

UCS2701 DISTRIBUTED SYSTEMS

Unit 2

Logical Time and Global State

Logical Time : Physical clock synchronization ; NTP - A framework for a system of logical clocks - scalar time - Vector time : Message ordering and group communication - Message ordering paradigms - Asynchronous execution with synchronous communication - synchronous program order on an asynchronous system - Group communication - Causal Order (CO) , Total order : Global state and snapshot recording algorithms: Introduction - system model and definitions - Snapshot algorithms for FIFO channels.

* Physical Clock Synchronization

- In centralized systems, there is only a single clock . A process gets the time by simply issuing a system call to the kernel.
- In distributed systems , there is no global clock or common memory .
- Clock synchronization is the process of ensuring that physically distributed processors have a common notion of time . Due to different clock rates, the clocks at various sites may diverge with time . Clock synchronization is used to correct clock skew .

→ Clocks that must not only be synchronized with each other but also have to adhere to physical time are termed physical clocks.

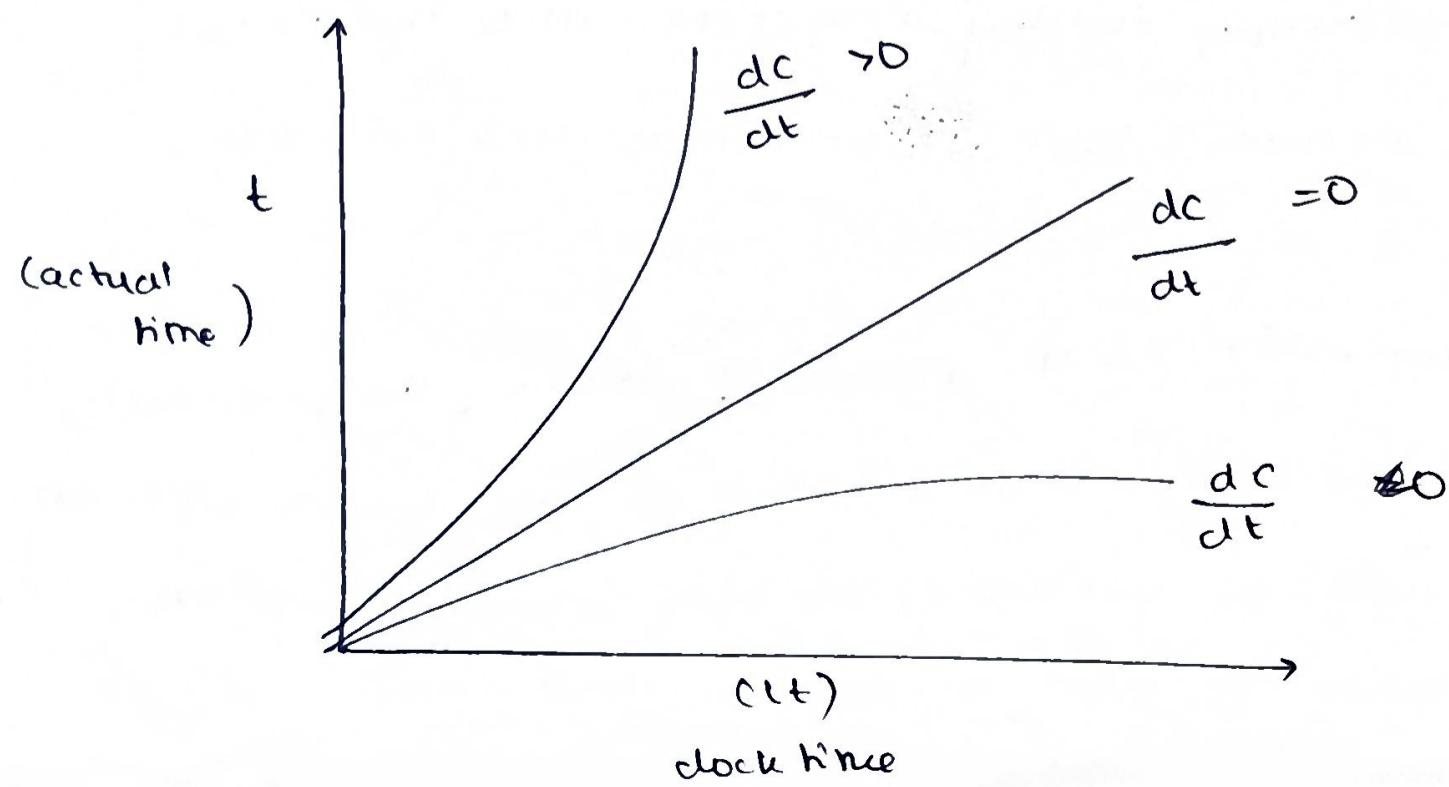
Physical Clock Inaccuracies

- ① Offset: Difference between the time reported by a clock and the real time.

$$\text{offset}(\text{clock}_a) = c_a(t) - t$$

- ② Skew: The skew of a clock is the difference in the frequencies of the clock and the perfect clock. The skew of a clock c_a relative to clock c_b at time t is $(c'_a(t) - c'_b(t))$

- ③ Drift: The drift of clock c_a is the second derivative of the clock value with respect to time, namely $c''_a(t)$. Relative to another clock c_b , it would be $c''_a(t) - c''_b(t)$.

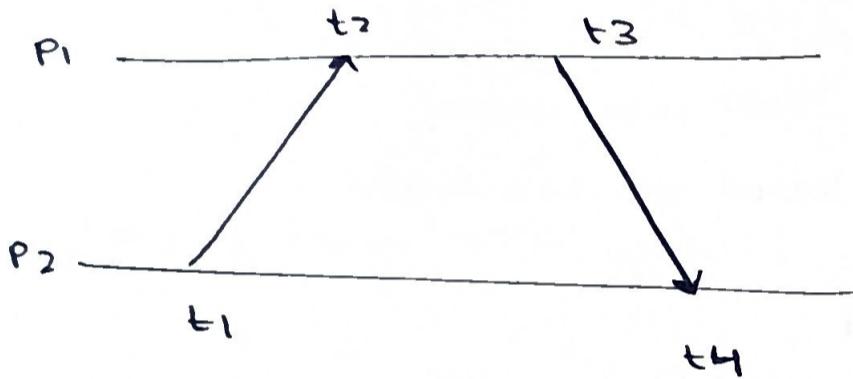


Algorithms for Synchronizing Physical Clocks

1. Christian's algorithm
2. Berkeley's algorithm
3. Network Time Protocol (NTP)

A. Christian's Algorithm

Consider two clocks P_1 and P_2



$$t_4 - t_1 = \text{round trip time (RTT)}$$

$$t_3 - t_2 = \text{computation delay}$$

$$\left. \begin{array}{l} t_2 - t_1 \\ 2 \\ t_4 - t_3 \end{array} \right\} = \text{communication delay}$$

Step 1 : P_2 computes RTT

Step 2 : P_2 sends the computation delay

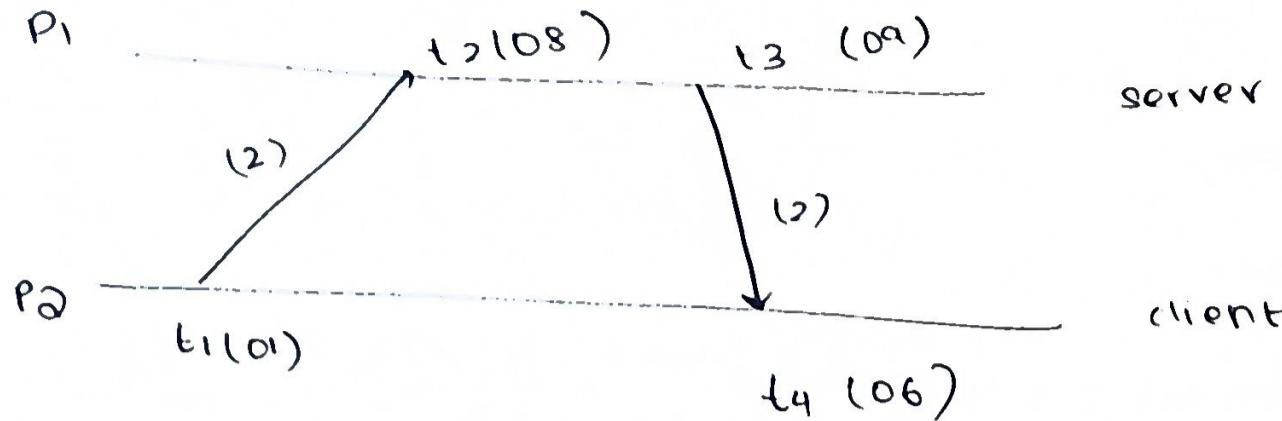
Step 3 : P_2 calculates $\frac{(RTT - \text{comp. delay})}{2}$ = one-way delay

Step 4 : Adjust t_4 's time :

$$t_4' = t_3 + \text{one way delay}$$

Example

Consider the following scenario:



Step1 : compute $RTT = t_4 - t_1 = 5$

Step2 : compute computation delay $= t_3 - t_2 = 9 - 8 = 1$

Step3 : compute $\frac{(RTT - \text{comp. delay})}{2} = \frac{5-1}{2} = 2$ one way delay

Step4 : t_4 needs to adjust its time based on t_3 's time

$$t_4' = t_3 + \text{one-way delay}$$

$$t_4' = 9 + 2$$

$$\boxed{t_4 = 11}$$

Limitations of Christian's Algorithm \rightarrow algo will fail when the two communication delays are not the same.

B. Berkeley Algorithm

\rightarrow The Berkeley algorithm works by averaging the clocks of all participating nodes to minimize discrepancy.

\rightarrow One node (called the master) is responsible for initiating and coordinating the synchronization process

Step1 : Master polls the nodes - The master node polls all participating nodes for their current time. Each node replies with its local clock time.

Step2 : Adjust for Round Trip Delay - When the master receives the replies, it estimates the round trip time (RTT) for each node, and adjusts the received time values to account for the communication delay.

$$T_{adjusted} = T_{received} + \frac{RTT}{2}$$

Step3 : Calculate average time

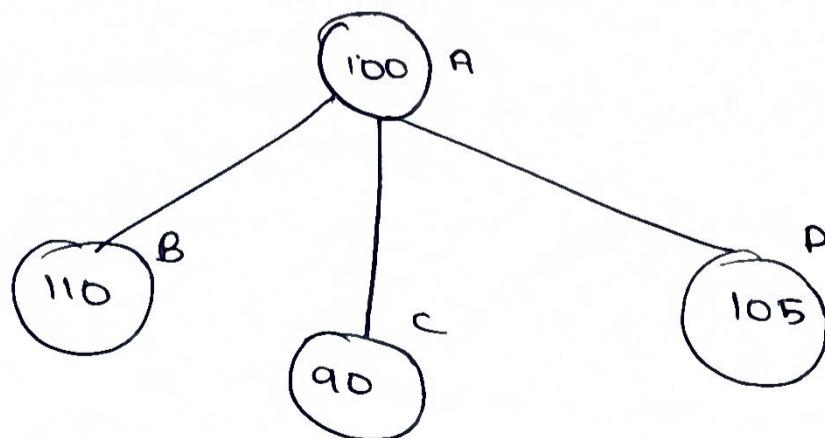
$$T_{average} = \frac{\sum \text{Adjusted Times of all nodes}}{\text{number of nodes}}$$

Step4 : Send adjustments to nodes

$$\text{Offset} = T_{average} - T_{local}$$

Example

Consider the following scenario



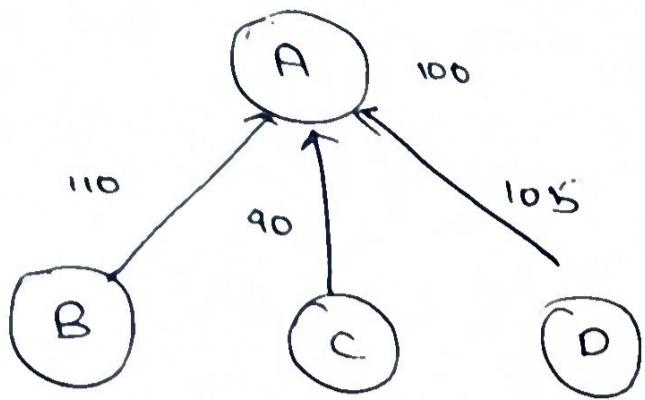
RTT

$$B = 4$$

$$C = 6$$

$$D = 8$$

Step1 : A polls times from each node



Step2 : Adjust for RTT

one way delay \rightarrow node B = 2
node C = 3
node D = 4

Adjusted values: B = 112

C = 93

D = 109

Step3 : calculate the average time

$$\text{average time} = \frac{100 + 112 + 93 + 109}{4} = 103.5$$

Step4 : Compute offsets

$$\text{offset} = \text{average time} - \text{node's time}$$

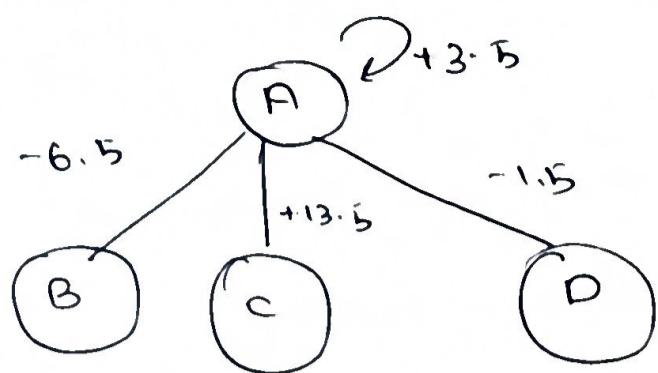
$$A: 103.5 - 100 = +3.5$$

$$B: 103.5 - 110 = -6.5$$

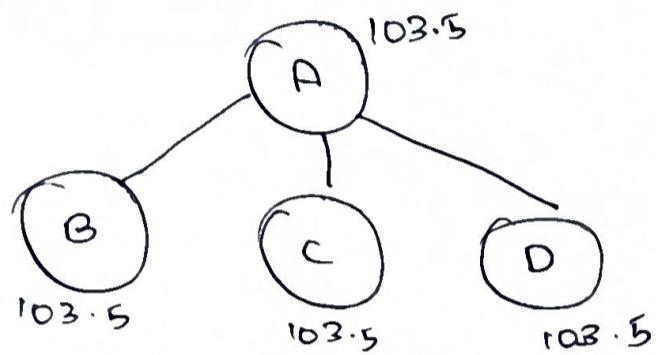
$$C: 103.5 - 90 = +13.5$$

$$D: 103.5 - 105 = -1.5$$

Send offsets to each node



Step 5 : Compute the new time



c. Network Time Protocol (NTP)

- Enables clients across the Internet to be accurately synchronized to UTC despite message delays.
- Works in a hierarchical fashion
- Primary servers are connected directly to a time source such as a radio clock; secondary servers are synchronized with primary servers.
- The servers are connected in a logical hierarchy called a synchronization subnet whose levels are called strata.
- Stratum 2 servers are synchronized directly with the primary servers.
- Stratum 3 servers are synchronized with stratum 2 servers and so on.

Working

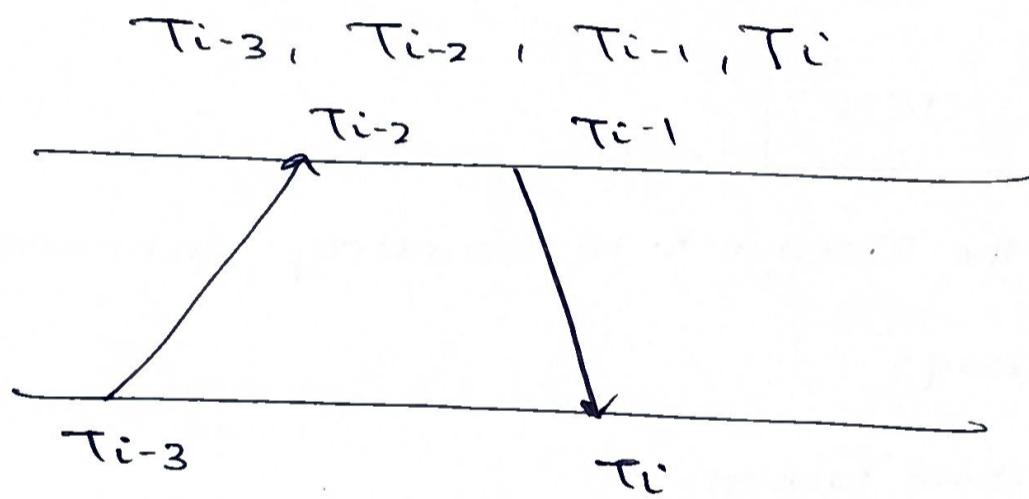
→ Both servers exchange timing messages to compute the Offset (O_i) and Delay (D_i) between their clocks.

$O_i \rightarrow$ how much the 2 clocks differ

$D_i \rightarrow$ measures the time taken for messages to travel between the servers.

