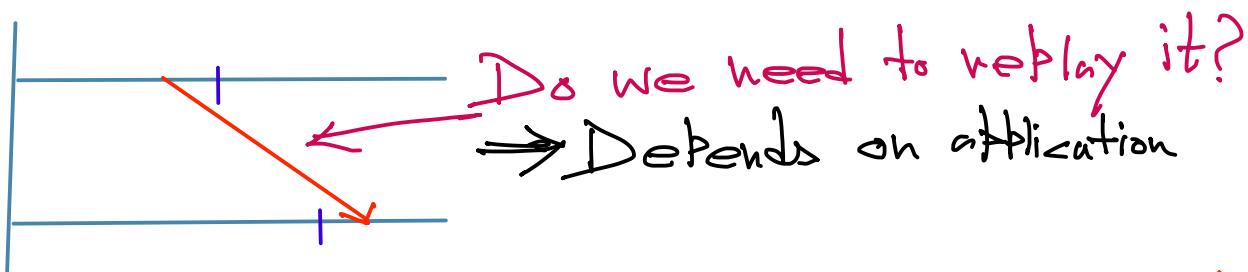
Q) Is consistent cut useful?

→ "Replay" from the cut.



Independent / Async
Checkpointing

Domino
Effect

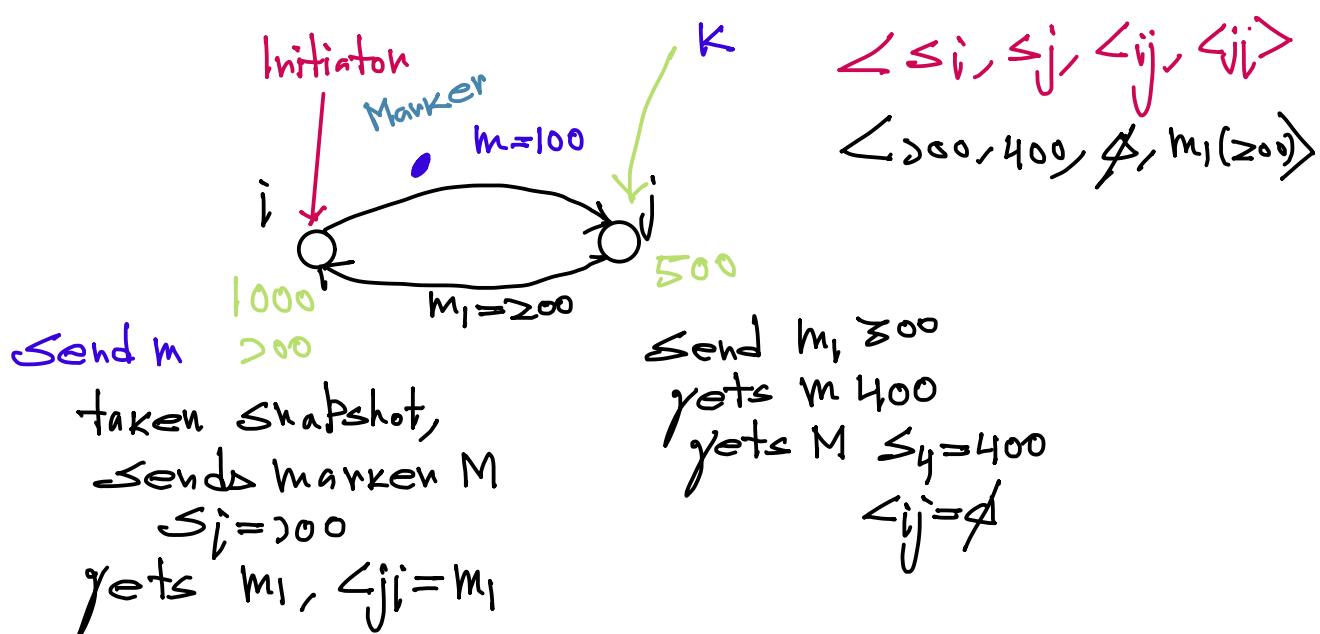


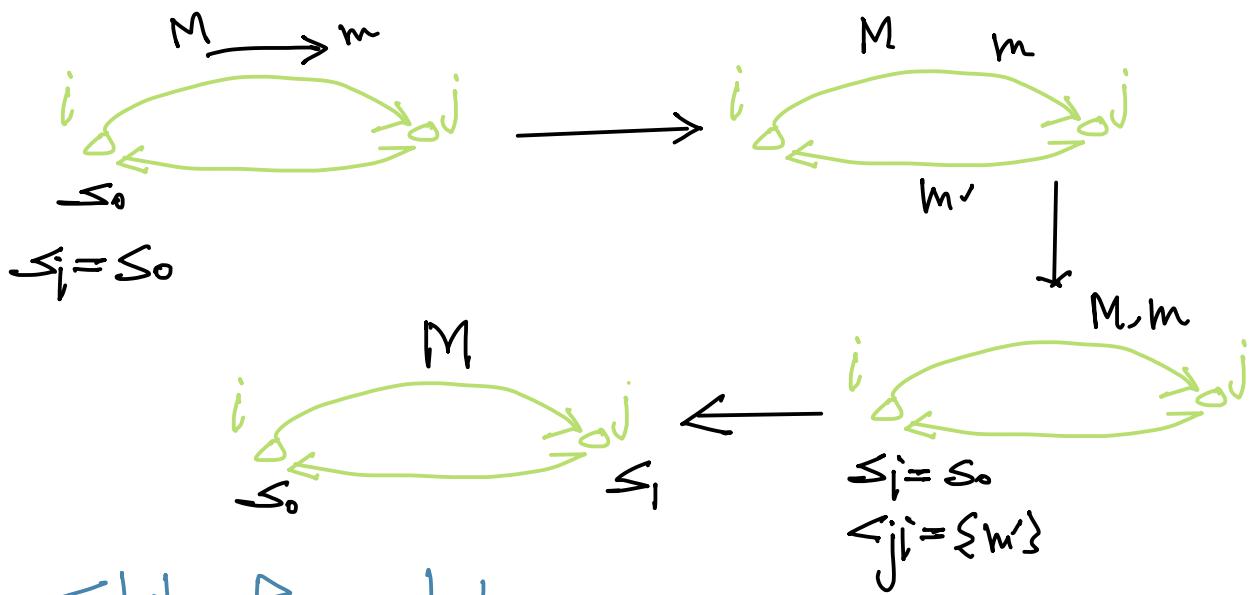
Q) Is consistent cut any good or should we always take strongly consistent cut?

Depends,

- if comm mode is unreliable communication, consistent cut is fine as the msgs in transit can get lost anyway, application should take care of it.
- if model is reliable comm, consistent is still OK but msg in transit has to be replayed first if recovery done from the cut.

