

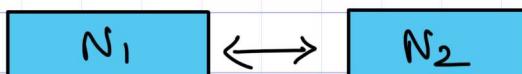
## Synchronization (Absence of global clock)

$N_3 \rightarrow N_1, N_2 : 01$

$N_1 \rightarrow N_3 : 01$        $N_2 \rightarrow N_3 : 46$

$N_2 \rightarrow N_1, N_3 : 47$

$N_1 \rightarrow N_2 : 03$        $N_3 \rightarrow N_2 : 04$



13:25:00

13:25:45

N3

13:26:01

Ordering of msgs w/ clock order (of  $N_i$ , say)  
leads to wrong ordering!

- Clock syncing isn't practical.
- Physical clocks aren't possible.
- Geographically distributed nodes  $\Rightarrow$  different timezones
- Cannot have a common clock.

### Assumptions:

$\rightarrow$  Comm. channel is error free.

$\rightarrow$  Comm. channel is FIFO.

Clock syncing with message passing: based on  
[rcvd. timestamp + RTT/2] assumptions

In case of delay & not FIFO:

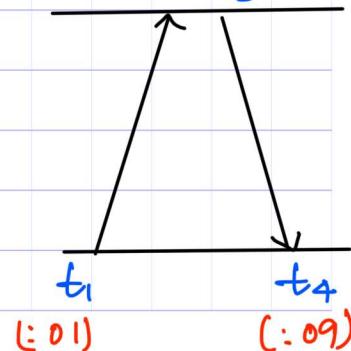
$$(t_4 - t_1) - (t_3 - t_2) = 8 - 4 = 4$$

$\downarrow$   
RTT

$\downarrow$   
processing overhead

true  
 $\downarrow$   
RTT

(:02) (:06)  
 $t_2$   $t_3$



(:01)

(:09)

$$\text{True RTT}_{\frac{1}{2}} = 4_{\frac{1}{2}} = 2$$

$\downarrow$   
almost accurate.

Delay is unpredictable.

$\therefore$  Cannot sync 2 physical clocks.

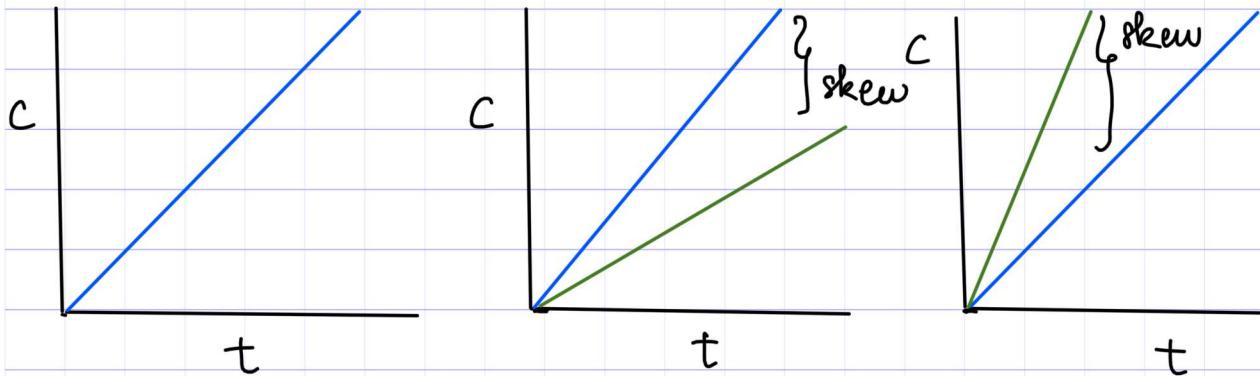
### Clock Synchronization

Clock drift : clocks tick @ diff. rates.