→ The key relationship is: $A(t) = B(t) + O$

→ The servers exchange 4 time stamps:



Step 1: @calculate offset

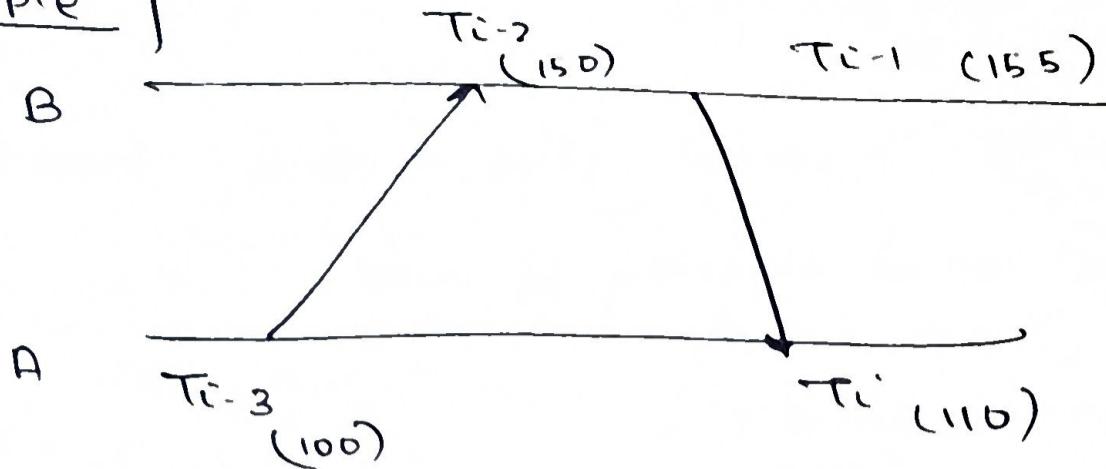
$$O = \frac{(T_{i-2} - T_{i-3}) + (T_{i-1} - T_i)}{2}$$

Step 2: @calculate round trip delay

$$D_i = (T_i - T_{i-3}) - (T_{i-1} - T_{i-2})$$

Step 3: The algorithm stores several pairs of (O_i, D_i) . Choose the offset corresponding to the minimum delay D_i .

Example



Step 1 : Compute D

$$D = \frac{(T_{i-2} - T_{i-3}) + (T_{i-1} - T_i)}{2}$$

$$= \frac{(150 - 100) + (155 - 110)}{2} = \frac{50 + 45}{2} = \underline{\underline{47.5}}$$

Step 2 : Compute Round Trip Delay

$$D_i = (T_i - T_{i-3}) + (T_{i-1} - T_{i-2})$$

$$= (110 - 100) + (155 - 150)$$

$$= 5$$

Step 3 : Adjust A's clock value : $T_i + D_i$

$$= 110 + 47.5 = \underline{\underline{57.5}}$$

* Physical Clock vs. Logical Clock

Physical Clock - internal clock present in a computer

Logical Clock - keeps track of event ordering among related (causal) events.

* Happened Before Relation

- Lamport defined the 'happened before' relation - denoted as \rightarrow , which describes the causal ordering of events.
- There are 2 possible cases:
- (i) a and b are in the same process and a occurs before b
then
 $a \rightarrow b$.
 - (ii) a and b are in different processes, and if a is the event of sending a message 'm' in one process and b is the event of receiving that message m in another process,
then
 $a \rightarrow b$.

$a \rightarrow b \Rightarrow a$ causally affects b

* Strict Partial Order

- Lamport's causal ordering also has the following properties:
- (i) irreflexivity $\Rightarrow a \not\rightarrow a$
 - (ii) antisymmetric $\Rightarrow a \rightarrow b \Rightarrow b \not\rightarrow a$
 - (iii) transitive $\Rightarrow a \rightarrow b, b \rightarrow c \Rightarrow a \rightarrow c$

* Lamport's Logical Clock Algorithm

If $a \rightarrow b$, then $c(a) < c(b)$

1. If a and b belong to the same process and $a \rightarrow b$

$$c(b) = c(a) + d$$

d is usually 1

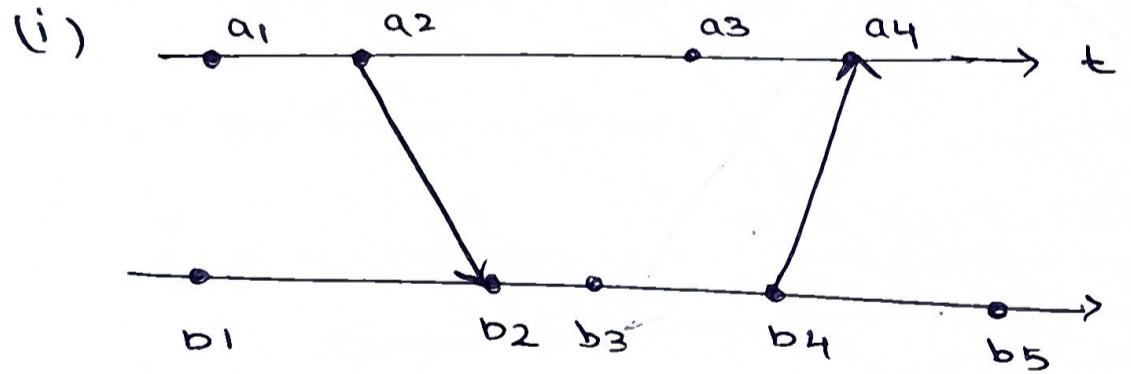
2. If a and b belong to different processes and $a \rightarrow b$, i.e. a is a send event and b is a receive event, m is the message

$$tm = c(a)$$

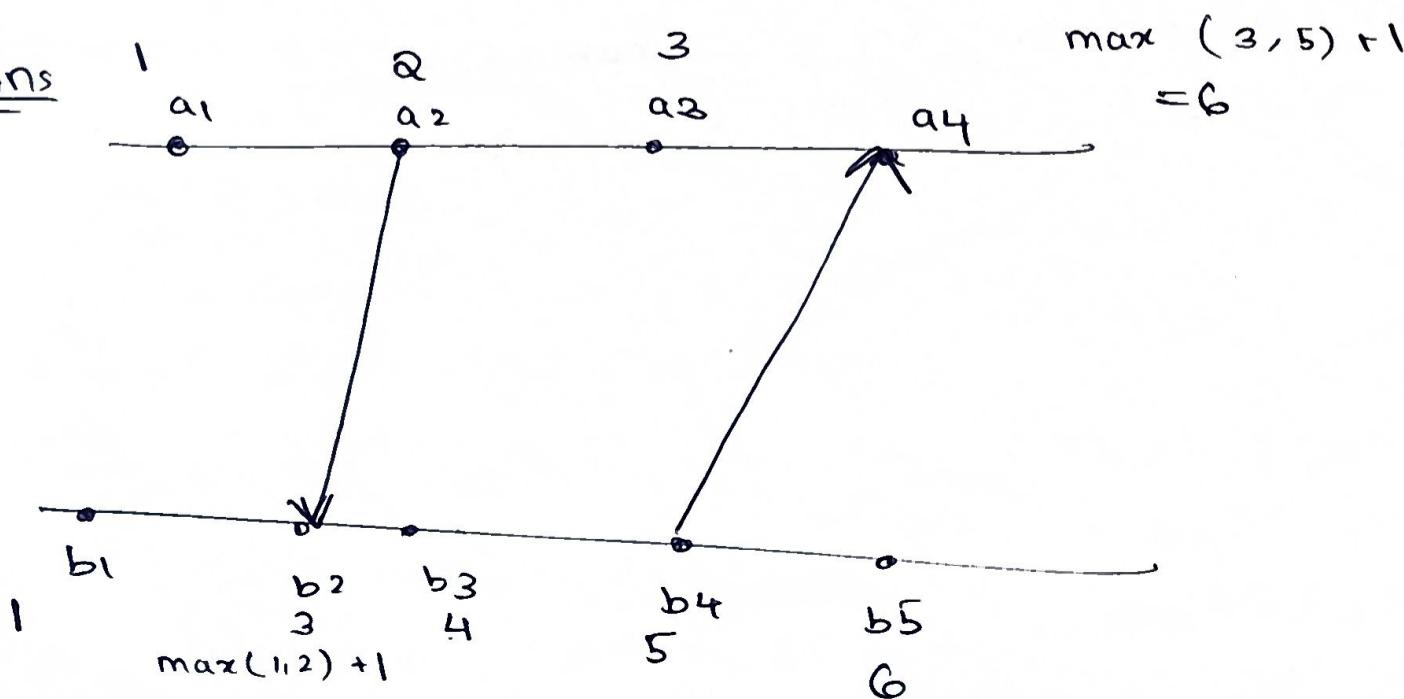
$$c_j(b) = \max(c_j(b), tm) + 1$$

Example

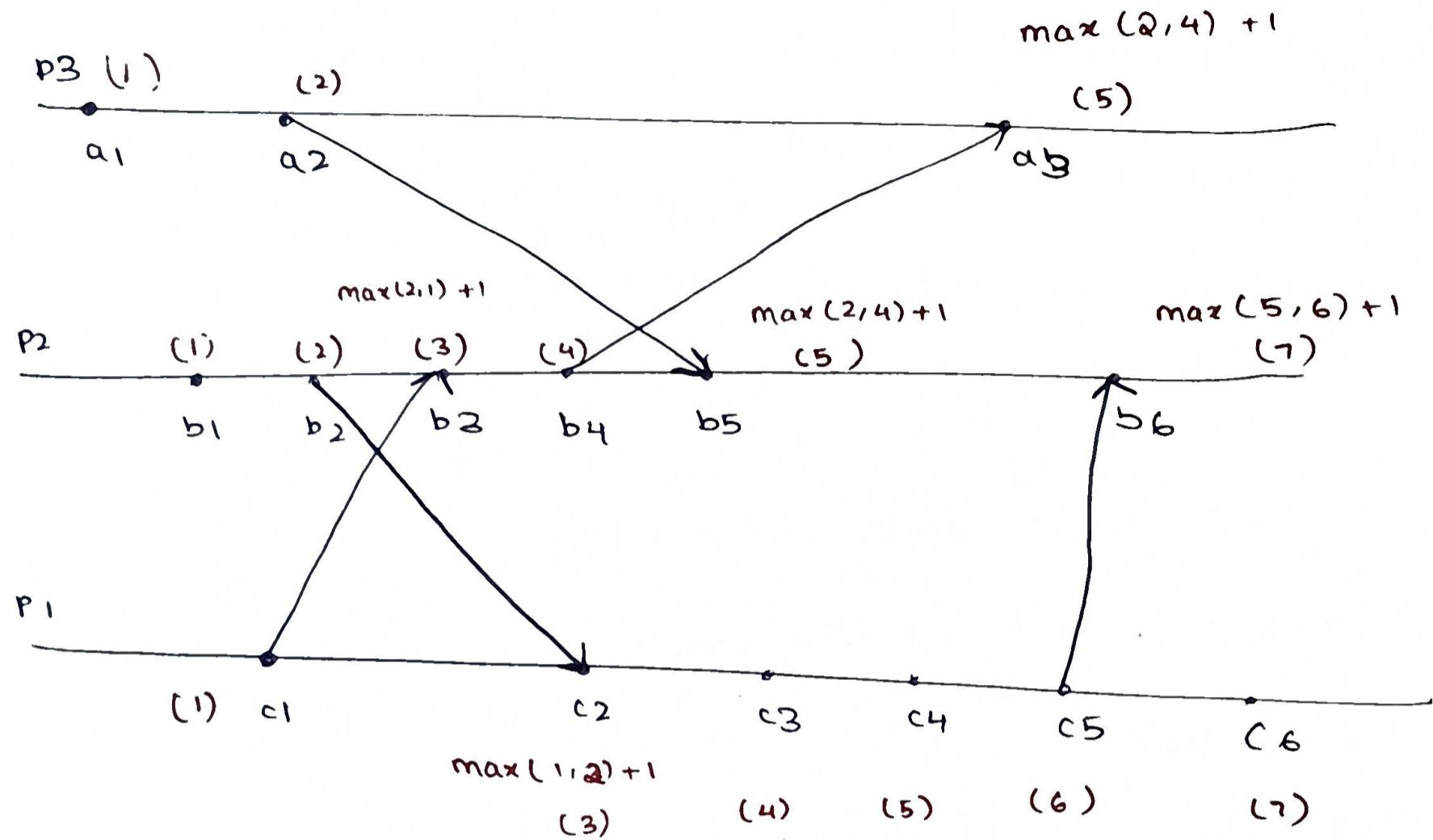
Calculate the Lamport's logical clock values



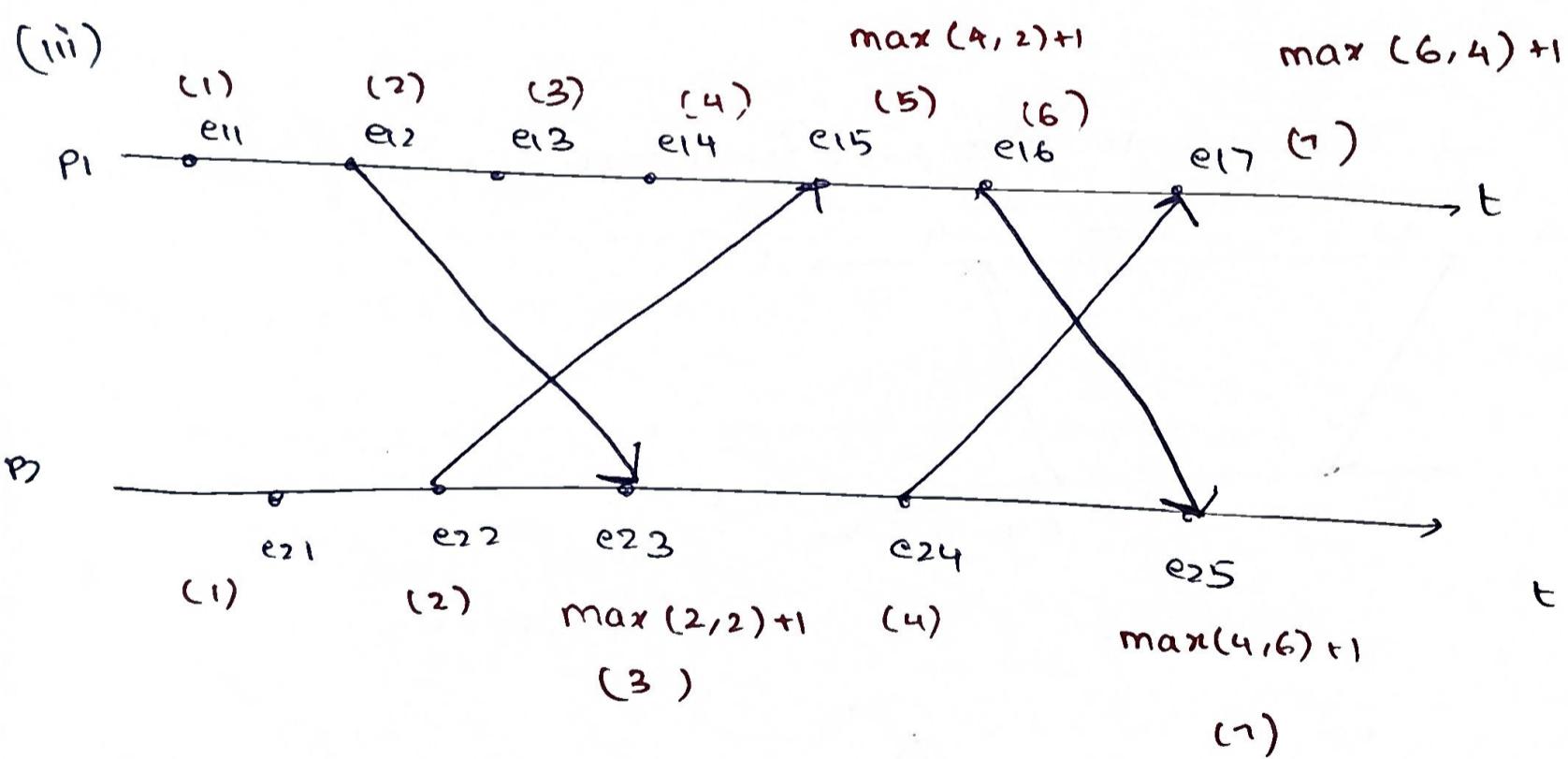
Ans



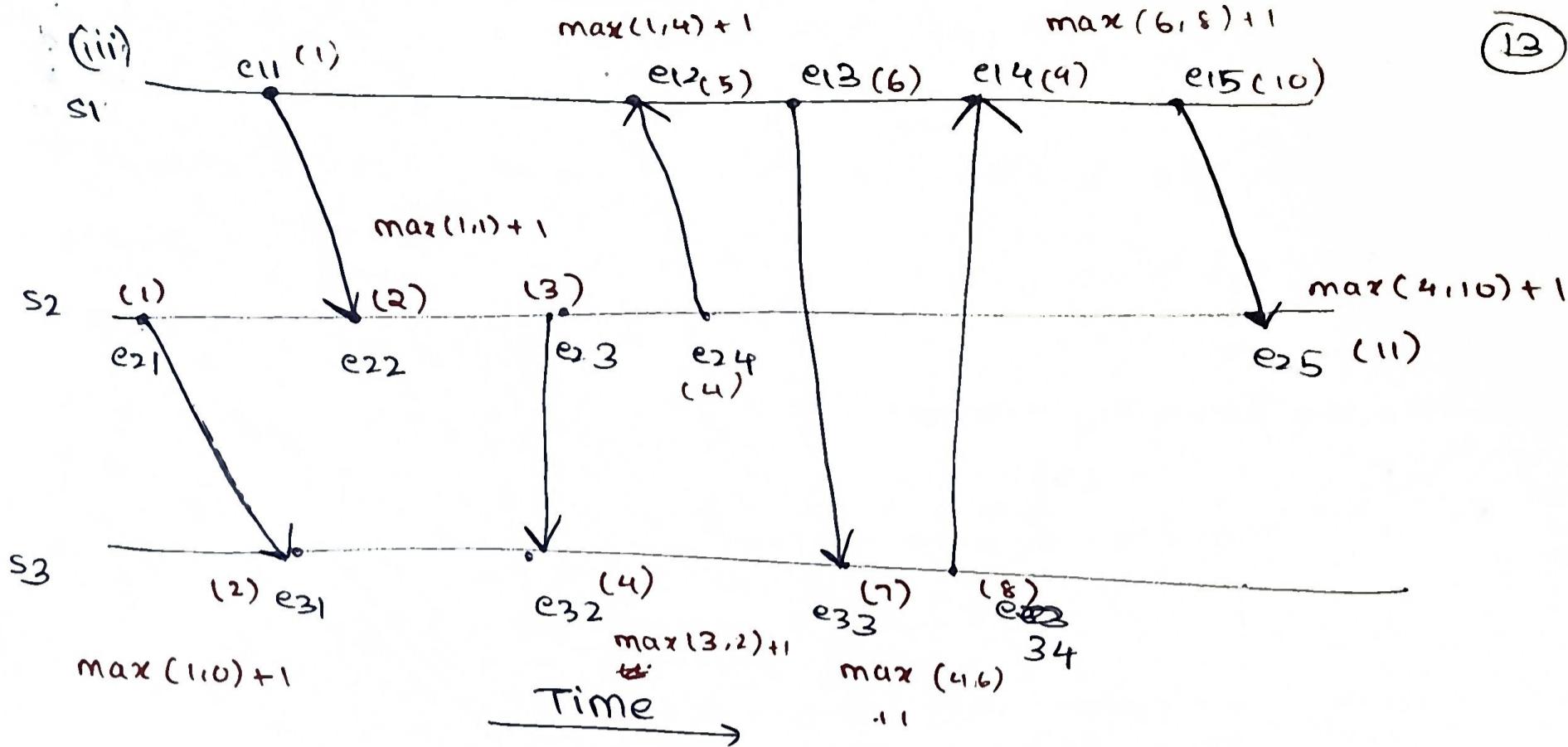
(ii)