C.L! - Keeps channel states explicitly





State Records!

$$s_i = s_0, s_j = s_1, \angle_{ij} = \emptyset, \angle_{ji} = \{m'\}$$

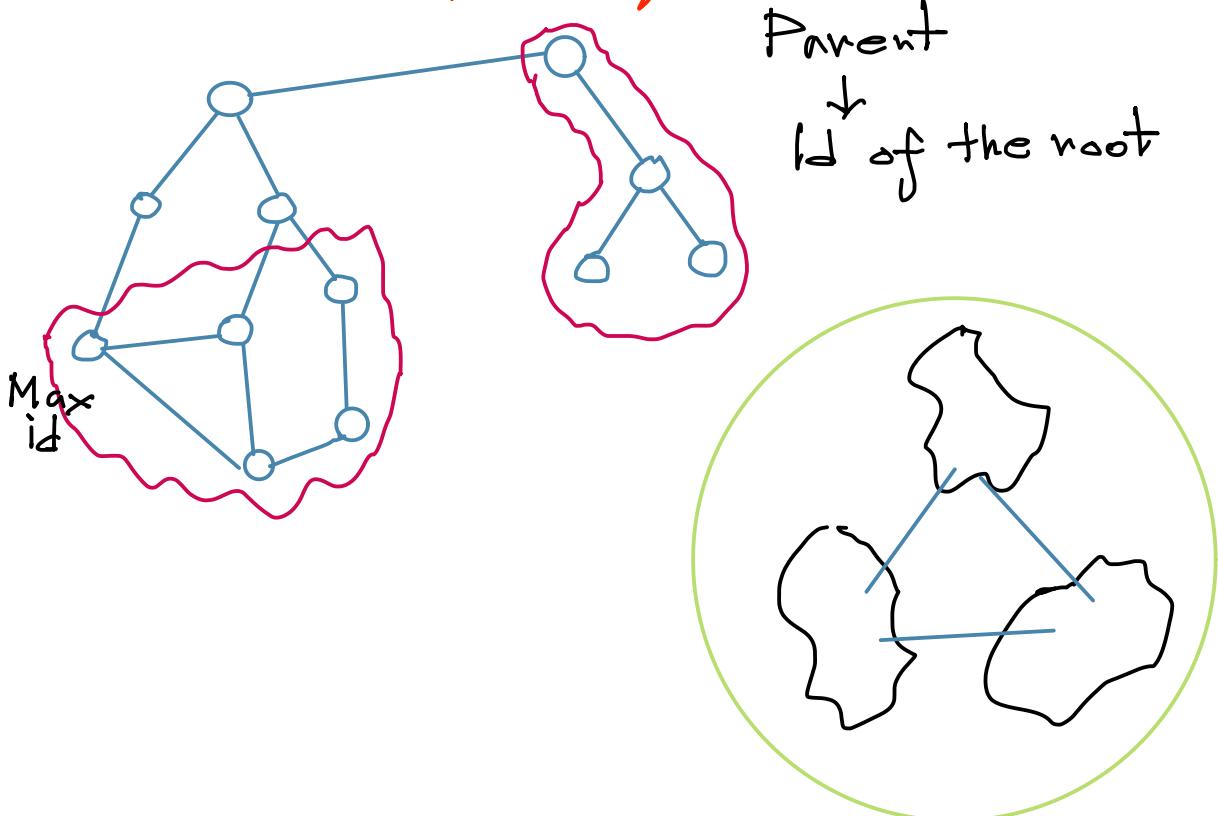
||

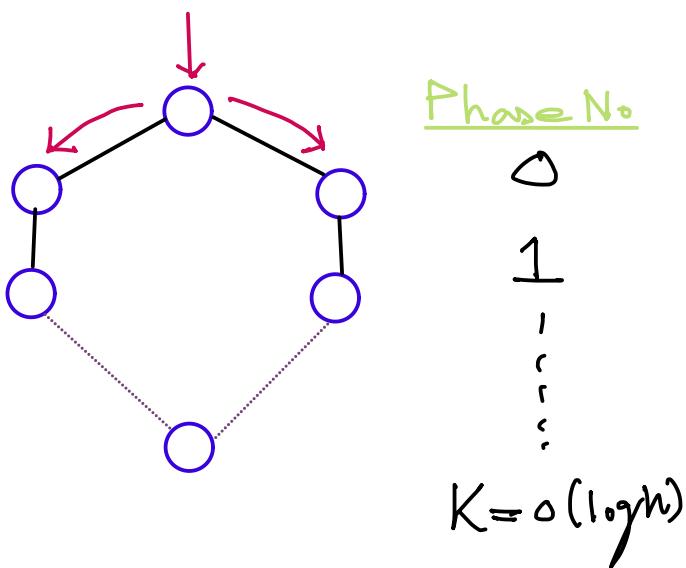
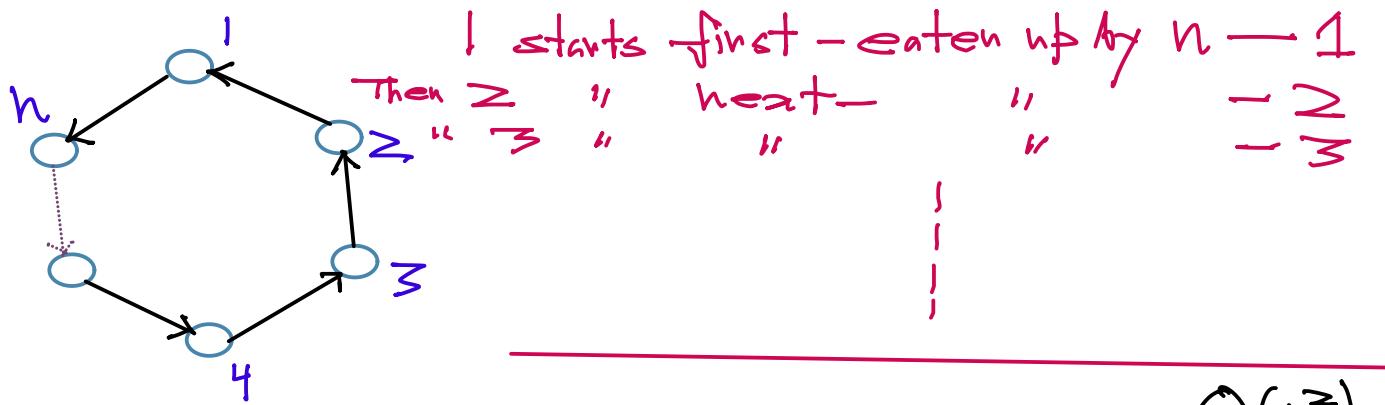


Q) What if no FIFO?

- Piggyback a flag to distinguish msg sent before & after taking state (The same work done by Marker in L.L)

⇒ How do I know how many messages?





$$\text{Total} = \Delta (n \log n)$$

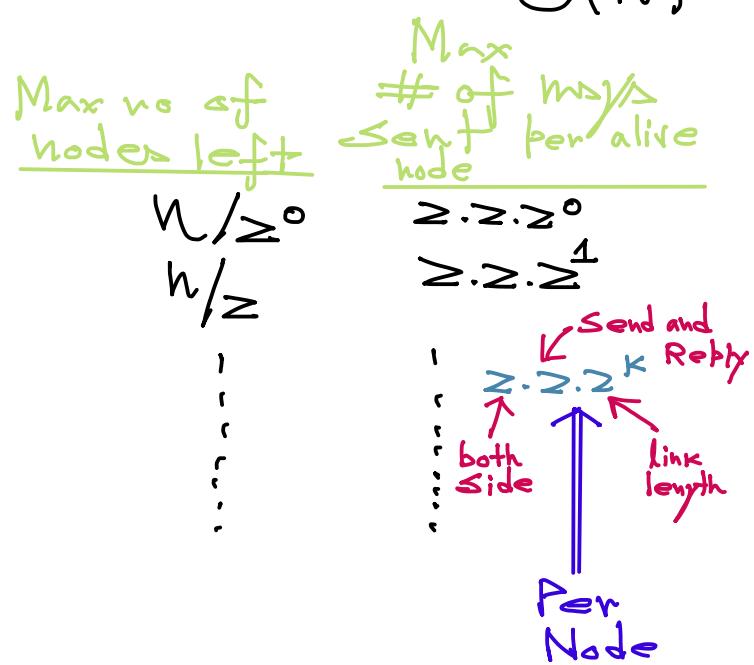
- Q) Study Unidirectional Case (Not Comparison)
 Q) Can we do better than $O(n \log n)$? (based)