quartz clock :  $10^{-6}$  seconds

diff. of 1 sec / 11.6 days

Clock skew : diff. between 2 clocks @ one point in time



$$\frac{dc}{dt} = 1$$

$$\frac{dc}{dt} < 1$$

$$\frac{dc}{dt} > 1$$

make clock  
run faster  
till sync

make clock  
run slower  
till sync

## Mechanisms to sync:

Cristian's Algorithm

$$T_{\text{new}} = T_{\text{server}} + \text{RTT}/2$$

Berkeley Algorithm

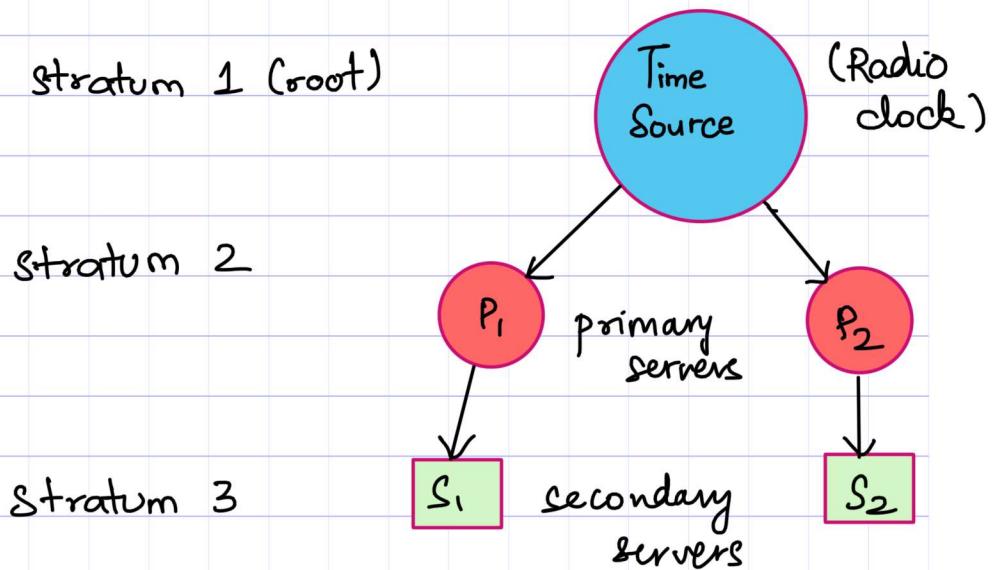
(Fault tolerant org.

Network Time Protocol

timestamping)

## Network Time Protocol :

Clients across Internet sync accurately  
to UTC despite message delays



## Modes :

multicast mode - sends time to all other nodes .

procedure call mode - similar to Cristian's algorithm

symmetric mode - for master servers pair of servers exchange messages .

UDP protocol.

## Logical clocks (causality: as per chronology)

By Leslie Lamport

"happened before" relationship  $\rightarrow$  causally affects.

- If 'a' happened before 'b', then

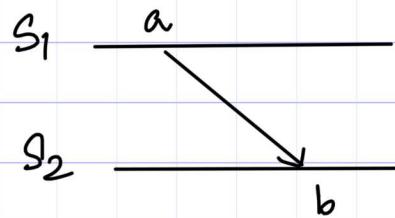
$$\text{Clock}_a < \text{Clock}_b.$$

$$a \rightarrow b$$

- If two processes interact with 'a' in one to 'b' in another, then

$$\begin{aligned} \text{send}(m) &\rightarrow \text{recv}(m) \\ a &\rightarrow b. \end{aligned}$$

$$C_a < C_b$$



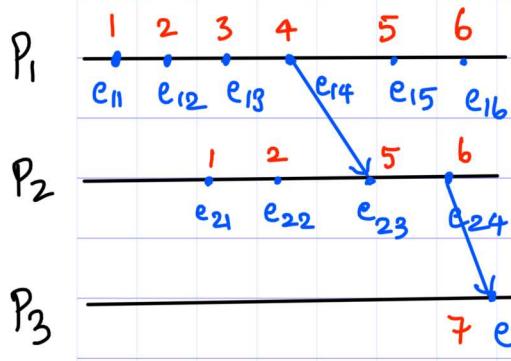
- Rules for incrementing:

$$C_i = C_i + d \quad \text{for processing events} \quad (d = 1)$$

$$C_{j+1} = \max(C_j, t_m) + 1 \quad \text{for communication events.}$$

$\downarrow$   
(message timestamp)

## Partial Orders



• Transitive  $e_{11} \rightarrow e_{22}$

• Antisymmetric  $e_{11} \not\rightarrow e_{21}$

• Irreflexive  $e_{11} \not\rightarrow e_{11}$

Cons : If  $a \rightarrow b$ ;  $C_a \leq C_b$ .

but,

(not strongly consistent) If  $C_a \leq C_b$ ;  $a \rightarrow b$ ? FALSE.

i.e. cannot determine chronological order given any two clock values.