(iii)



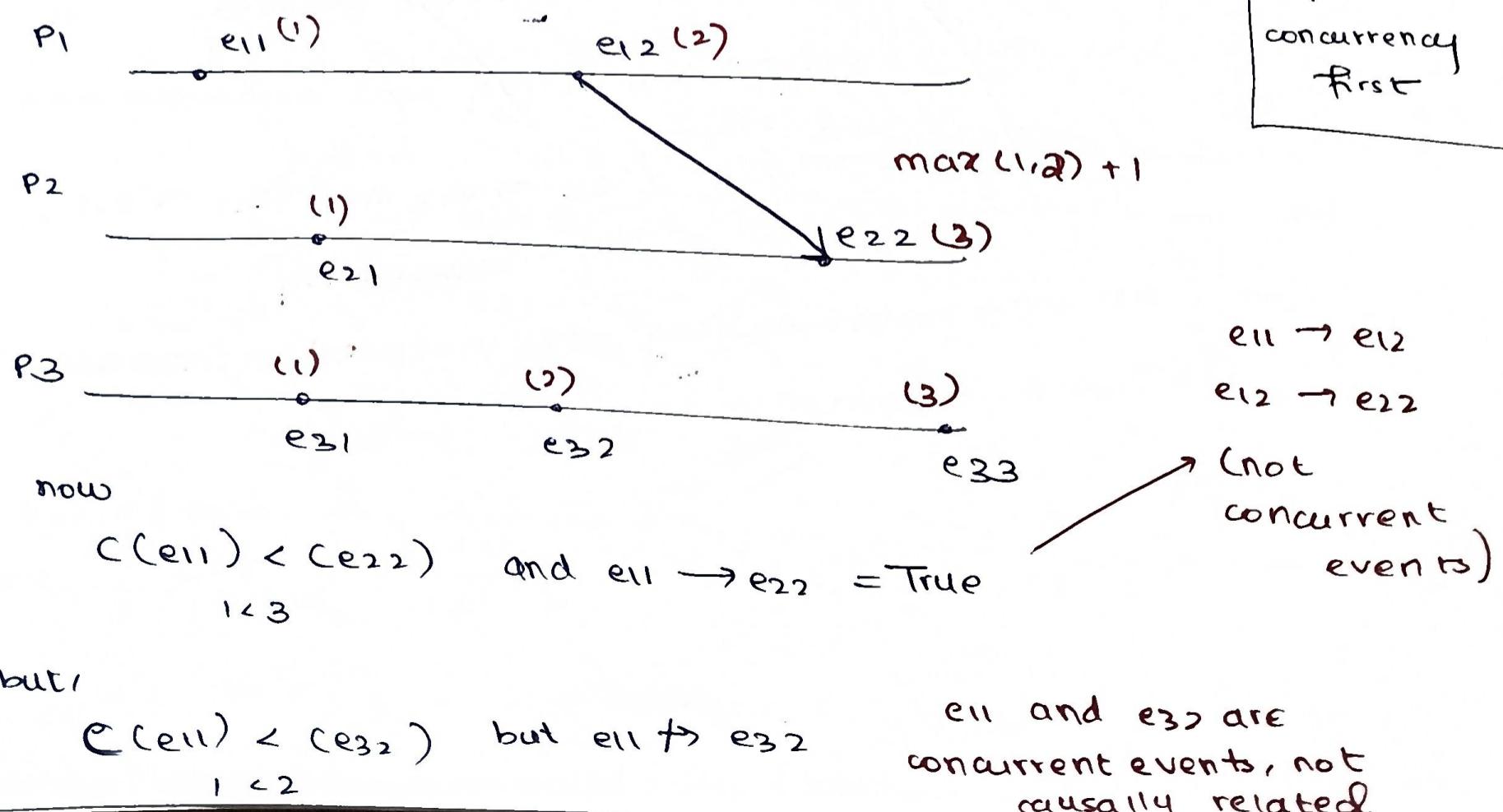
(13)



* Limitations of Lamport's Logical Clock

- With Lamport's logical clocks, if $a \rightarrow b$, then $c(a) < c(b)$
- However, if $c(a) < c(b)$, it is not necessarily true that $a \rightarrow b$
 i.e. **CANNOT DETERMINE CAUSALITY FROM TIMESTAMPES**

Consider the following example:

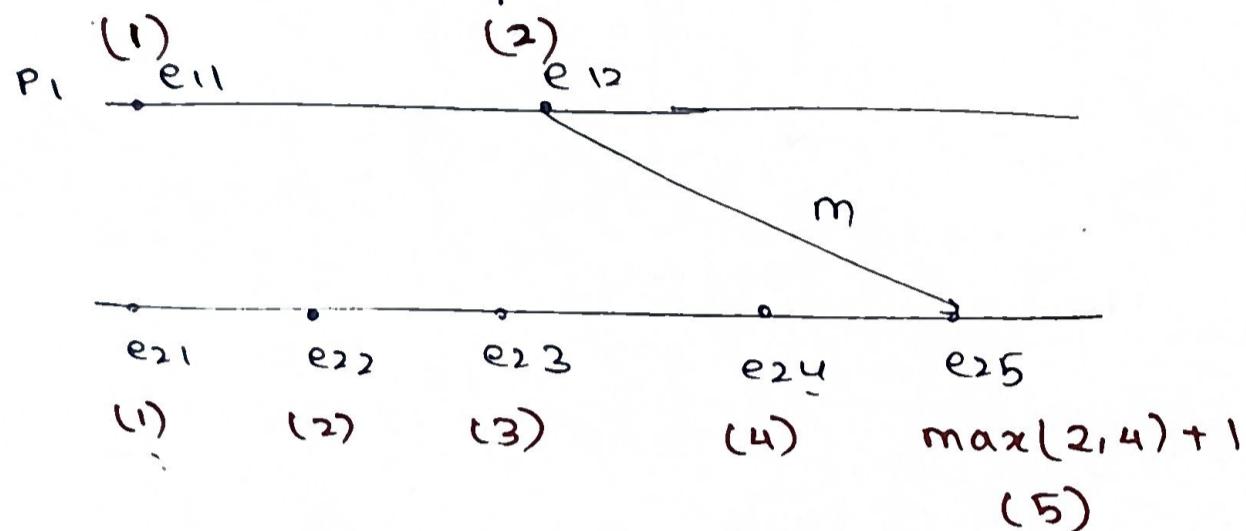


* Concurrency of Events

If $a \rightarrow b$ and $b \rightarrow a$

then all b

Consider the following example:



Here, $e_{11} \parallel e_{21}, e_{22}, e_{23}, e_{24}$

$e_{12} \parallel e_{21}, e_{23}, e_{24}$

$e_{11} \nparallel e_{12}, e_{25}$

$e_{12} \nparallel e_{25}$

Determine concurrency by checking for causality first

$e_{11} \rightarrow e_{12}$ (same process)

$e_{12} \rightarrow e_{25}$ (msg sent across processes)

$e_{11} \rightarrow e_{25}$ ($e_{11} \rightarrow e_{12}$ and $e_{12} \rightarrow e_{25}$, transitivity)

All other combinations are independent and are hence concurrent

* Vector Clocks

→ maintains a vector of values for every event that happens in all processes

→ In vector clocks, if $a \rightarrow b$

$$\forall k \quad c_a[k] < c_b[k]$$

and if $a \not\rightarrow b$

$$\exists k \quad c_a[k] \neq c_b[k]$$

Rules to calculate timestamps

$$c_i[k] = c_i[k] + 1 \quad \text{if } i=k \quad (\text{same process})$$

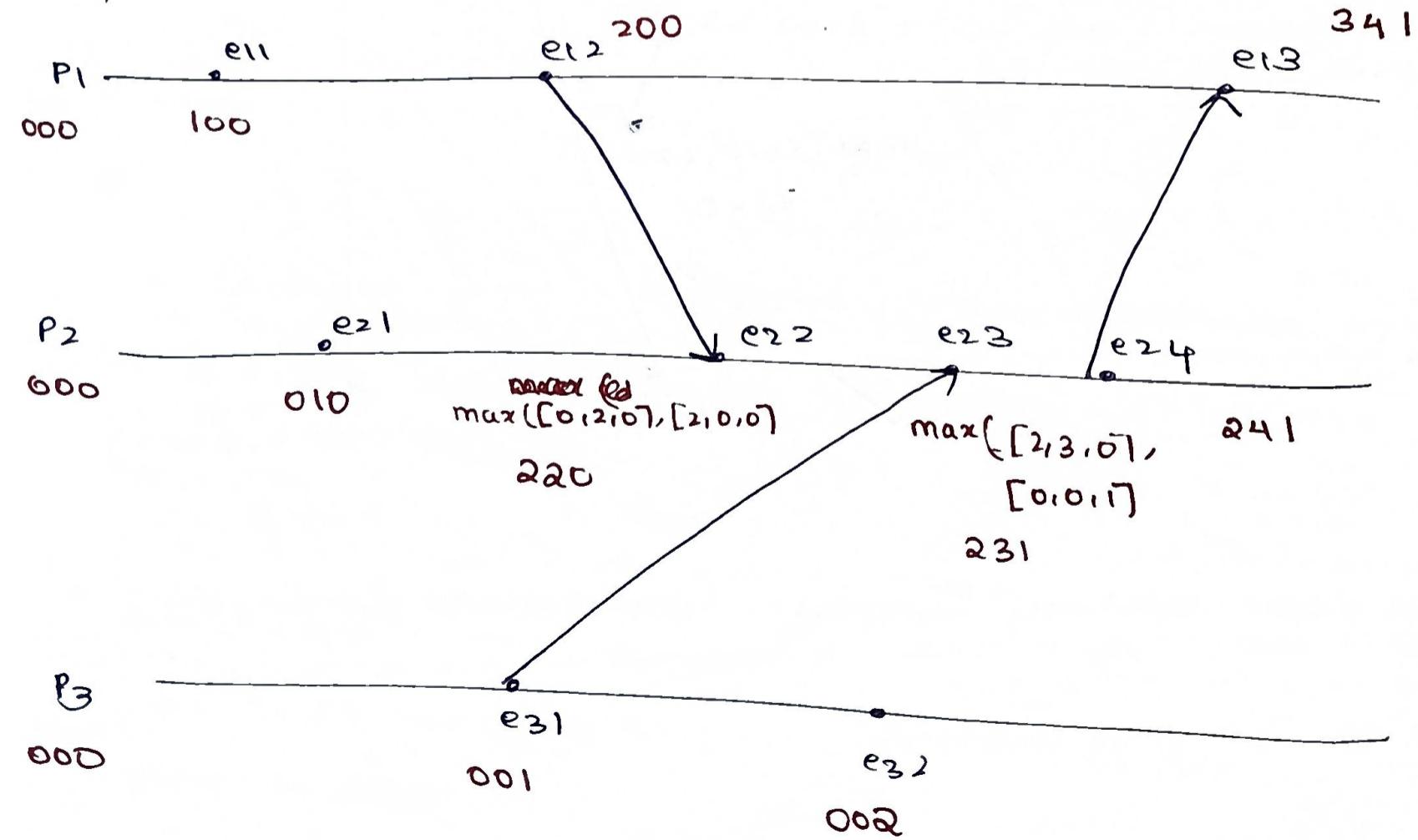
$$\forall c_j[k] = \max(c_j[k], e_m[k]) \quad (\text{across processes})$$