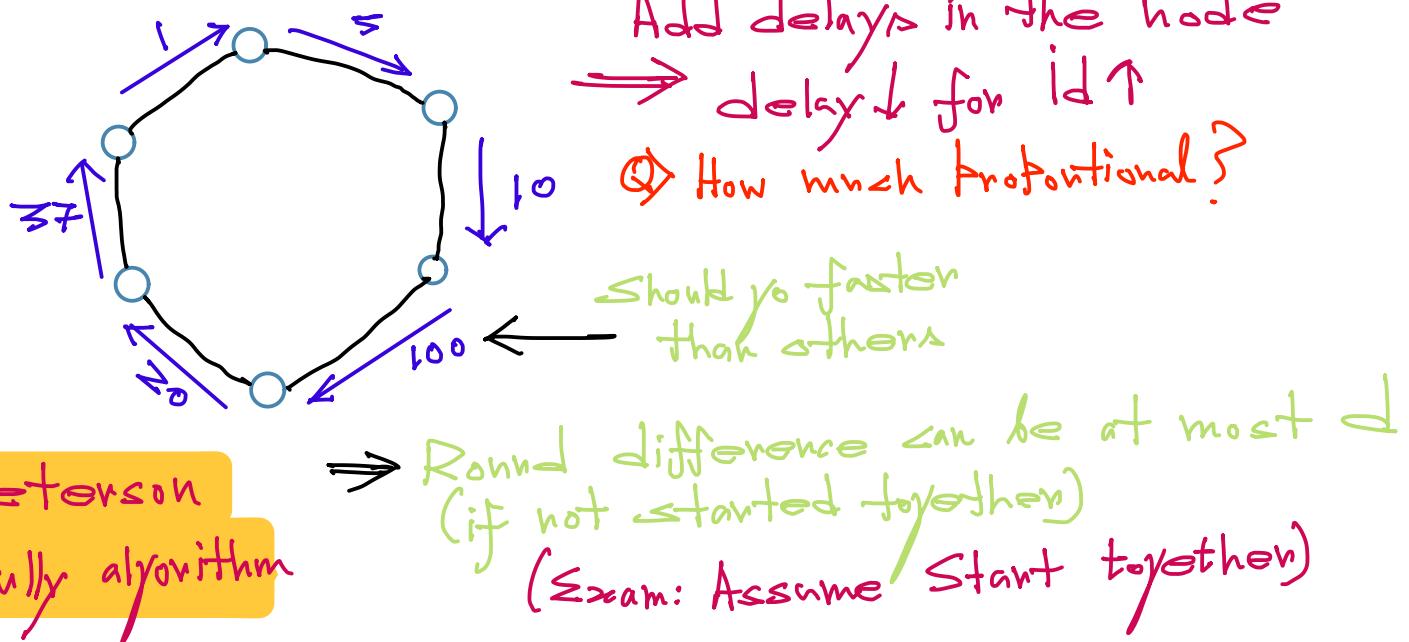
Round in a synchronous system!

Can be defined in several (slightly diff ways)

Here's how we will use it

Round = Max time for processing (α)
 + Max time for sending a msg (if needed to send)
 (& the msg reading the receiver) (β)





Renewable Resource:

- Resource can be renewed
 - ex: A lock, A semaphore

Consumable Resource:

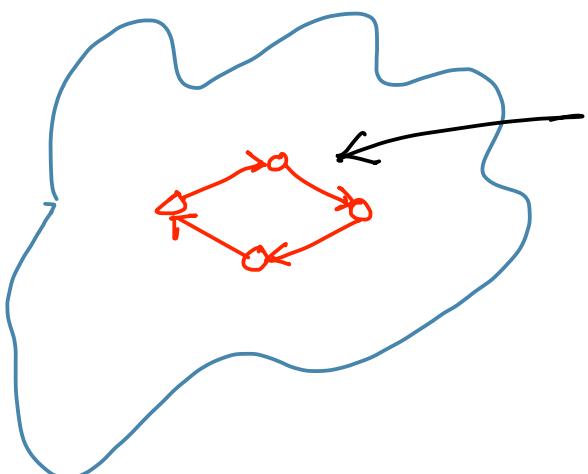
- Resource that can not be renewed
 - ex: A message

AND deadlock:

Process wait for a set of resources and all of them must be acquired before the process can start.

OR deadlock:

Any one can wake up the process.



No need for whole snapshot

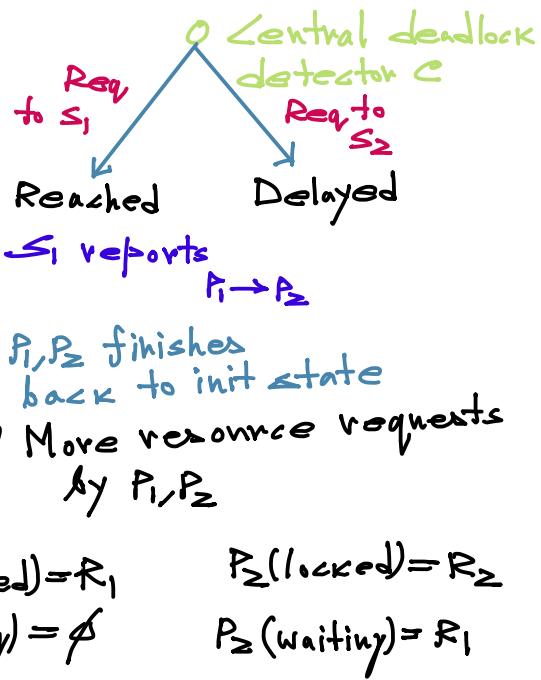
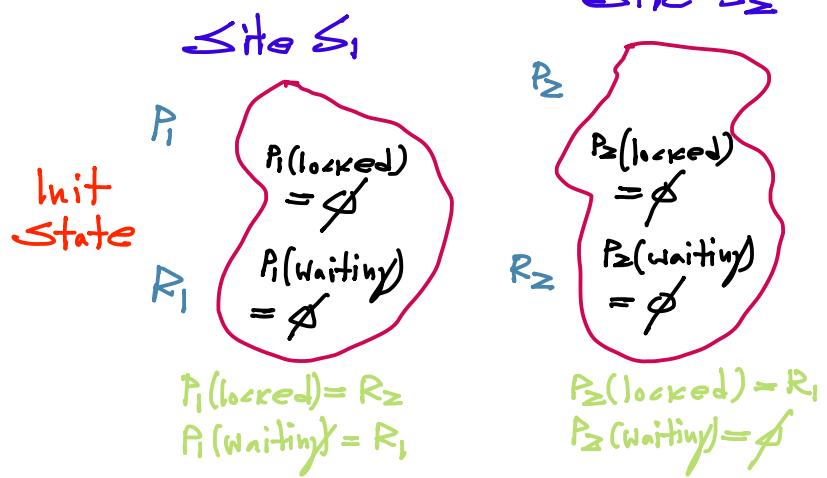
AND deadlock:

- Cycle in WFS \equiv deadlock

OR deadlock:

- Knot in WFG \equiv deadlock

Ho-Ramamirthy Two-Phase algo



Now, a request reaches S_2

S_2 report $P_2 \rightarrow P_1$

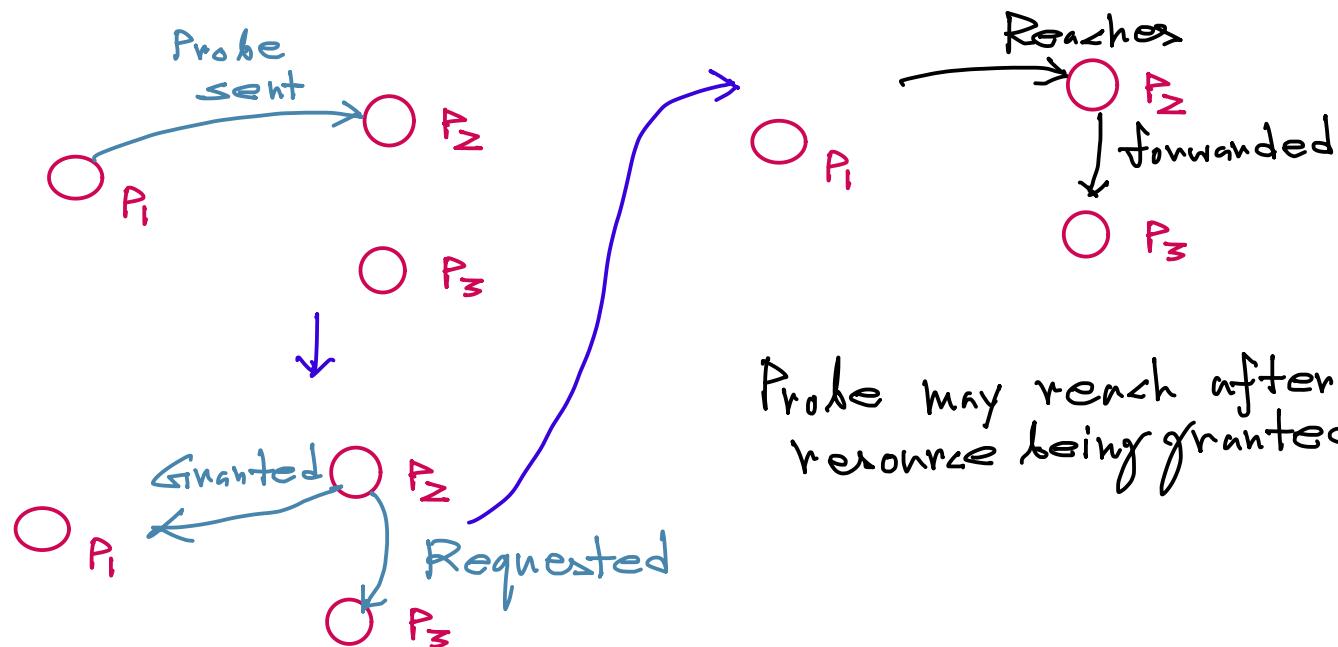
$\xleftarrow{P_1} P_2$

says deadlock = Phantom Deadlock

Q) Do we have

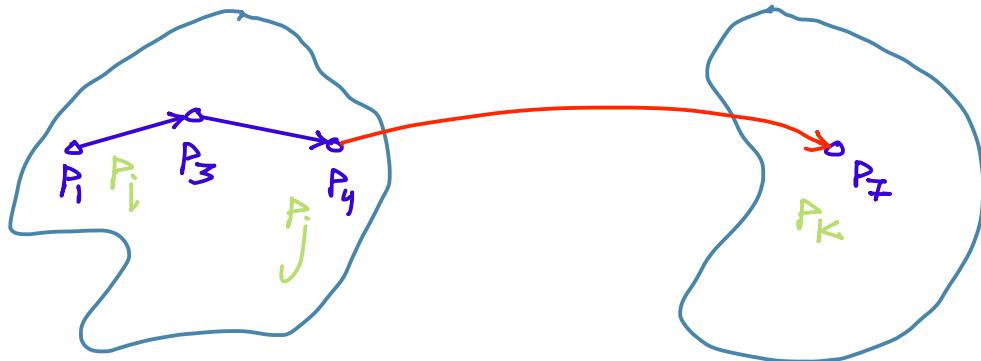
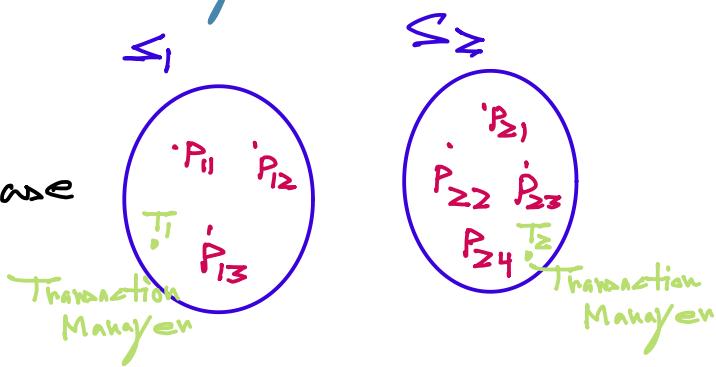
- ⇒ Keep two tables
 - > What resource every process locked and waiting for
 - > What process has locked a request and who are waiting for it.
- Idea: The table will report that no deadlock through inconsistency

Problem with simple probe for WFG



Background

Think local
Part of a
distⁿ database
(DDS)



XT< Algo!

Δu

$L(-----)$

$V(-, w, -, u, --)$

$NUM = \{ - \underline{w} - - \}$

Q) $V \sqcap u$ both add each other (Even for arbitrary ordering)
or don't

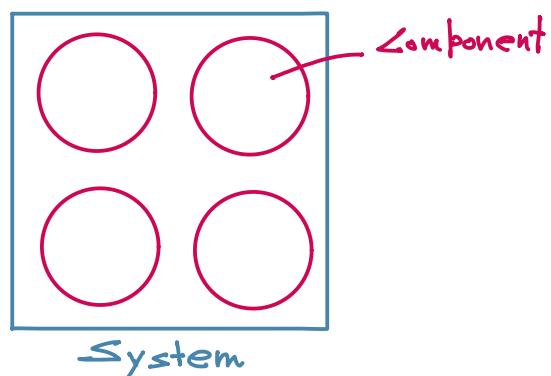
Fault— Some component has failed

Error— Manifestation of fault, may or may not cause the system to fail

Failure— System fails

Reliable \Rightarrow Available

Available \nRightarrow Reliable



Stable Storage

- Storage that is available even after the crash of the component.

Building Blocks

1. Stable Storage

- storage that persists even when a system crashed
- Allows state to be accessible even after crash

Q) How do you make your h/w more reliable

1. Servers
 - CPV
 - memory
 - disk
 - Power supply
 - fan
2. H/W
 - Switches
 - Cables
3. VPS
4. AC

How do you make a server reliable?

- CPV - Nothing specific
- Memory - ECC RAM
- Disk - RAID
- Power Supply / Fan - dual redundant

N/w?