Examples

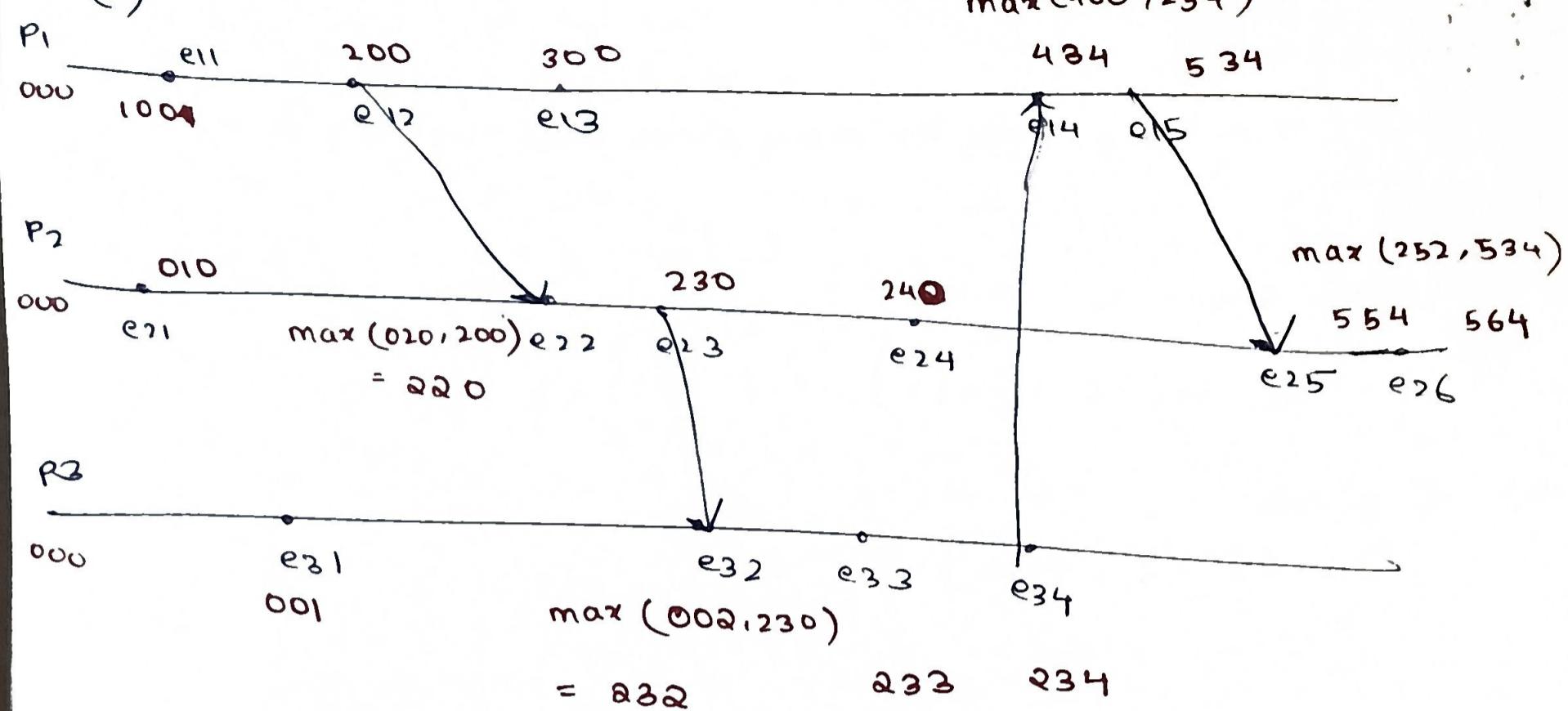
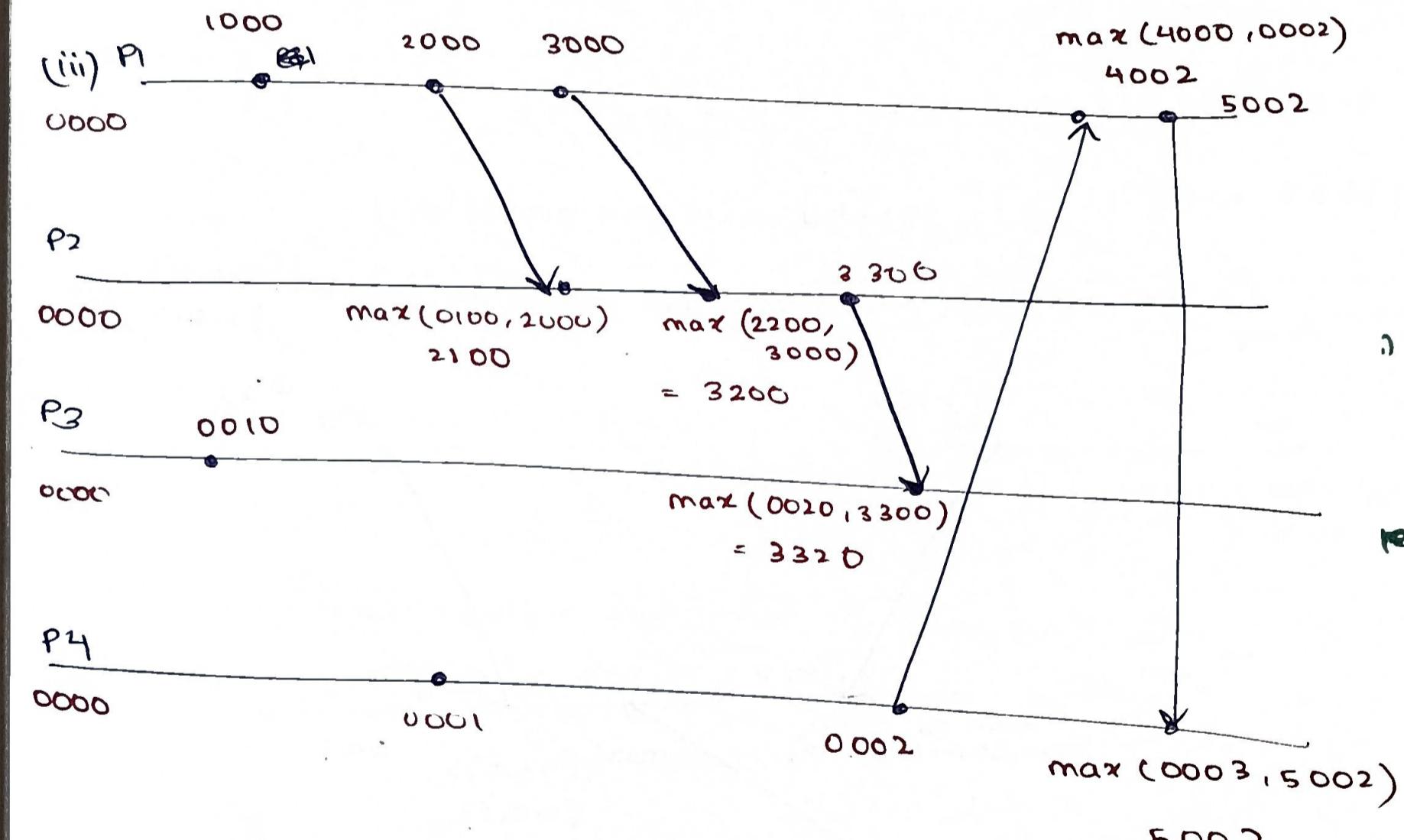
Calculate vector clock timestamps

$$\max([3, 0, 0], [2, 4, 7])$$

(1)



(ii)

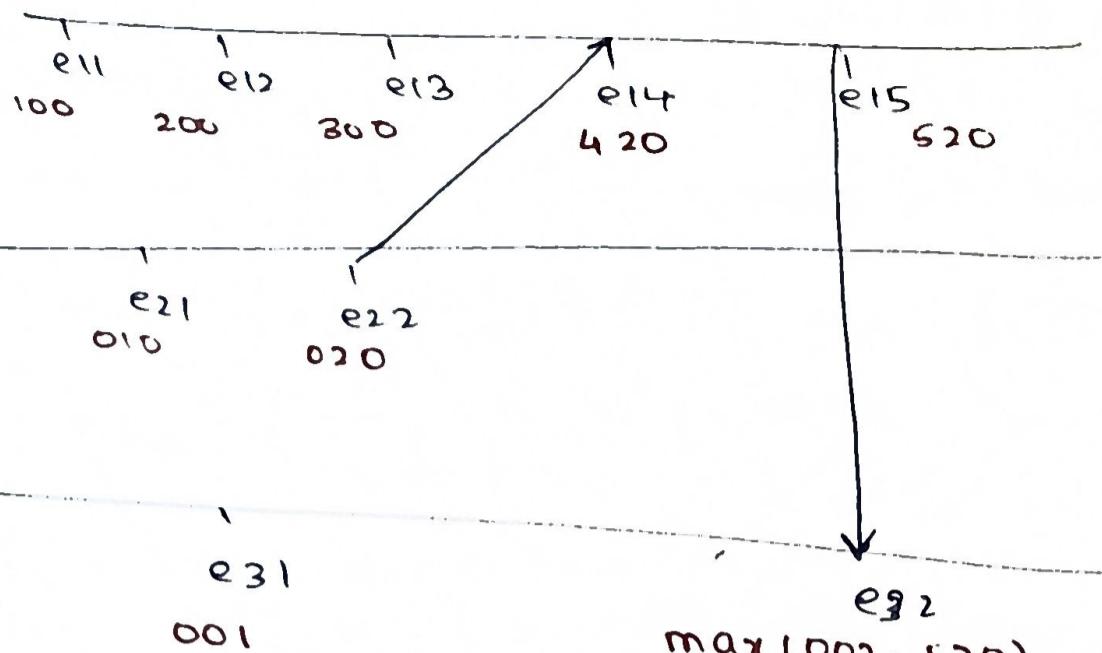
(iii) P₁

* Vector @locks resolving Lamport's logical @lock limitations

Consider the following scenario:

max(400, 020)

P1

consider e_{13} and e_{31}

(300) (001)

 $3 > 0$ $0 \leq 0$ $0 \leq 1$ \Rightarrow we cannot say that $e_{13} \rightarrow e_{31}$ \Rightarrow they are concurrent events

11104

consider ~~e_{21}~~ and e_{32}

(420) (522)

 $4 < 5$ $2 \leq 2$ $\Rightarrow e_{14} \rightarrow e_{32}$ $0 \leq 2$

→ Vector clocks resolve the issue of Lamport's Logical clock,

because by looking at the timestamp, we can determine

causality

* Drawbacks of vector clocks

→ need to maintain memory units for tracking all events in all processes as vectors.

→ The total no. of processes may not be known in advance.

Applications

1. Distribute debugging
2. causal distributed shared memory
3. Establish global breakpoints

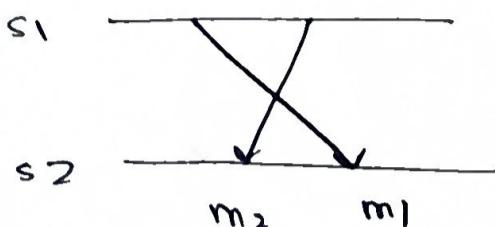
* Causal Ordering of Messages

If $\text{send}(m_1) \rightarrow \text{send}(m_2)$ then

$\text{rec}(m_1) \rightarrow \text{rec}(m_2)$



ideal behaviour



non-ideal behaviour

A. Causal Ordering of Messages in Broadcast Condition

Sender - i

increment $c_i[i]$

$t_m = c_i[i]$

broadcast (t_m)

Receiver - j

Case 1 - Broadcasts are not mixed with other concurrent sends / concurrent broadcasts

Case 2 - Concurrent broadcasts

Case 3 - missing receives,

single buffer item what the rules mean

case 4: multiple buffer items for a single system

1. If $c_j[i] = t_m[i] - 1$

2. If $\forall k c_j[k] \geq t_m[k]$, where $k \neq i$

IF ① and ② satisfy, deliver

else buffer the message

Rule 1: Have I

missed a message from the current sender?

Rule 2: Have I missed a message from any other sender other

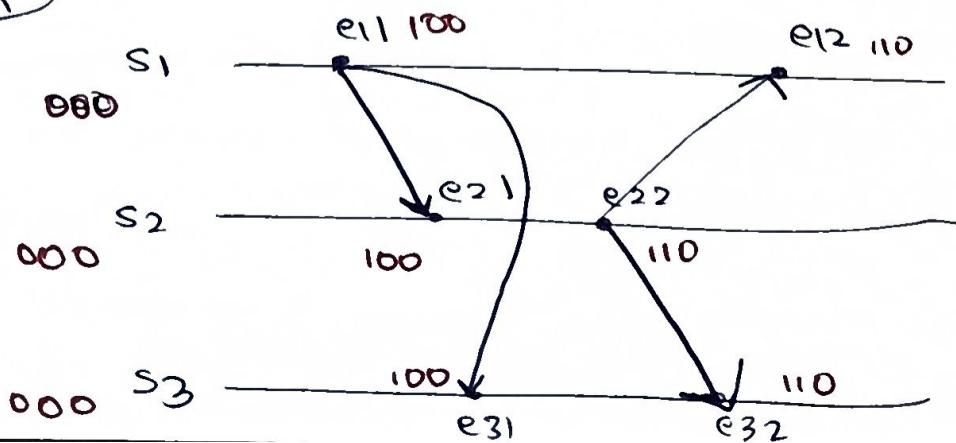
than the current sender,

following broadcast scenario

Example 1

Type
Case 1
mm

Apply the causal ordering of messages on the following broadcast scenario



(i) at e_{21}

$c_2 = 000$

for that process

$t_m = 100$

at sender's side

Rule 1

$c_j[i] = t_m[i] - 1$

$\Rightarrow 0 = 1 - 1 \quad \checkmark \text{satisfied}$

Rule 2

$c_2[2,3] \geq t_m[2,3]$

$a \geq 0$

$b \geq 0$

 $\checkmark \text{satisfied}$ calculate e_{21} as

$\max[(0,0,0), (1,0,0)]$

 \Rightarrow deliver the message \Rightarrow update the vector clock

$e_{21} = 100$

For the vector clock update, only increment for ~~the process alone~~ the sender's process

i.e $(0,0,0) \rightarrow (1,0,0)$

This is because we should consider only communicating events and not computing / processing events.

(ii) at e_{31}

$c_3 = 000$

$t_m = 100$

Rule 1

$c_j[i] = t_m[i] - 1$

$\Rightarrow 0 = 1 - 1 \quad \checkmark \text{satisfied}$

(-X-) The most important thing is - don't

increment the recipient's timestamps when processing

 c_0 or t_m . This is how

this differs from the

OG vector

clock algo.

Rule 2

$c[2,3] \geq t_m[2,3]$

$a \geq 0 \text{ and } b \geq 0 \Rightarrow \text{satisfies } \checkmark$

 \Rightarrow deliver the message \Rightarrow update the vector clock

$e_{31} = 100$

(ii) $e_{22} = 110$ (basic increment)

(iv) at e_{12} $c_1[100] \text{ tm } [110]$

Rule 1 $c_1[0] == \text{tm}[0] - 1$ \curvearrowright satisfies
 $0 == 1 - 1$

Rule 2 $c_1[1,3] \geq \text{tm}[1,3]$

$$(1,0) \geq (1,0)$$

\Rightarrow satisfies \curvearrowright

\rightarrow deliver the message

\rightarrow update the vector clock $e_{12} = 110$

(v) at e_{32}

$$c_3[100] \text{ tm } [110]$$

Rule 1 $c_3[0] == \text{tm}[0] - 1$
 $0 == 1 - 1$ \curvearrowright satisfies

Rule 2 $c_3[1,3] \geq \text{tm}[1,3]$
 $(1,0) \geq (1,0)$ \rightarrow satisfies

\rightarrow deliver the message

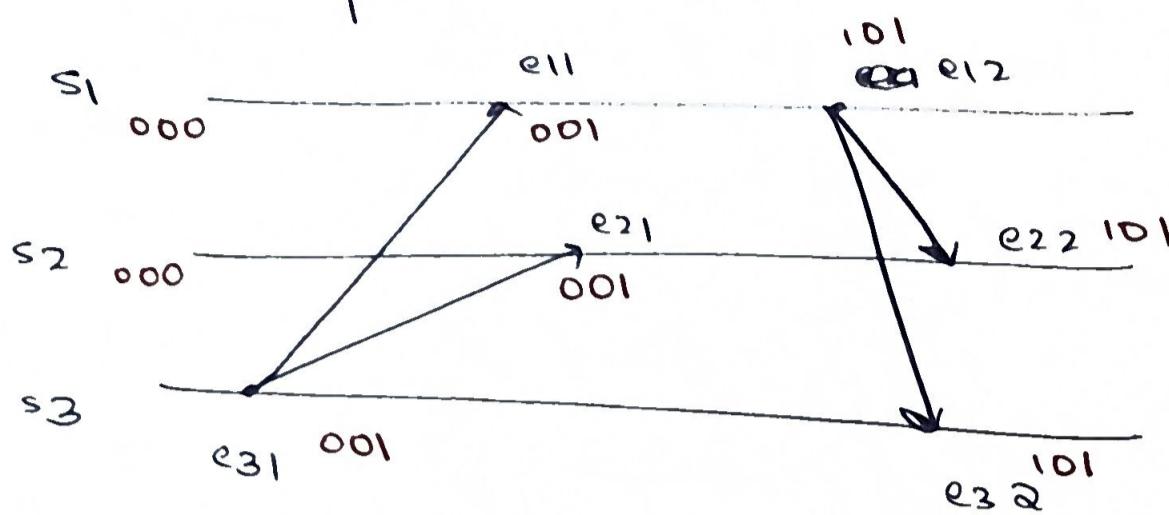
\rightarrow update the vector clock $e_{32} = 110$