- Switch - redundant power supply / fan
- extra ports
- redundant control cards
- for chassis based switch

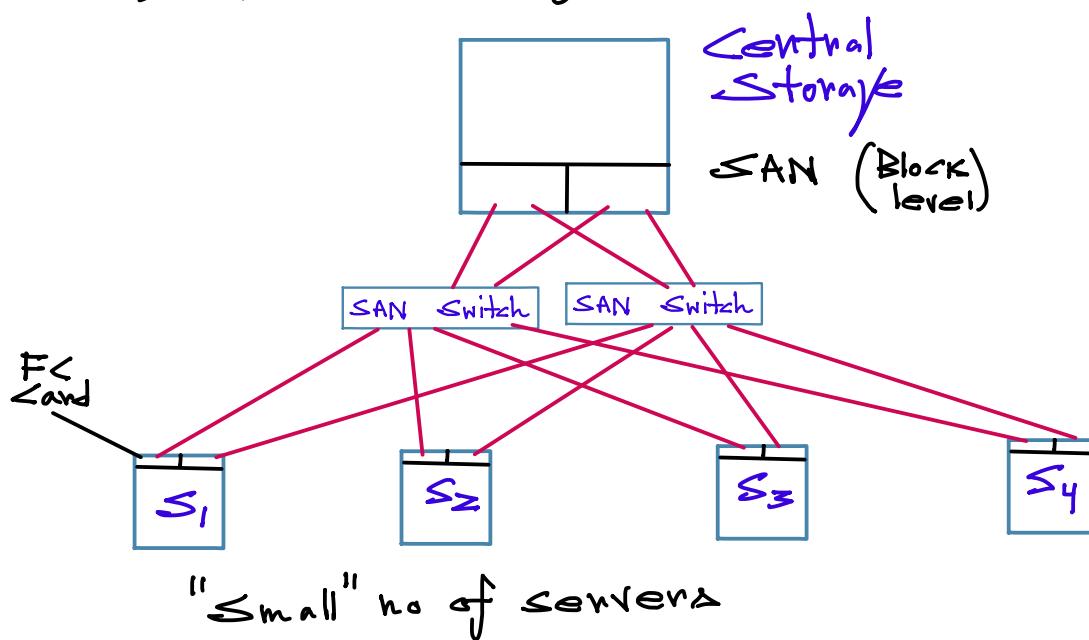
Q) How do you make storage reliable?

→ RAID "hot swap" disks

Still if the box fails as a whole (say powersupply failed)

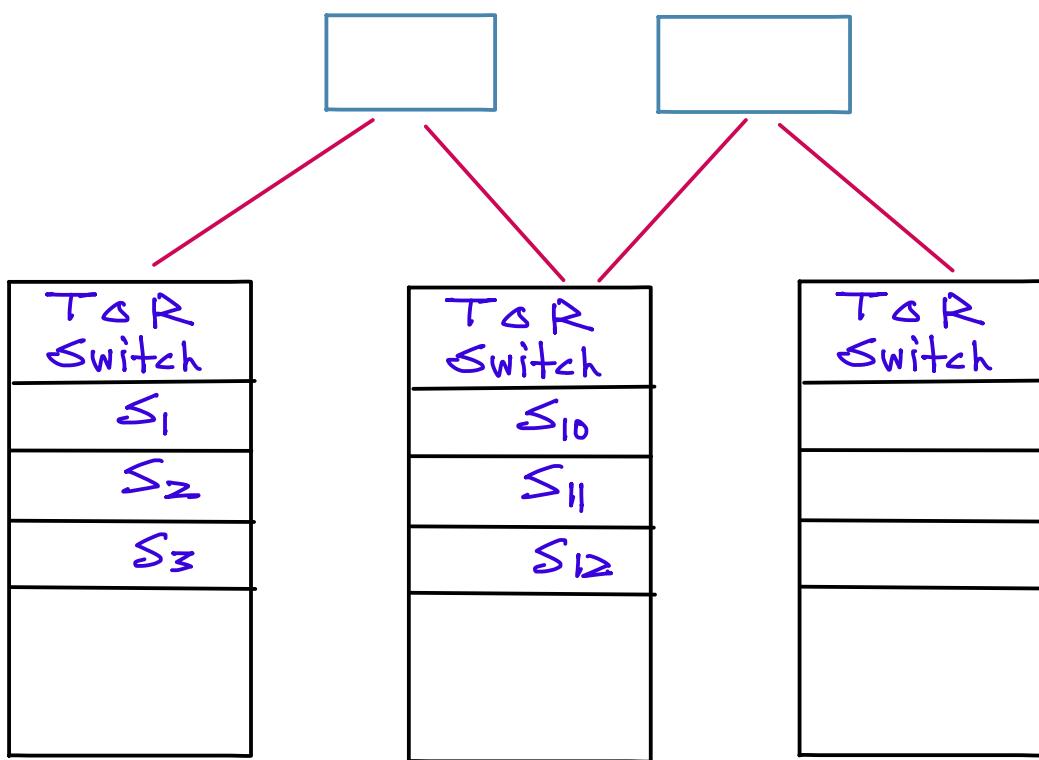
→ Storage is not accessible.

So, take it out of the box.

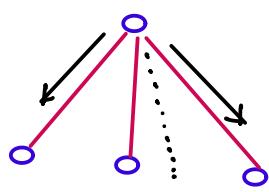


Data Center

TOR = Top of Rack



Byz Gen



	I/P	O/P
B Agree	the value	the value
Consensus	many values	the value
I/C	many values	many value (vector)

Q) If I know one, can I do the other?

Third (Implied) Condition

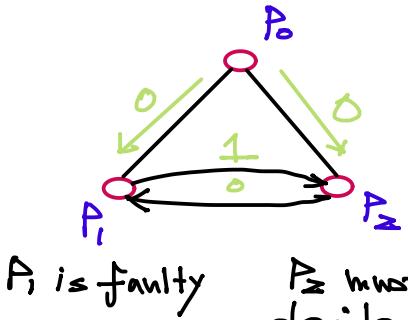
in all three version

- All non-faulty processes eventually decide a value

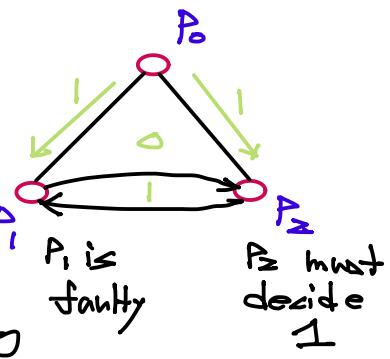
Agreement under Byzantine Fault! why assume only \geq values

Impossible with ≥ 3 process & one fault. $n=3, m=1$ $0 \& 1$

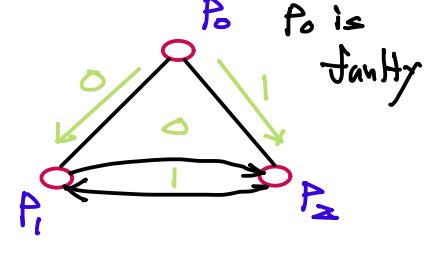
Case 1:-



Case-2:-



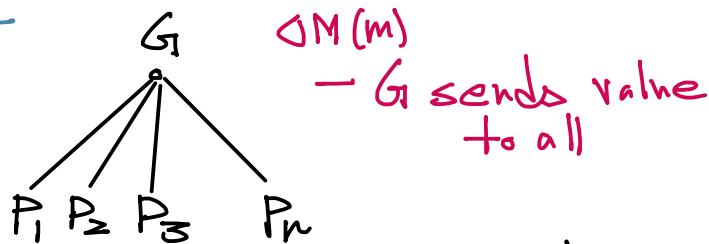
Case-3:-



P_1 must decide 0 (from Case 1 decision of P_2)

P_2 must decide 1 (from Case 2)

$L \leq P_1^1$



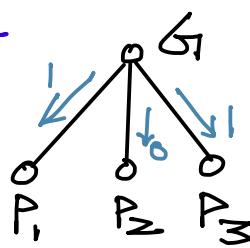
Say value received by P_i from G_1 is v_i

Each P_i starts a $\Delta M(m-1)$ with P_i as general & v_i as the value sent by it. $P_i^x = P_i$'s decision on what G_1 sent to P_x .