Example 2

Type case 1

Apply the causal ordering of messages on

the following broadcast scenario

(i) at e_{11} $c_1[000] \geq tm[001]$ **Rule 1**

$$c_1[3] = tm[3] - 1$$

$$0 = 1 - 1 = 0 \quad \text{satisfies } \checkmark$$

Rule 2

$$c_1[110] \geq tm[110]$$

$$(0,0) \geq (0,0) \quad \text{satisfies}$$

 \Rightarrow deliver message

$$\Rightarrow e_{11} = 001$$

(ii) at e_{21}

$c_2[000]$ at receiver $tm[001]$ message

Rule 1

$$c_2[3] = tm[3] - 1$$

$$0 = 1 - 1 = 0 \quad \text{satisfies } \checkmark$$

Rule 2

$$c_2[110] \geq tm[110]$$

$$(0,0) \geq (0,0) \quad \text{satisfies } \checkmark$$

\rightarrow deliver msg
 $\rightarrow e_{21} = 001$

~~(iii) $e_{12} = 101$~~

(iii) $e_{12} = 101$

(iv) at e_{22} $c_2[001] \geq tm[101]$

Rule 1 $c_2[1] = tm[1] - 1$

$$0 = 1 - 1 \quad \text{satisfies } \curvearrowright$$

Rule 2 $c_2[2,3] \geq tm[2,3]$

$$(0,1) \geq (0,1) \quad \text{satisfies}$$

\Rightarrow deliver

$e_{22} = 101$

(v) at e_{32} $c_3[001] \geq tm[101]$

Rule 1 $c_3[1] = tm[1] - 1$

$$0 = 1 - 1 \quad \Rightarrow \text{satisfies } \curvearrowright$$

Rule 2 $c_3[2,3] \geq tm[2,3]$

$$(0,1) > (0,1) \quad \curvearrowright \text{satisfies}$$

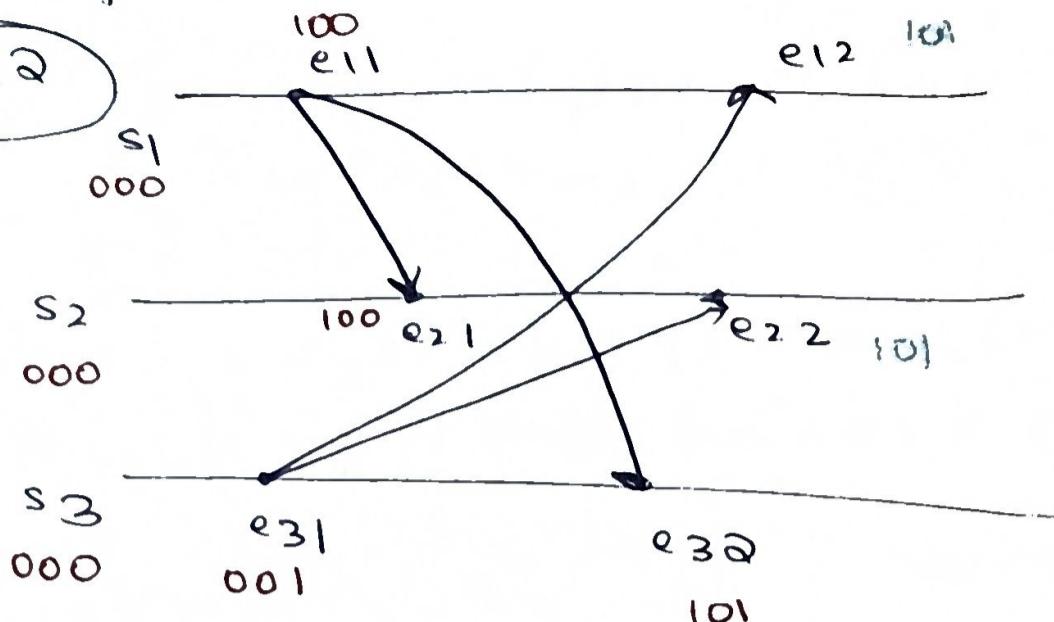
\Rightarrow deliver

$e_{32} = 101$

Example 3

Type : Base 2

broadcast are concurrent



(i) at e_{21} : $c_2[000]$ $tm[100]$

Rule 1 $c_2^r[1] = tm[1] - 1$
 $0 = 1 - 1 \Rightarrow \text{satisfies}$

Rule 2 $c_2[2,3] \geq tm[2,3]$
 $(0,0) \geq (0,0) \Rightarrow \text{satisfies}$

$\Rightarrow \text{deliver}, e_{21} = 100 \leftarrow e_{21} = \max(000, 100)$

~~000, remember don't do~~

(ii) at e_{32} : $c_3[001]$ $tm[100]$

Rule 1 $c_3^r[1] = tm[1] - 1$
 $0 = 1 - 0 \Rightarrow \text{satisfies}$

Rule 2 $c_3[2,3] \geq tm[2,3]$
 $(0,1) \geq (0,0) \Rightarrow \text{satisfies}$

$\rightarrow \text{deliver}$

$e_{32} = \max(001, 100) = 101$

(iii) at e_{12} $c_1[100]$ $tm = [001]$

Rule 1 $c_1[3] = tm[3] - 1$
 $0 = 1 - 1 \Rightarrow \text{satisfies}$

Rule 2 $c_1[1,2] \geq tm[1,2]$

$(1,0) \geq (0,0) \Rightarrow \text{satisfies}$

$\Rightarrow \text{deliver}$

$\Rightarrow e_{12} = \max(100, 001) = 101$

(iv) at e_{22} $c_2[100]$ $tm = [001]$

Rule 1 $c_2[3] = tm[3] - 1$
 $0 = 1 - 1 = 0 \Rightarrow \text{satisfies}$

Rule 2 $c_2[1,2] \geq tm[1,2]$

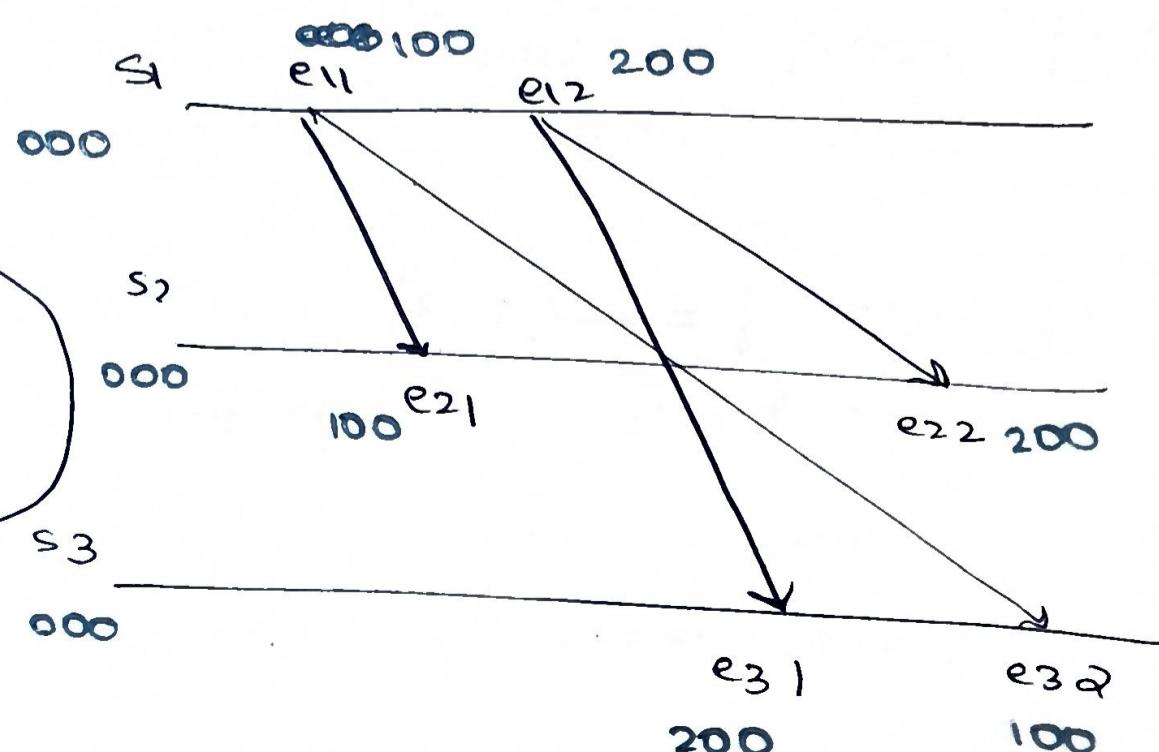
$(1,0) \geq (0,1) \Rightarrow \text{satisfies}$

$\Rightarrow \text{deliver}$

$e_{22} = \max(100, 001)$
 $= 101$

Example 4

Type: case 3
includes
missing receives



$$(i) \quad \underbrace{e_2}_{\text{eq}} \quad C_2 = [000] \\ t_m = [100]$$

Rule 1 $C_2[0] = t_m[1] - 1$
 $0 = 1 - 1 \Rightarrow \text{satisfies}$

Rule 2 $C_2[2,3] \geq t_m[2,3]$
 $(0,0) \geq (0,0) \Rightarrow \text{satisfies}$

\Rightarrow deliver

$$\Rightarrow e_{21} = \max(000, 100) = 100$$

$$(ii) \quad \underbrace{e_2}_{\text{m}} \quad C_2 = [100] \\ t_m = [200]$$

Rule 1 $C_2[1] = t_m[1] - 1$
 $1 = 2 - 1 \Rightarrow \text{satisfies}$

Rule 2 $C_2[2,3] \geq t_m[2,3]$
 $(0,0) \geq (0,0) \Rightarrow \text{satisfies}$

$$\Rightarrow e_{22} = \max(100, 200) = (2,00)$$

$$(iii) \quad \underbrace{e_3}_{\text{m}} \quad C_3 = [000] \\ t_m = [200]$$

Rule 1 $C_3[1] = t_m[1] - 1$
 $0 \neq 2 - 1 \Rightarrow \text{error}$

\Rightarrow there is a message that has been missed

\Rightarrow push (sender, tm) onto S₃'s buffer (queue)

(1, 200)

$$(iv) \text{ e}_{32} \quad c_3 = [000]$$

$$tm = [100]$$

Rule 1 $c_3[1] == tm[1] - 1$
 $0 == 1 - 1 \Rightarrow \text{satisfies}$

Rule 2 $c_3[2,3] \geq tm[2,3]$
 $(0,0) \geq (0,0) \Rightarrow \text{satisfies}$

\Rightarrow deliver

$$\Rightarrow e_{32} = \max(000, 100) = 100$$

Now that there has been a successful case, handle the case in the buffer

$$(v) \text{ e}_{31} \quad c_3 = [100]$$

$$tm = [200]$$

Rule 1 $c_3[1] == tm[1] - 1$
 $1 == 2 - 1 \Rightarrow \text{satisfies}$

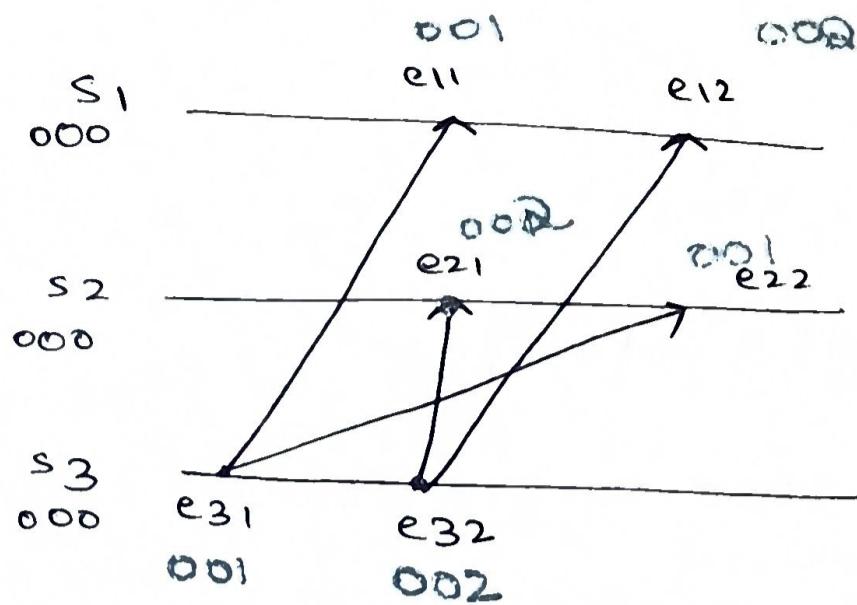
Rule 2 $c_3[2,3] \geq tm[2,3]$
 $(0,0) \geq (0,0) \Rightarrow \text{satisfies}$

\Rightarrow deliver

$$\Rightarrow e_{31} = \max(100, 200)$$
$$= (200)$$

Example 5

Type - case 3



(i) e_{11} $c_1 = [000]$
 $tm = [001]$

Rule 1 $c_1[3] == tm[3] - 1 \quad \checkmark$

Rule 2 $c_1[1..2] \geq tm[1..2] \quad \checkmark$

$$\Rightarrow e_{11} = \max(000, 001) = 001$$

(ii) e_{12} $c_1 = [001]$
 $tm = [002]$

Rule 1 $c_1[3] == tm[3] - 1 \quad \checkmark$

Rule 2 $c_1[1..2] \geq tm[1..2] \quad \checkmark$

$$\Rightarrow e_{12} = \max(001, 002) = 002$$

(iii) e_{21} $c_2 = [000]$
 $tm = [002]$

Rule 1 $c_2[3] == tm[3] - 1 \quad X \text{ Fails}$
 $\Rightarrow \text{missing receive}$

push onto buffer

$s_2 \boxed{(3, 002)}$

(N) e_{22} $c_2 = [000]$
 $t_m = [001]$

Rule1 $c_2[3] = -t_m[3] - 1 \checkmark$

Rule2 $c_2[1,2] \geq t_m[1,2] \checkmark$

$$\rightarrow e_{22} = \max(000, 001) = 001$$

(v) e_{21B} $c_2 = [001]$
 $t_m = [002]$

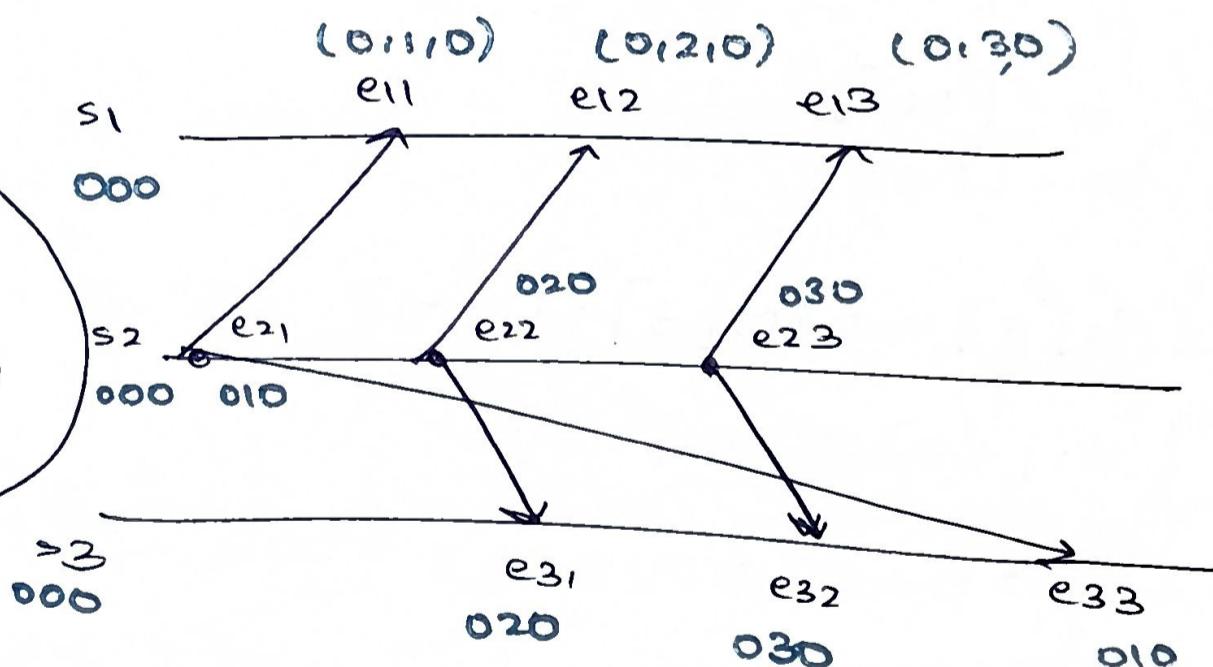
Rule1 $c_2[3] = t_m[3] - 1 \checkmark$

Rule2 $c_2[1,2] \geq t_m[1,2] \checkmark$

$$\rightarrow e_{21B} = \max(001, 002) = 002$$

Example 6

Type: Case 4
Multiple buffer items in a single system



Ans e_{11}, e_{12} & e_{13} will not have any problem of messages received out of order

consider (i) e_{31} $c_3 = [000]$
 $t_m = [020]$

$c_3[2] = -t_m[2] - 1 \times \rightarrow$ push $(2, 020)$ onto buffer

(ii) Try e_{32} $c_3 = [000]$

$$tm = [030]$$

$c_3[2] == tm[2] - 1 \times \rightarrow$ push $(2, 030)$ onto the buffer

(iii) Try e_{33} $c_3 = [000]$

$$tm = [010]$$

$$c_3[2] == tm[2] - 1 \checkmark$$

$$c_3[1, 3] \geq tm[1, 3] \checkmark$$

$$(010) > (010)$$

\Rightarrow deliver

$$\Rightarrow e_{33} = \max(000, 010) = 010$$

(iv) By means of FIFO, now try e_{31_B}

$$c_3[010]$$

$$tm = [020]$$

$$c_3[2] == tm[2] - 1 \checkmark$$

$$c_3[1, 3] \geq tm[1, 3] \checkmark$$

$$(0,0) > (0,0)$$

\Rightarrow deliver

$$\Rightarrow e_{31_B} = \max(010, 020) = 020$$

(v) Now by means of FIFO, attempt e_{32_B}

$$c_3 = [020]$$

$$tm = [030]$$

$$c_3[2] == tm[2] - 1 \checkmark$$

$$c_3[1, 3] \geq tm[1, 3]$$

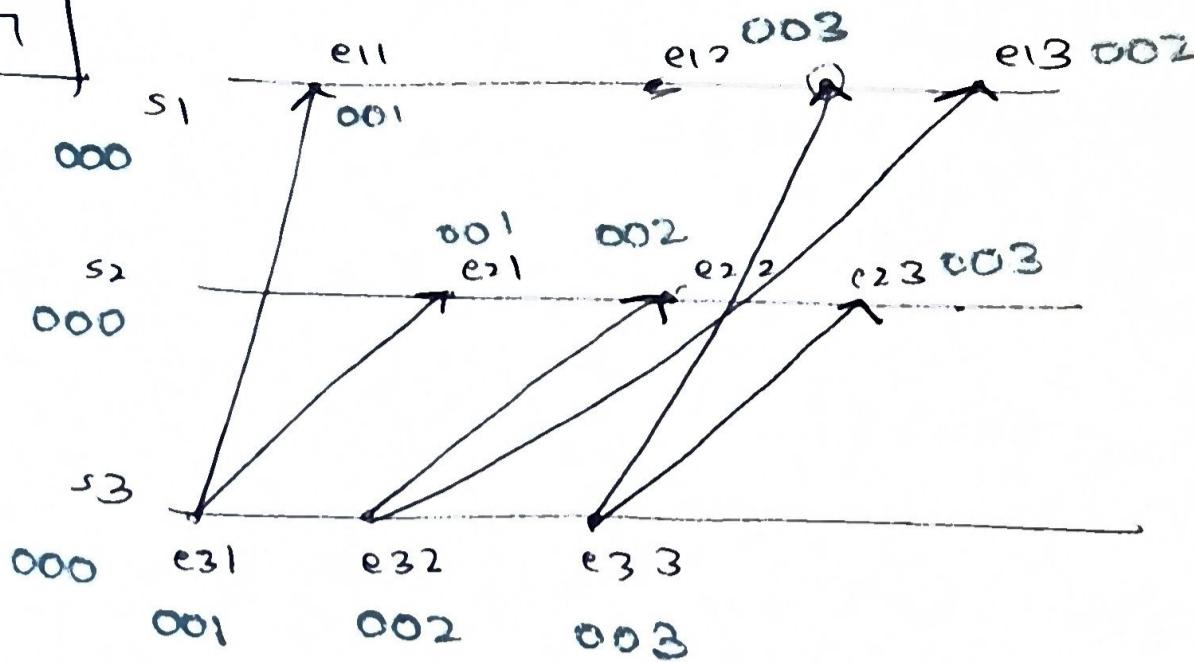
$$(0,0) > (0,0)$$

\Rightarrow deliver

$$\Rightarrow e_{32_B} = \max(020, 030) = 030$$

Example 7

(Type = case 3)