Each process maintains a vector

So $\Delta M(m-1)$ for P_i is for all other $n-2$ nodes to decide on what value is sent to P_i in first round.

$n=4, m=1$



	1	\geq	3
P_1	1	0	1
P_2	1	0	1
P_3	1	0	1

$\Delta M(1) \rightarrow$ by G_1

$\Delta M(0) \rightarrow$ by P_1

→ Sends 1 to P_2 & P_3

$\Delta M(1) \rightarrow$ by P_2

→ Sends 0 to P_1 & P_3

$\Delta M(2) \rightarrow$ by P_3

→ Sends 1 to P_1 & P_2

\geq in $\Delta M(n)$	No of invocation	No of msgs/ Invocation
m	1	$n-1$
$m-1$	$n-1$	$n-2$
$m-2$	$(n-1)(n-2)$	$n-3$
\vdots	1	1
Δ	1	.

$$Total = O(n^m)$$

Q) $n=7, m \geq$

Q) Given Byzantine Agreement sol^h, can you do consensus?

\Rightarrow Each process has a value
 " " starts agreement br with its value

At the end, each process has a vector

Where i-th component = Value P_i started Byz
agreement with

(So you solved I_G)

↓ Take majority

Consensus Solved

Q) Other way round?

1/p: one general with a value, $n-1$ lieutenants.

General send to all lieutenants.

Lieutenants ($n-1$) solve consensus to agree on value.

Notes about $\Delta M(n)$ also!

For the $n=3, m=1$ impossibility

notes that the problem case was if the general is faulty.

If general is not (known), problem is trivial.

But the problem is that you do not know whether the general is faulty or not.

$\Delta M(m)$ - Started by 1 node

↓ Causes

(n-1) $\Delta M(m-1)$ to start.

We reasoned if $\Delta M(m-1)$ is so what they are supposed to, $\Delta M(m)$ will work.

Case 1: General is faulty

So rest of the nodes,

$n_1 = n-1$ nodes, at most $m_1 = m-1$ faulty nodes.

$n_1 > 3m_1$ holds \rightarrow So OK

Case 2: if general is non-faulty,
all m-faults can be in (n-1) node

So, $n_1 > 3m$ may not hold.

So, is $\Delta M(m-1)$ supposed to work?

Works since general has sent the same value to
"most" of the Processors

Assignment ≥ 1 :

To be submitted hand work on paper

Work out L-S-P if $n=7, m \geq$

When (1) General is faulty (Δ sends Δ to P_1, P_2, P_3)
 Δ to P_4, P_5, P_6)

(2) General is non-faulty (Choose P_1, P_2 as faulty)

Rank tolerant agreement!

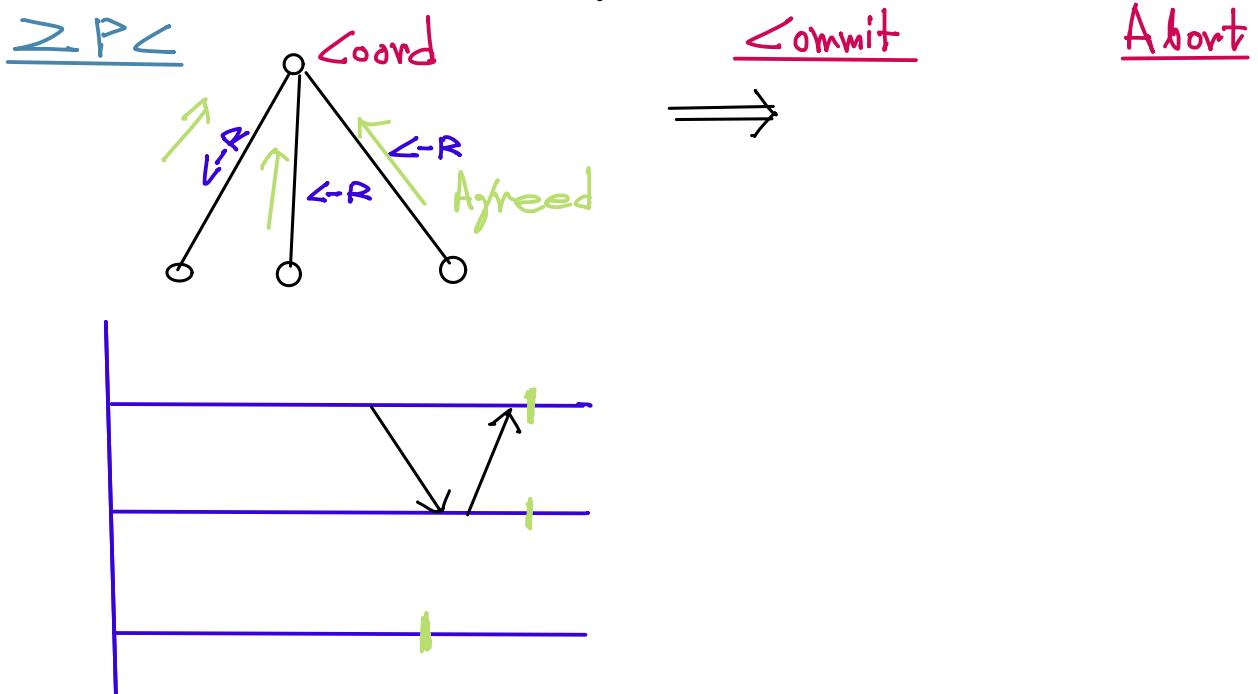
Suppose at end, there are \geq Processors P_i & P_j which
decides on \geq diff minimum. Say V_i and V_j

Δ say $V_i < V_j$

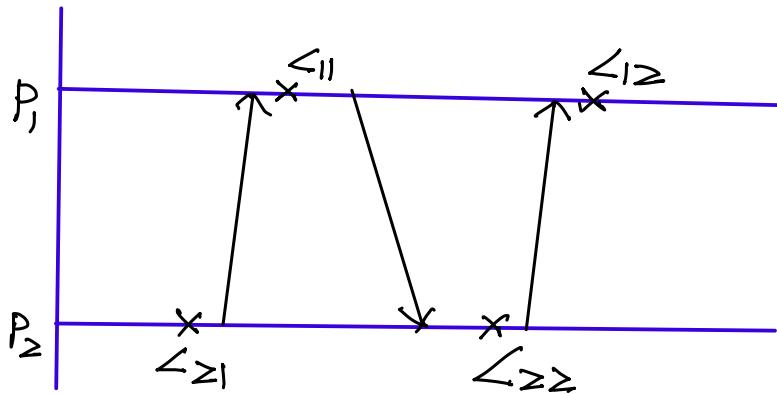
So, P_j does not have V_i in its set. V_i must have been
"hidden" for P_j

Can happen only if at any round, whoever knows about v_i , fails before sending to p_j (1 failure)
 Can be hidden for a maximum of m rounds (per round)
 Still has a failure free round at end of which p_j will know v_i

Adv: — You can check convergence after every round.
 \Rightarrow No new process failed means convergence

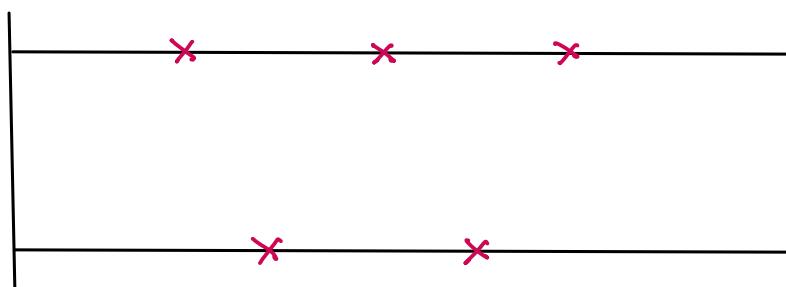


Async / Uncoord checkpoint



Dominos Effect

Communication Induced



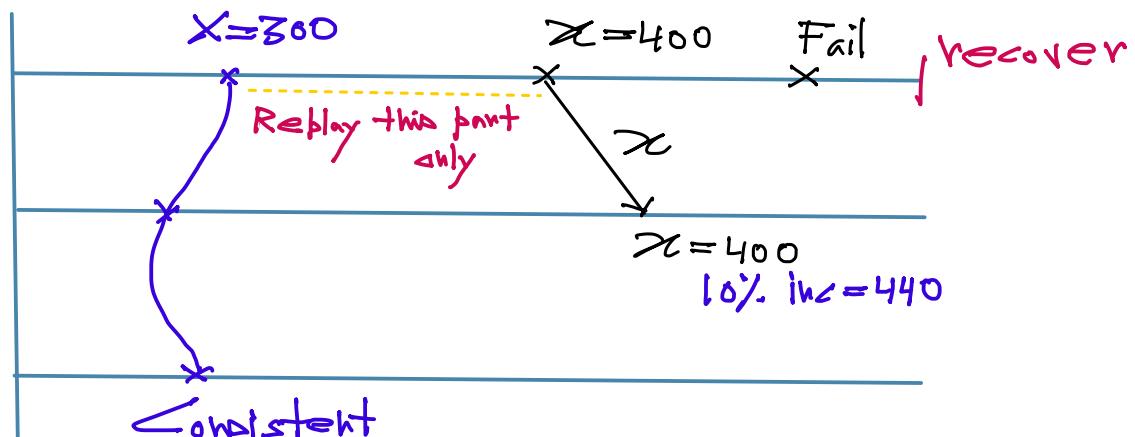
(\leq_{12}, \leq_{22})

→ Not Consistent

Need to coordinate which checkpoints to take from each process to get a consistent recovery line.

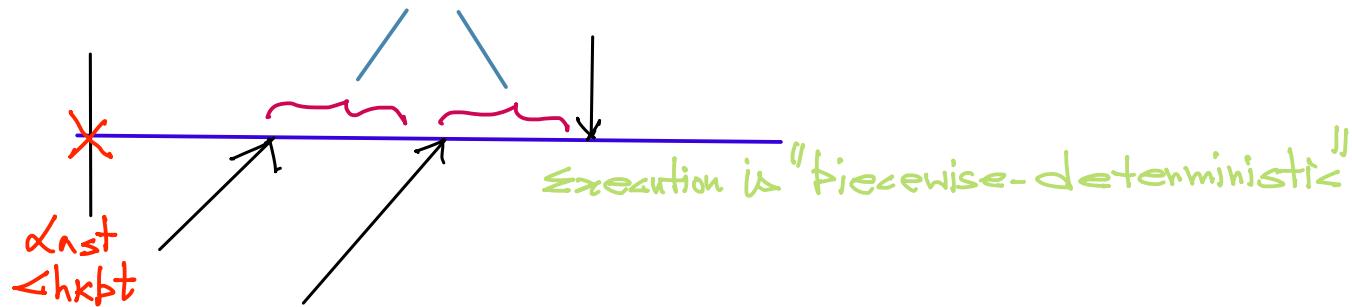
No communication so far
 So could have done uncoord checkpoints

Name of Async, Sync, Comm-Induced handle the output commit
problem



① only the failed process has to recover (Has to do lot of things although)

These parts
are deterministic



Koo-Teng:-

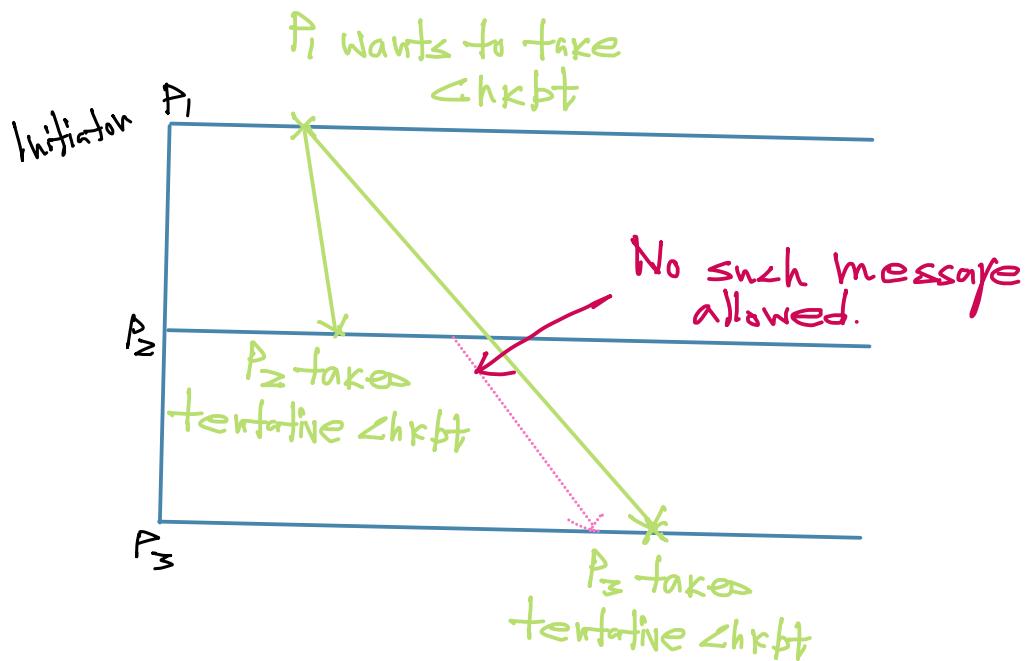
Checkpointing Protocol!