e_{21} and e_{11} are just regular updates

$$e_{11} = \max(001, 000) = 001$$

$$e_{21} = \max(001, 000) = 001$$

111th for e_{22} and e_{23}

$$e_{22} = \max(002, 001) = 002$$

$$e_{23} = \max(003, 002) = 003$$

(i) For e_{12} $c_1^s = [001]$
 $tm = [003]$

$$c_1[3] = tm[3] - 1 \quad X$$

Place (3, 003) in the buffer $\otimes S_1$

(ii) For e_{13} $c_1 = [001]$
 $tm = [002]$

$$c_1[3] = tm[3] - 1 \quad \checkmark$$

$$c_1[1,2] \geq tm[1,2] \quad \checkmark$$

\Rightarrow deliver $\Rightarrow e_{13} = \max(002, 001) = 002$

(iii) From the buffer e_{12} B

$c_1 = [002]$
 $tm = [003]$

$$c_1[3] = tm[3] - 1 \quad \checkmark$$

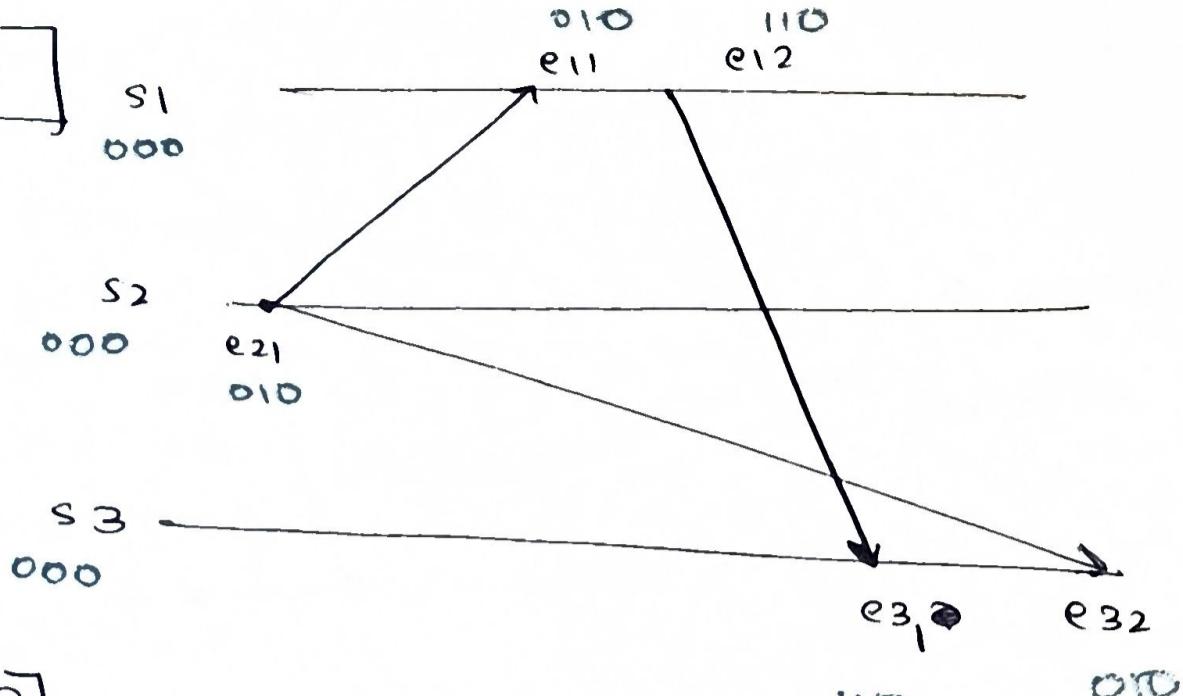
$$c_1[1,2] \geq tm[1,2] \quad \checkmark$$

\Rightarrow deliver $\Rightarrow e_{12} = \max(002, 003) = 003$

Example 8

Type : (case 5)

Missed from
other sender



(i) For e_{11} $c_1 = [000]$
 $tm = [010]$

$$c_1[0] == tm[0] - 1 \quad \checkmark$$

$$c_1[1,3] \geq tm[1,3] \quad \checkmark$$

$$e_{11} = 010$$

(ii) For $e_{12} = 110$

(iii) For e_{31} $c_3 = [000]$
 $tm = [110]$

$$c_3[1] == tm[1] - 1 \quad \checkmark$$

$$c_3[2,3] \geq tm[2,3]$$

$$(0,0) \geq (1,0) \times$$

$\rightarrow S_3$ has missed a msg from a system that is not the current sender.

push $(1,110)$ onto $buff_0$

(iv) For e_{32} $c_3 = [000]$
 $tm = [010]$

$$c_3[2] == tm[2] - 1 \quad \checkmark \quad \Rightarrow e_{32} = \max(000, 010) = 010$$

$$c_3[1,3] \geq tm[1,3] \quad \checkmark$$

(v) Now e_{31} $c_3 = [010]$
 $tm = [110]$

$$\Rightarrow e_{31} = \max(010, 110) \\ = 110$$

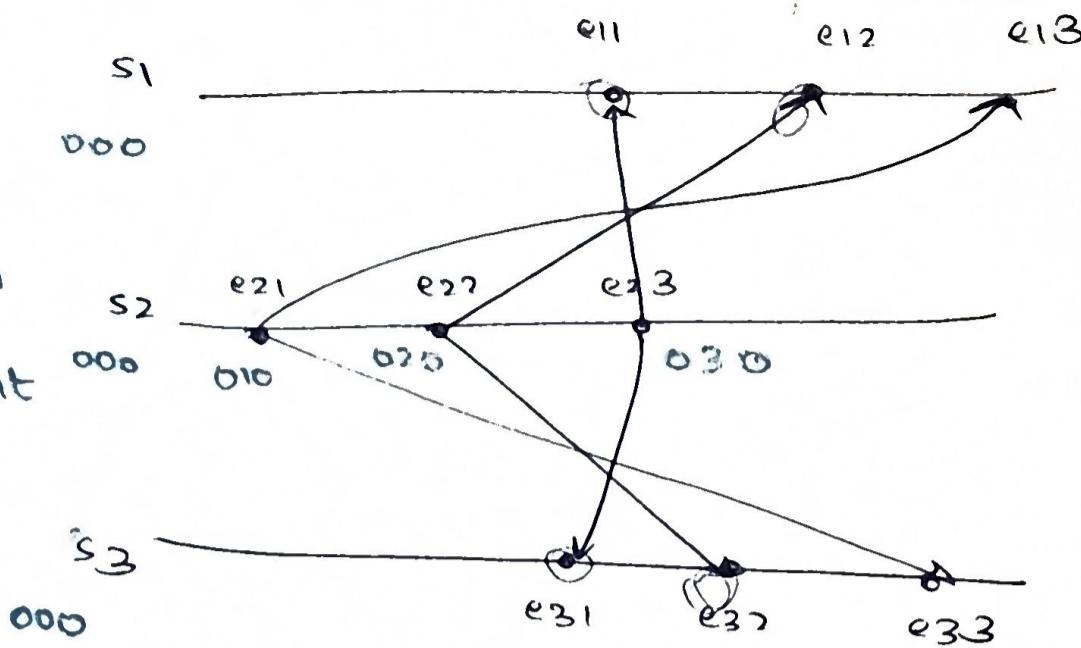
$$c_3[1] == tm[1] - 1 \quad \checkmark$$

$$c_3[2,3] \geq tm[2,3] \quad (110) \geq (110) \quad \checkmark$$

Example 9

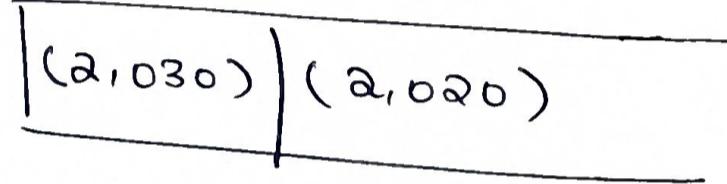
Type : Case 6

multiple elements in
buffers of different
systems



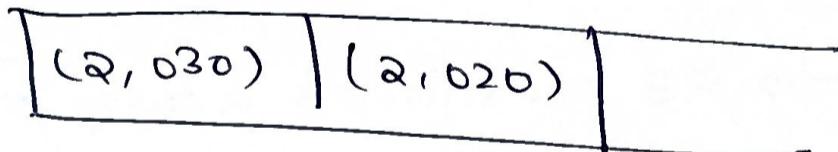
In this case e_{31} and e_{32} will fail since the message from e_{21} has not reached e_{33}

push onto buffer
of s_3



by e_{11} and e_{12} will fail because the message from e_{21} has not reached e_{13}

push onto buffer of s_1



process $e_{13} \Rightarrow 010$

and $e_{33} \Rightarrow 010$

then ~~pop~~ dequeue each process from the queues of s_1 & s_3

* Causal Ordering of Messages - Multicast

→ uses an extra data structure : Vector v_i

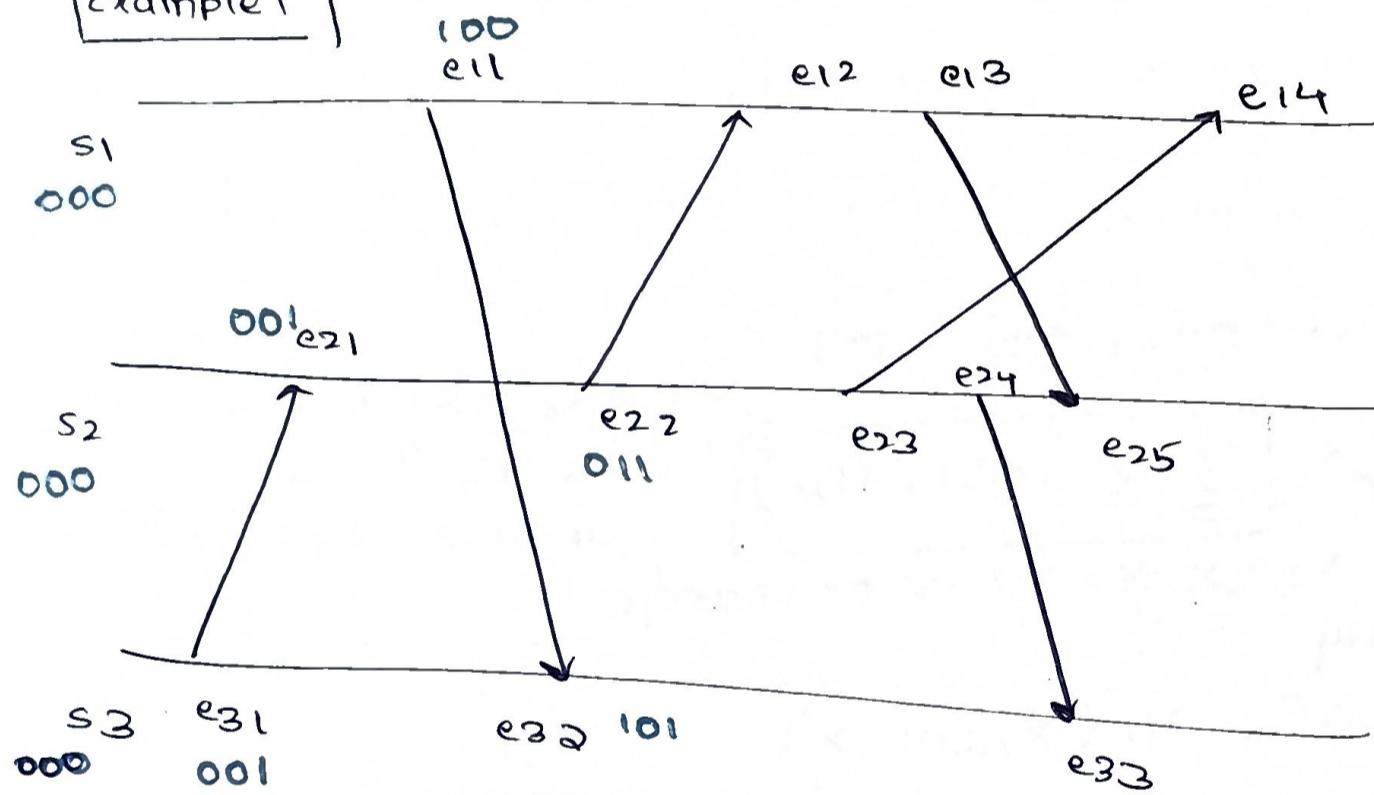
→ Algorithm for send

```

update clock  $c_i$ 
send  $(c_i, v_i)$  to j
update  $v_i$ 

```

Example 1



(i) Initially

① at e_{31}

$c_3(0,0,0)$ $v_3(X,X,X)$

(ii) update clock = $c_3(0,0,1)$

Send

(iii) send $c_3(0,0,1)$ $v_3(X,X,X)$ to S_2

(iv) update v_3 :

$v_3(X, 001, X)$

sent to S_2

sent to S_2 ,
update second pos

② at e_{21}

(i) Initially $c_2(0,0,0)$, $v_2(X,X,X)$

Recv

(ii) Receive $c_3(0,0,1)$, $v_3(X,X,X)$ update received

ump

(iii) $\frac{e_{21} = \max(000, 001)}{\text{update}} = 001$ → clock value at
 $(X, X, X) \in | V_2 = (X, 001, X) |$ place received clock value receiver's own

③ at e_{11} : (i) Initially

$$c_1(0,0,0) \quad v_1(x, x, x)$$

(ii) update clock $c_1(1,0,0)$

send (i) send c, v to e_{32}

The msg is sent to

s_3 , update 3rd pos

(iv) update $v_1: (\cancel{x}, x, 100)$

compare with

sent to s_3

④ at e_{32}

(i) Currently

$$c_3(001) \quad v_3(x, 001, x)$$

recv (ii) receive c, v from e_{11}

$$(iii) e_{32} = \max(001, 100) = 101$$

(iv) update $v_3:$

$$v_3(x, 001, 101)$$

The msg is sent from
 $1 \rightarrow 3$

update v_3 at 3rd
pos

compare with $v_1(x, x, x) \Rightarrow$ no change

⑤ at e_{22} (i) Currently

$$c_2(001) \quad v_2(x, 001, x)$$

send (ii) update clock $= c_2(011)$

(iii) send c_2, v_2 to s_1

(iv) update $v_2: (011, 001, x)$

The msg is sent to from
 $2 \rightarrow 1$
update First pos

⑥ at e_{12} : (i) Currently

$$c_1(100) \quad v_1(x, x, 100)$$

recv

(ii) receive c_2, v_2 from s_2

The message is
sent from $2 \rightarrow 1$,
update 1

$$(iii) e_{12} = \max(100, 011) = 111$$

(iv) update v_1 : ~~prefer earlier~~

$$v_1(111, x, 100)$$

v_2 received : $(x, 001, x)$
update v_1 $(111, 001, 100)$

remember to add clock
value, received v_1 , old
current v

35

⑨ at e23

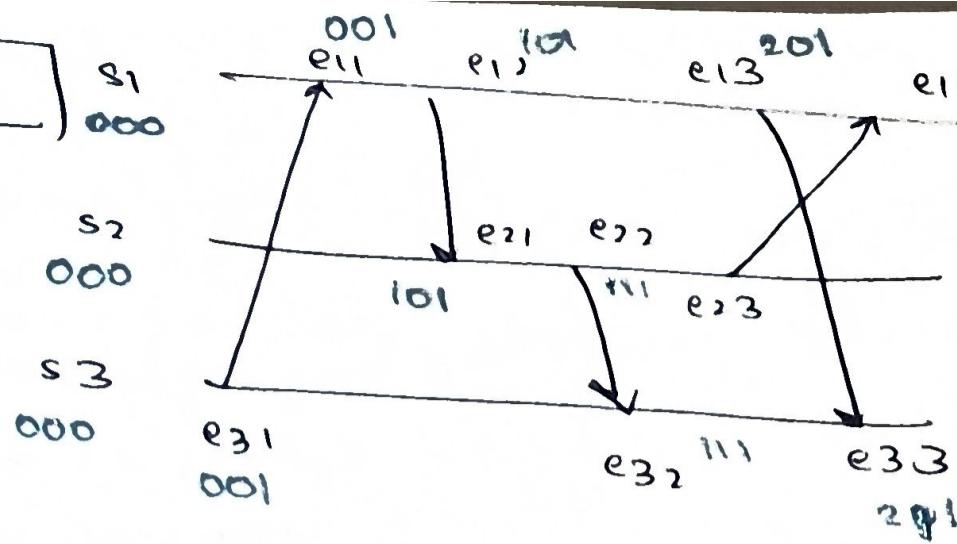
$c_2(011), v_2 (011, 001, x)$

$c_2(021), v_2 (021, 001, x)$

$s_2 \rightarrow s_1$ update s_1
with clock value

continue

Example 2



@ e₃₁

$c_3(0,0,0)$ $v_3(x,x,x)$
 update c, v $c_3(0,0,1)$, $v_3(001, x, x)$ $s_3 \rightarrow s_1$

send

@ e₁₁

$c_1(0,0,0)$ $v_1(x,x,x)$

receive here

recv

$c_1(0,0,1)$ $v_1(001, x, x)$ $s_3 \rightarrow s_1$

@ e₁₂

$c_1(0,0,1)$ $v_1(001, x, x)$

send

update c, v $c_1(1,0,1)$ $v_1(001, 101, x)$ $s_1 \rightarrow s_2$

@ e₂₁

$c_2(0,0,0)$ $v_2(x,x,x)$

recv

receive here
 $c_2(1,0,1)$ $v_2(001, 101, x)$ $s_1 \rightarrow s_2$

@ e₁₃

$c_1(101)$ $v_1(001, 101, x)$

send

$c_1(201)$ $v_1(001, 101, 201)$ $s_1 \rightarrow s_3$

@ e33

(3 (001) v3(001, ¹⁰¹~~01~~, ²¹~~11~~)

(37)

recv

c3 (001)

v3(001, 101, ²¹~~11~~)

s1 → s3

@ e22

c2 (101)

v2(001, 101, X)

send

c2 (1101)

v2(001, 101, ¹¹¹~~200~~)

s2 → s3

@ e32

c2 (001)

v3 (001, 101, 201)

recv

c2 (¹¹~~001~~)

v3 (001, 101, ¹¹¹~~201~~)

s2 → s3

@ e23

c2 (111)

v2 (001, 101, 111)

send

c2 (1121)

v2 (121, 101, 111)

s2 → s1

@ e14

c1 (201)

v1 (001, 101, 201)

recv

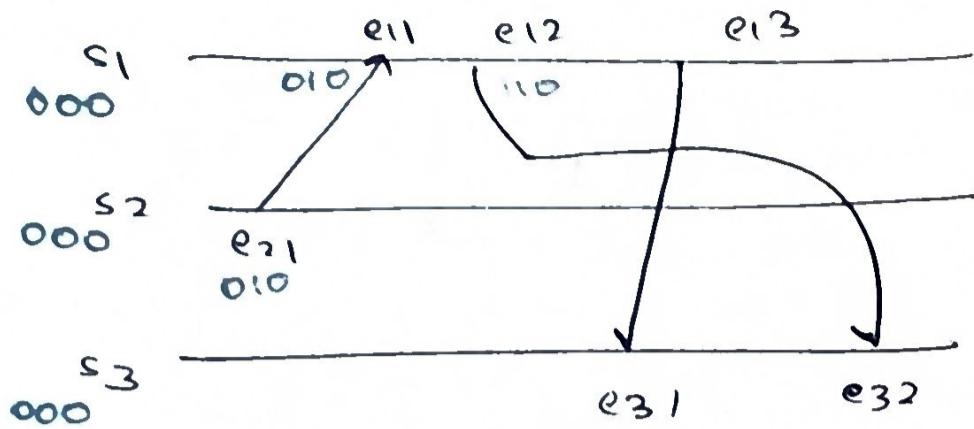
c1 (2121)

v1 (221, 101, 211)

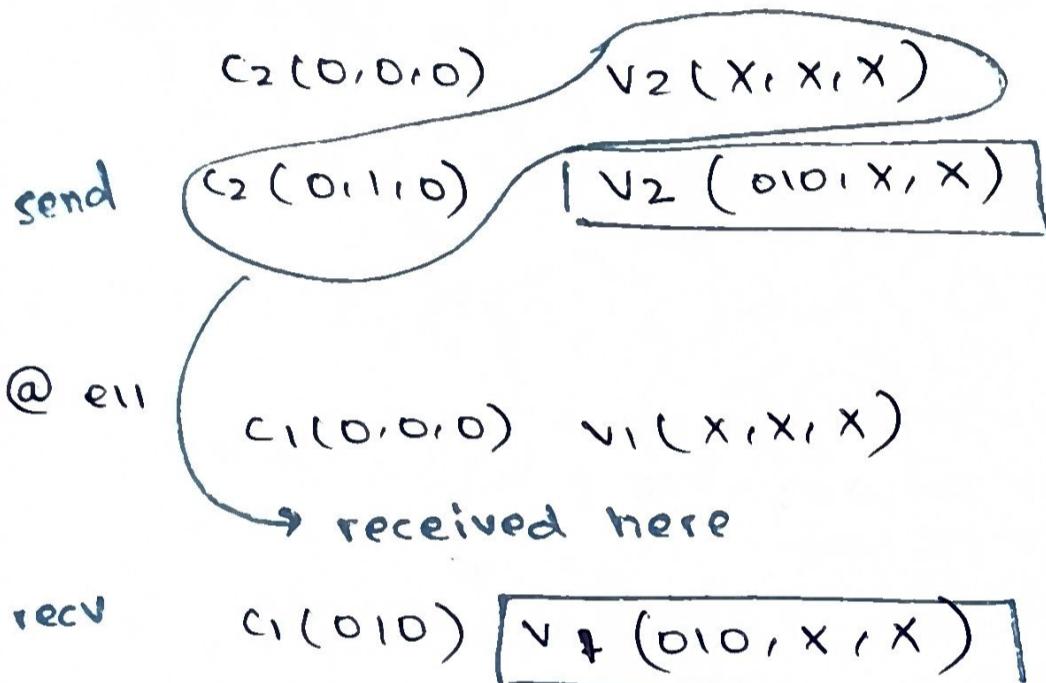
s2 → s1

Example 3

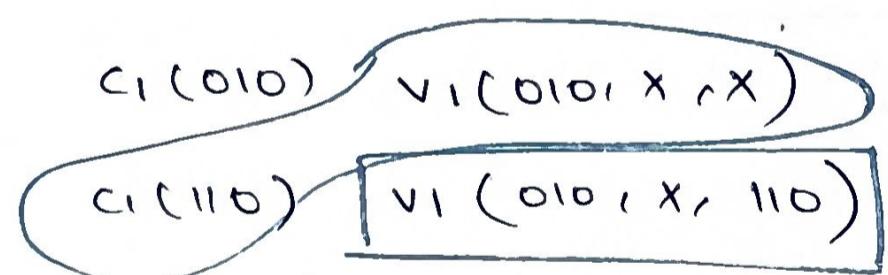
- with conflict



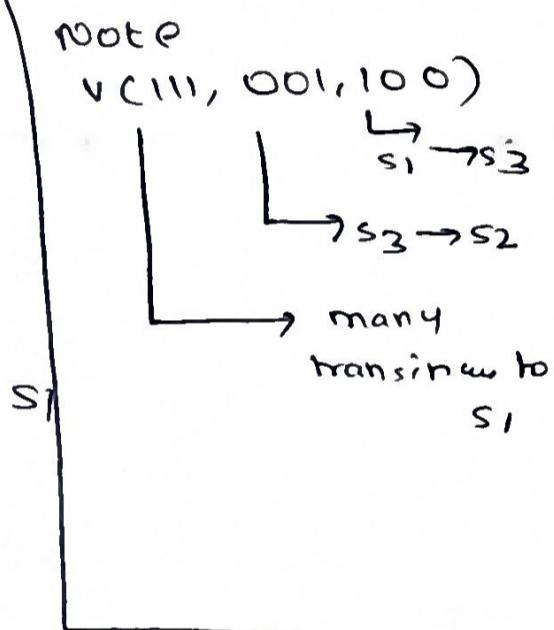
@e21



@e11



$s_2 \rightarrow s_1$

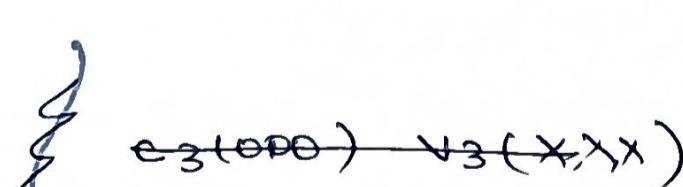


$s_2 \rightarrow s_1$

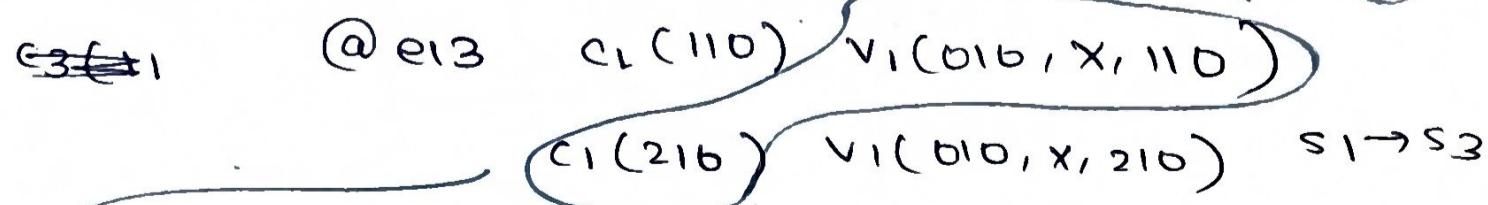
$s_1 \rightarrow s_3$

2

@e32



recv



@e31



$c_3(000)$ $v_3(x, x, x)$ → The receiver says nothing is received

but the received $v_1(010, x, 110)$ s_1 has sent to s_3

4

∴ push 210, (010, X, 110) onto the buffer

now @ e32:

$c_3(000)$ $v_3(XXX)$

→ recv from 2 here

$c_3(110)$ $v_3(010, X, 110)$ $s_1 \rightarrow s_3$

now reattempt @ e31B

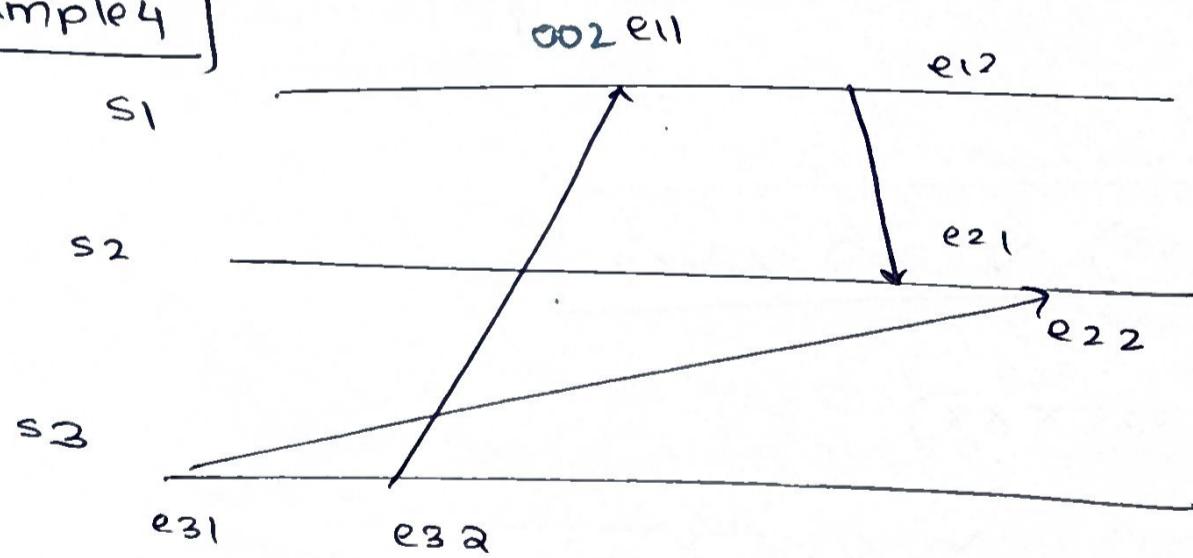
$c_3(110)$ $v_3(010, X, 110)$

→ receives from 4 (from buffer)

$s_1 \rightarrow s_3$

$c_3(010)$, $v_3(001, X, 210)$

Example 4



@ e31 : $c_3(000)$ $v_3(XXX)$

$c_3(001)$ $v_3(X, 001, X)$

$s_3 \rightarrow s_2$

①

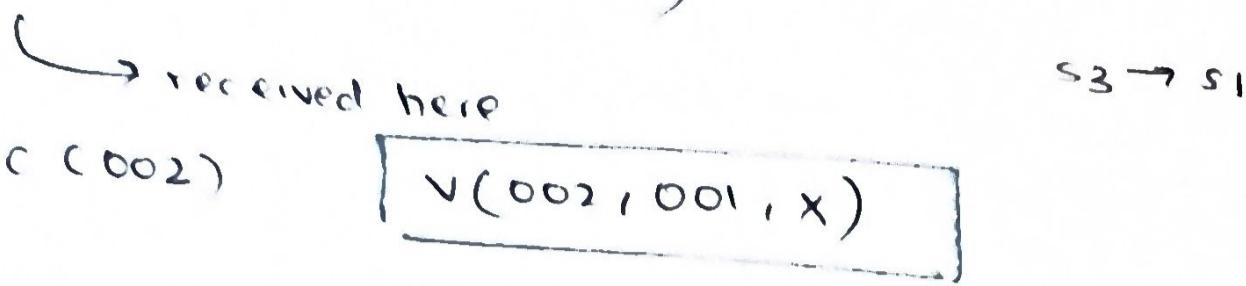
@ e32

$c_3(001)$ $v_3(X, 001, X)$

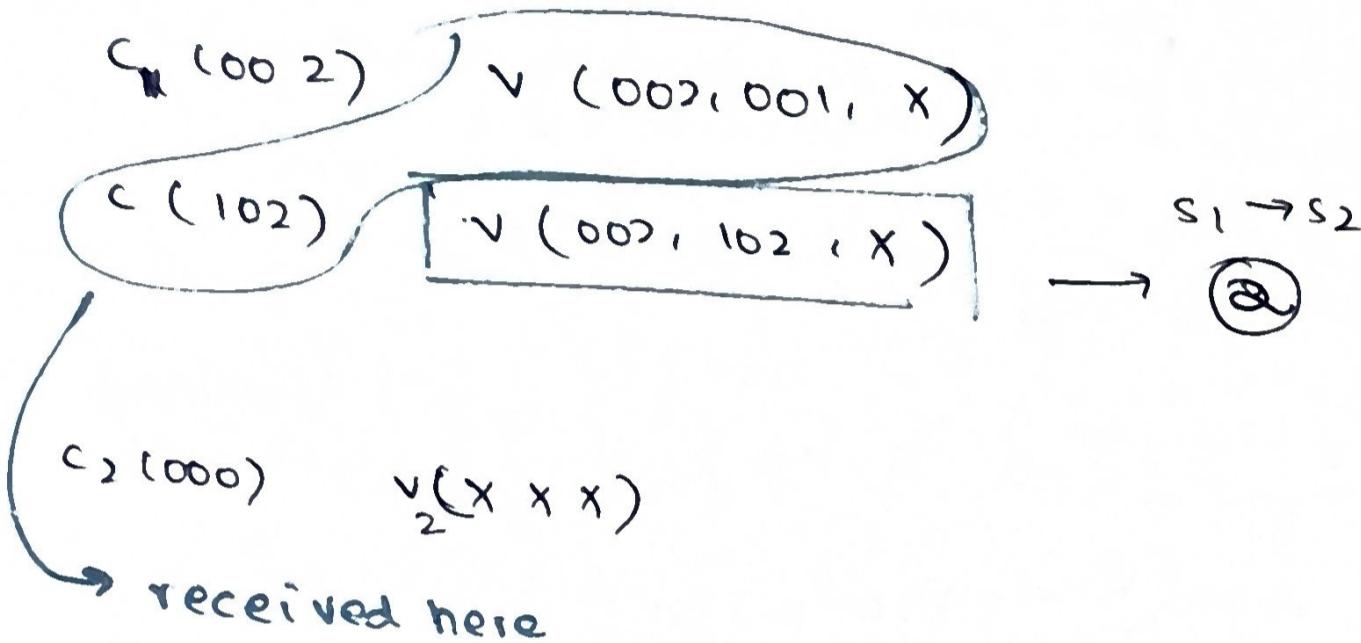
$c_3(002)$ $v_3(002, 001, X)$

$s_3 \rightarrow s_1$

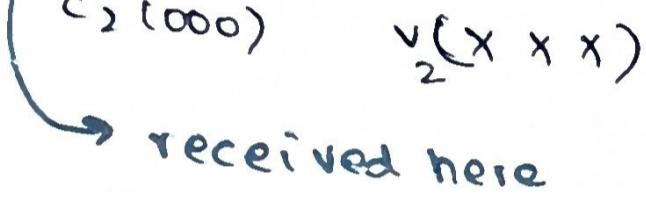
② e₁₁ c(000) v(xxx)



② e₁₂



② e₂₁



Then receiver says nothing is received by S₂ so far but

$v(002, 001, x) \Rightarrow S_3$ has sent a message to S₂ that has not reached.