Idea!-

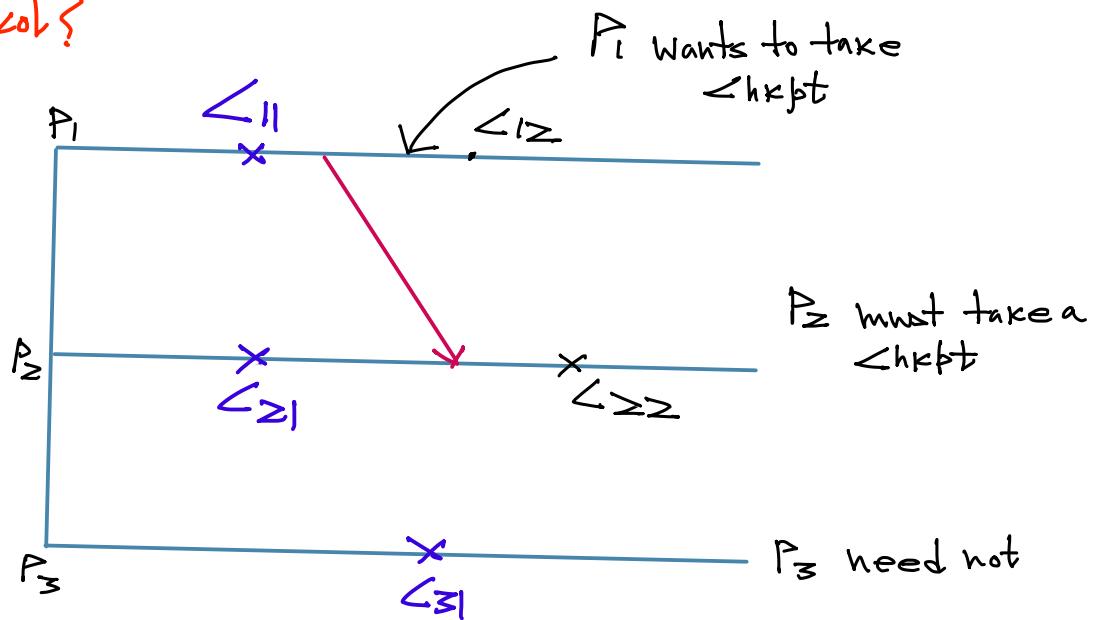
1. One process initiates
 - takes "tentative" checkpoint
 - asks other process to take checkpoint
2. Any other process takes a "tentative" checkpoint on receiving the request & answers "yes", else answers "no"
3. Coordinator / Initiator after getting all responses
 - If any one "no", ask everyone to discard the checkpoint
 - Else make its own checkpoint "permanent", ask others to do so.

4. If asked, other makes it permanent.

Between the time a $\langle \text{chkpt} \rangle$ is tentative & permanent,
do not send any other messages.

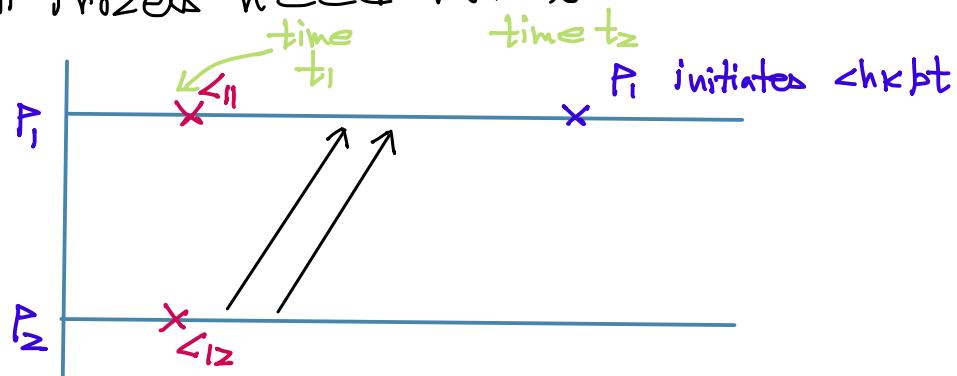


Q) Do all processes Need to be part of checkpointing protocol?



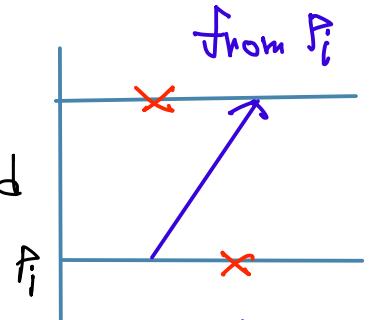
$$(\langle \text{I2} \rangle, \langle \text{Z2} \rangle, \langle \text{Z3} \rangle) = \text{Consistent}$$

So all boxes need not be involved. So, who all?



Note, between t_1 & t_2 ,

- ① If P_i has not read any msg from another process P_j ,
then P_i need not take $\leq h_{kpt}$
- ② If P_i has read a msg from P_j , but P_i has
already included it in its last $\leq h_{kpt}$, P_i need
not take another $\leq h_{kpt}$.



⇒ Each app may has a distinct, monotonically increasing label
Each process keep the following data structures:

$l_l_r_x[y]$ = label of last msg received by X from Y after
 X 's checkpoint

$f_l_s_x[y]$ = label of first msg sent by X to Y after
its last checkpoint

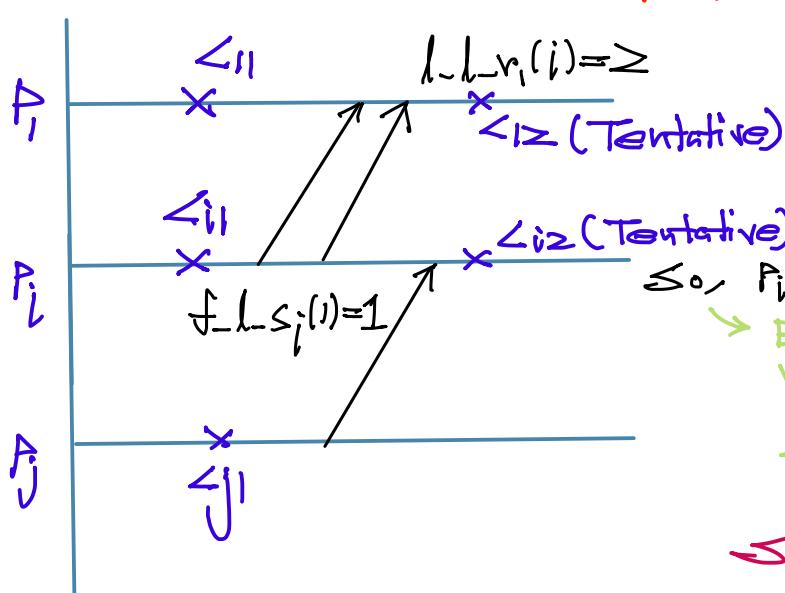
So, P_i = initiator

= sends only to process P_j such that
 $l_l_r(i) = \text{NULL}$

In receiving a req from P_i ,

P_j does this

P_j takes a $\leq kpt$ if $\underline{l_l_r(i)} \geq f_l_s_j(t)$
Sent in req msg for P_i to P_j



So, P_i must $\leq h_{kpt}$

→ But then by same logic P_j
must also take $\leq h_{kpt}$.

Even though the initiator has
not sent a request

So, how P_j starts as an initiator
on behalf of P_i

A non-communicating node will lose lost data.
So, it will itself initiate one

Rollback Recovery

Ideal!

- The initiator
- Everyone rolls back to its last checkpoint

⇒ Consistent but everyone loses somewhere
Q) Do everyone need to rollback?

P_i = initiator

= failed & now trying to recover