\Rightarrow buffer $c(102), v(002, 001, x)$

② e₂₂ c₂(000) v₂(xxx)

from ① c₂(001) v₂(x, 001, x) $s_3 \rightarrow s_2$

now reattempt for e_{21B}

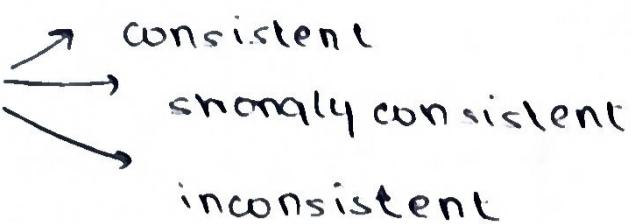
c₂(001) v₂(x, 001, x)

from ② in buffer

c₂(102) v₂(002, 102, x)

$s_3 \rightarrow s_2$

* Global States



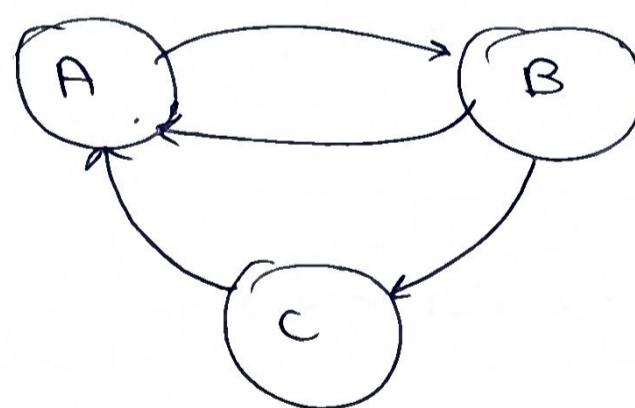
41

Consistent State - even some sends are recorded, but are not received - they are still in transit - will be recorded in the next global state

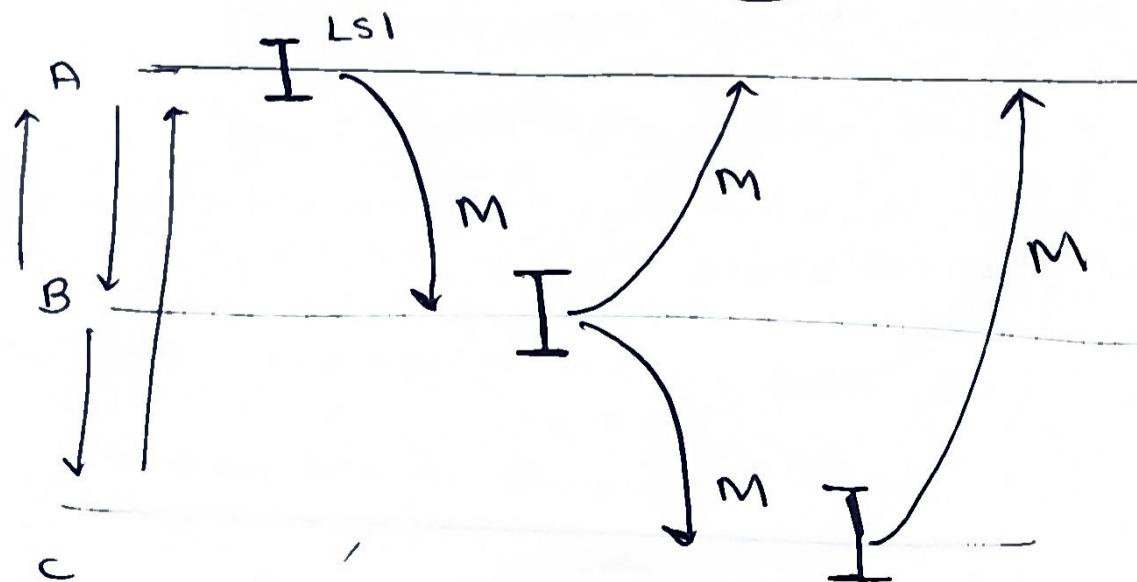
Strongly consistent - every send has a corresponding rec

* Global State Recording Algorithm (GISRA) - chandy lamport algorithm

Consider the following topology



assume A starts the algorithm



1. A records its local state

2. Immediately send a marker to its connections (here: B)

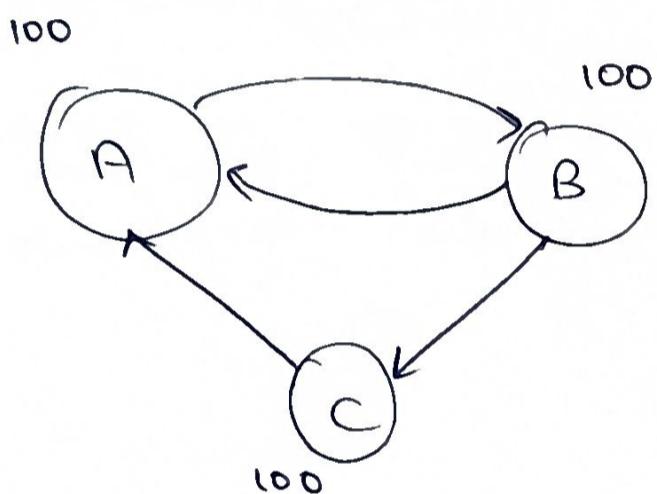
3. B takes a snapshot

4. Immediately send a marker to its connections

5. Now C takes a snapshot and sends a marker to its connection A.

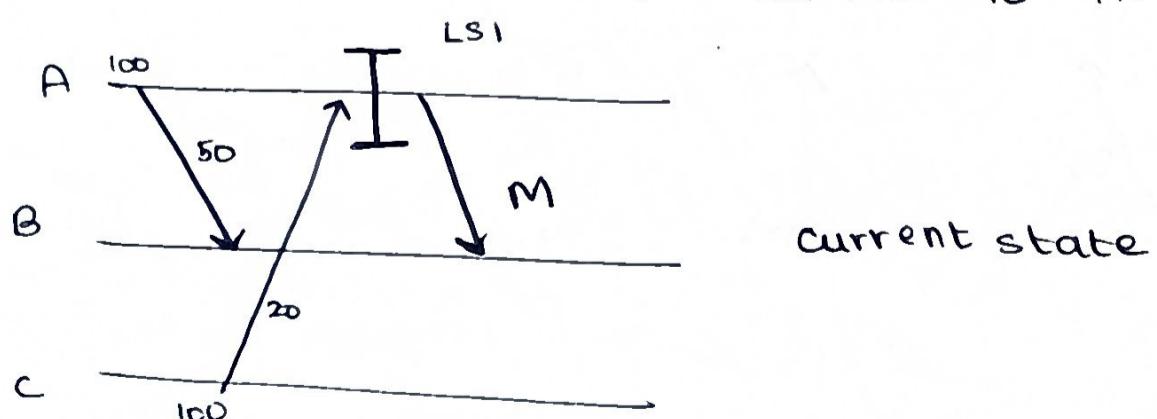
6. The algorithm ends when the initiator has received a marker from all the other connected states

Now, execution of the algorithm with values associated with it.



Start:

1. A sends 50 to B
2. C sends 20 to A
3. A takes a snapshot $\Rightarrow \text{LSI} = 50$
4. Now it has to send markers to its connections (B)



5. Now, before B receives M, it sends 30 to A and 10 to

43

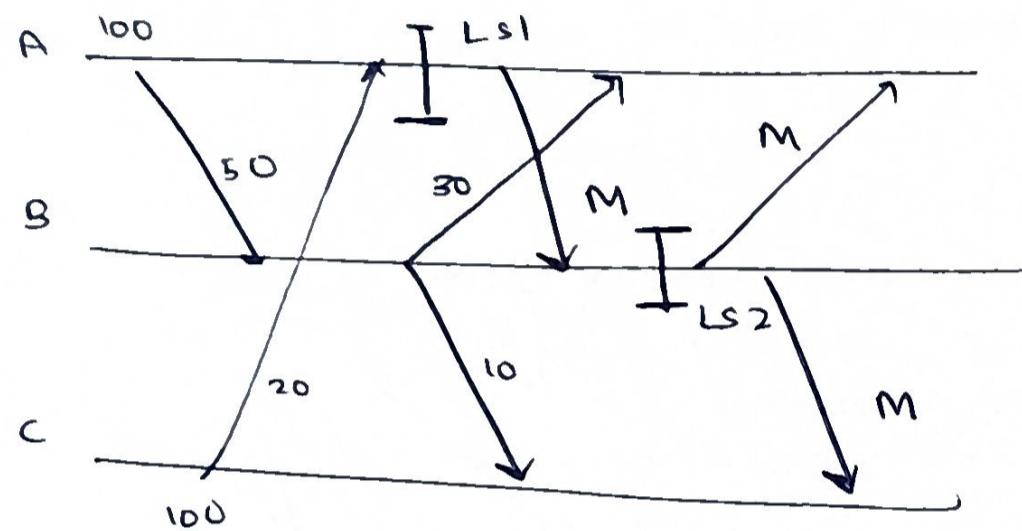
C.

(a)

6. Now it has received a marker

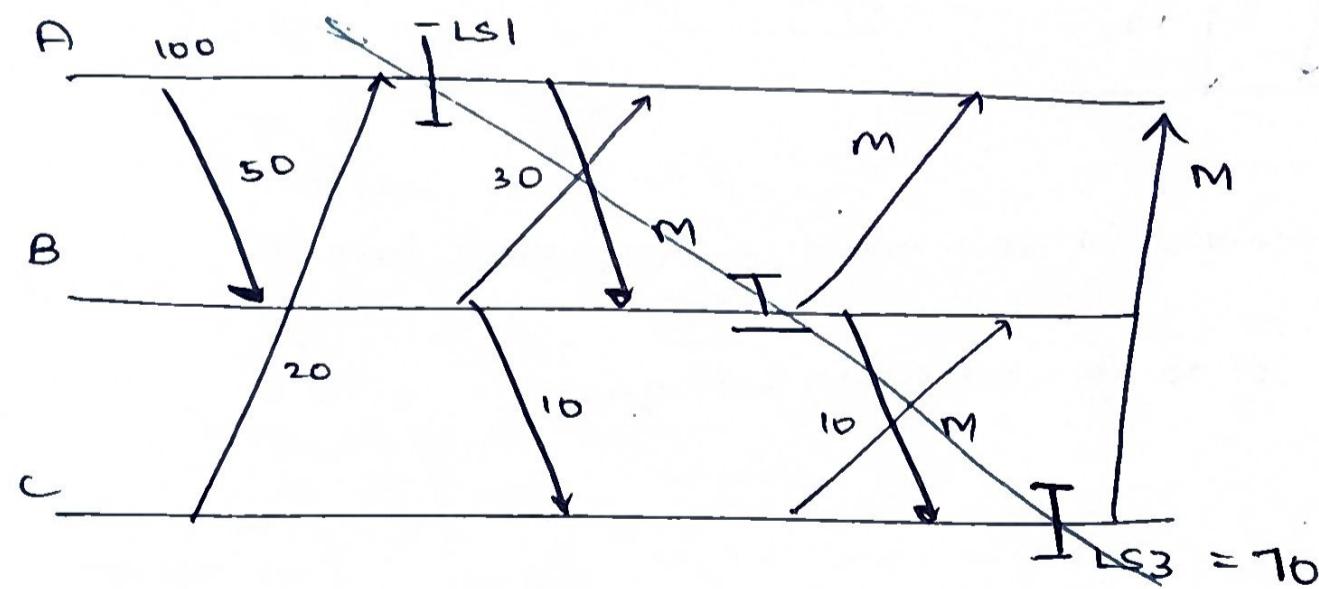
7. Takes a snapshot $LSA : B = 60$

8. Now B sends markers to its connections B & A



9. Before C receives its marker it sends 10 to B.

10. Then it takes a snapshot and sends a marker to its connections(A)

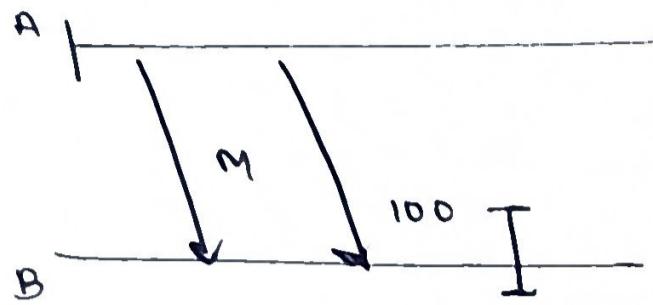


11. A has received markers from all its connected processes.

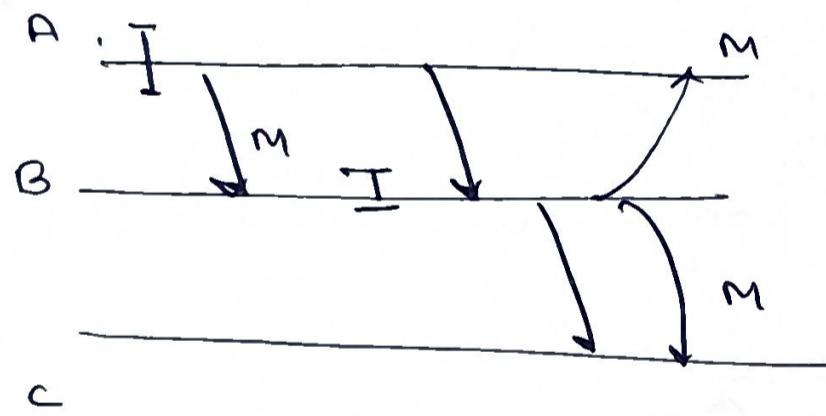
The algorithm is complete

→ Examples for Storing consistent,

* Avoiding Inconsistency in GSRA Algorithm



- B receives a marker, but doesn't take a snapshot (i.e it postpones) until it receives from A.
- Hence, in this case, the global state shows that B has received, but the send is not recorded
- This is inconsistent



- In this case B receives M and takes a snapshot, but it postpones the send of M to its outgoing channels

Basically

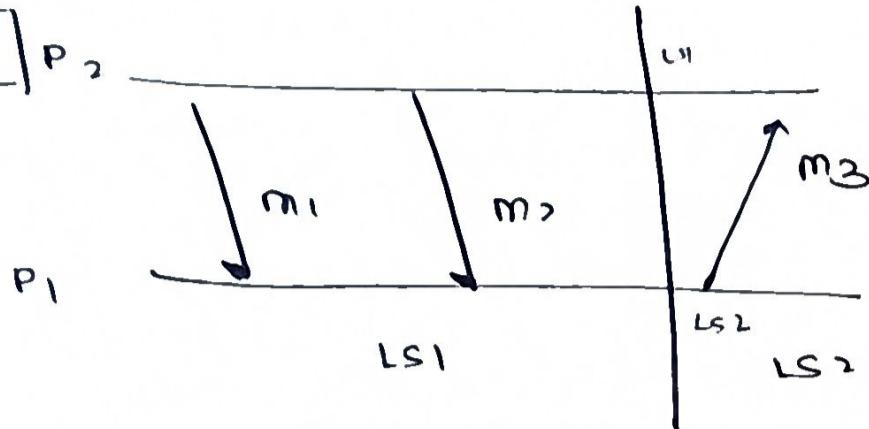
- (i) Receive M
- (ii) Take snapshot
- (iii) send M

⇒ There should be no other sends / receives after an M is received and until all Ms are sent out
a snapshot is taken

45

→ Examples for strongly consistent, consistent and inconsistent

Example 1



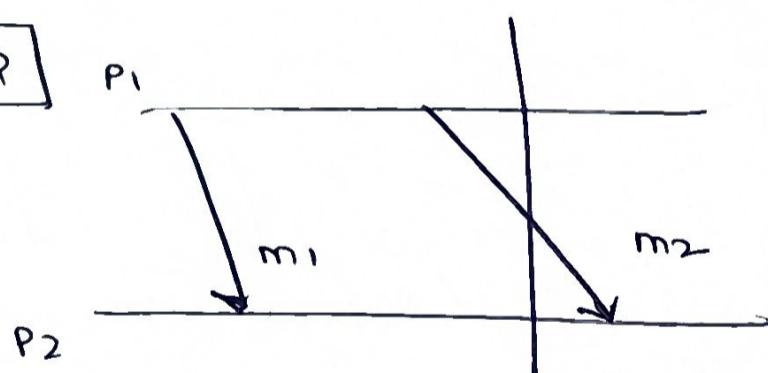
$$GS = LS_1 \cup LS_2$$

$$LS_1 = \{ \text{send}(m_1), \text{send}(m_2) \}$$

$$LS_2 = \{ \text{rec}(m_1), \text{rec}(m_2) \}$$

⇒ strongly consistent

Example 2

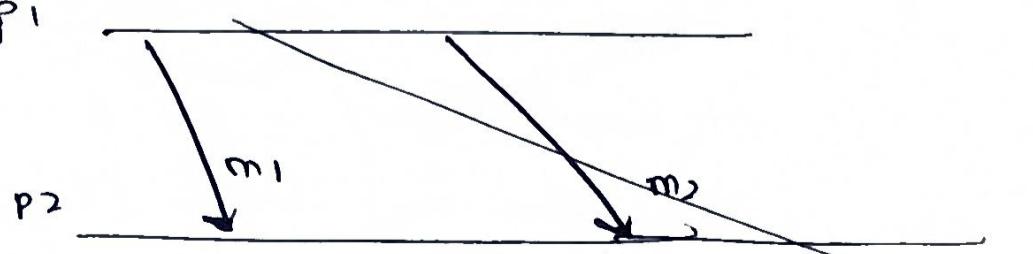


$$LS_1 = \{ \text{send}(m_1), \text{send}(m_2) \}$$

$$LS_2 = \{ \text{rec}(m_1) \}$$

⇒ consistent

Example 3



$$LS_1 = \{ \text{send}(m_1) \}$$

$$LS_2 = \{ \text{rec}(m_1), \text{rec}(m_2) \}$$

⇒ Inconsistent