## Vector Clocks

- Rules for incrementing

$$C_i[i] = C_i[i] + 1 \text{ for processing events.}$$

$$\forall_k C_j[k] = \max(E_m[k], C_j[k] + 1)$$

↓

Increment self clock  
on receive.

For some clock values:  $t_a, t_b$

$$t_a == t_b \Rightarrow \forall_k t_a[k] == t_b[k]$$

$$t_a != t_b \Rightarrow \exists_k t_a[k] != t_b[k].$$

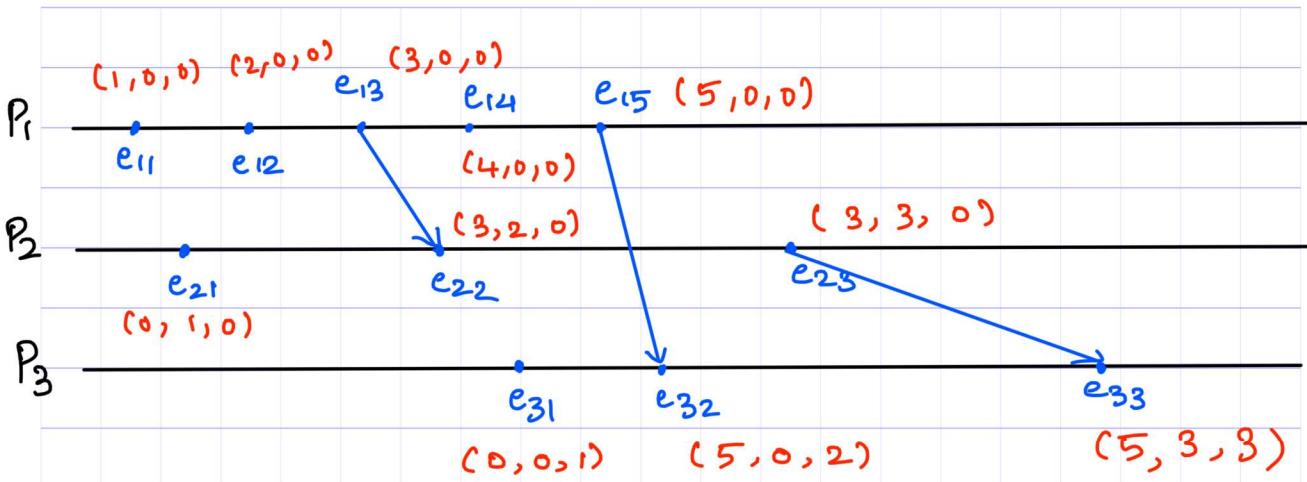
### Cons :

Overlooks concurrent events.

Cannot determine which event occurred  
in which process, in what order.

Need more memory.

No. of processes may not be known in advance.



Applications:

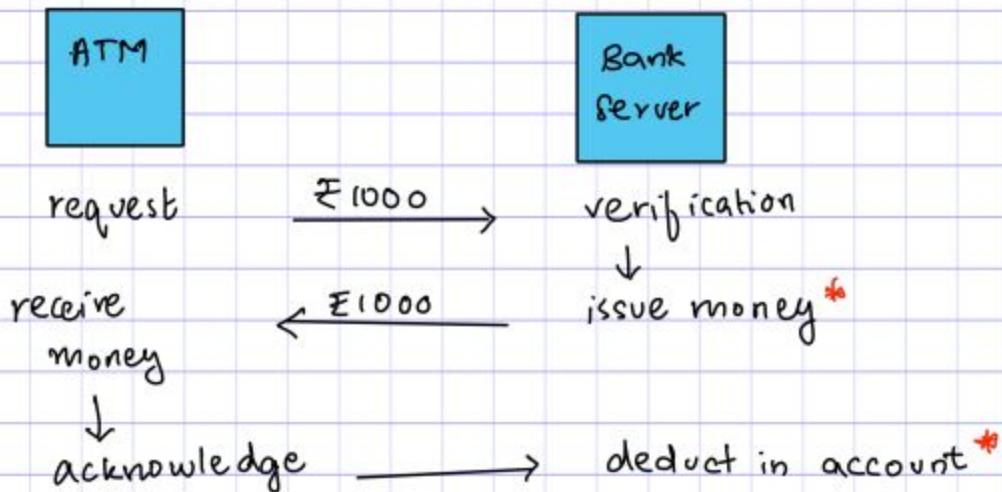
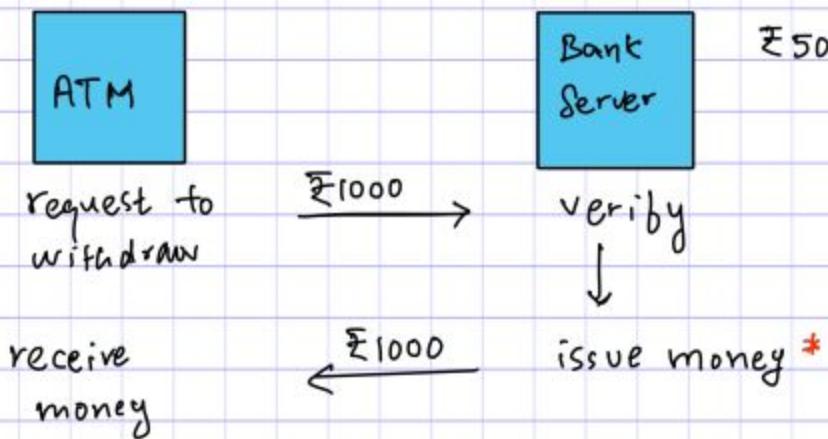
- distributed debugging
- causal distributed shared memory
- global breakpoints
- consistency of checkpoints in optimistic recovery.

Vector clocks are strongly consistent.

i.e. if  $\underbrace{t_a < t_b}_{\downarrow} \Rightarrow a \rightarrow b$ .

solves the Lamport clock problem.

## Absence of shared memory problem



If error happens in bank side, (in both cases)  
the account balance becomes inconsistent  
& bank loses money.

## Causal Ordering of Messages

if  $\text{send}(m_1) \rightarrow \text{send}(m_2)$   
 then  $\text{receive}(m_1) \rightarrow \text{receive}(m_2)$ .

### Assumptions:

Group communication (broadcast)

No loss in communication channel

Customized vector clocks

→ Communication events only

→ On receive, no increment

### Algorithm:

Sender's rule :

Increment  $i^{\text{th}}$  value of send timestamp.

Receiver's rule : ( $i$ -sender,  $j$ -receiver)

$\rightarrow c_j[i] = tm[i] - 1$

$\rightarrow c_j[k] \geq tm[k], k \neq i$

( $>$  because send events can be concurrent)

### Case (i) (Basic Case)

$e_{21} : (0, 1, 0) \rightarrow \text{timestamp}$

$e_{11} : (i=2, j=1) \in [0, 0, 0]$

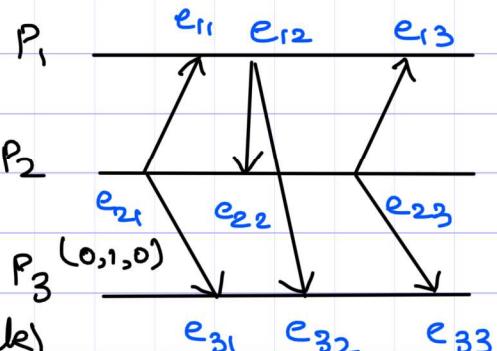
(current clock)

(i)  $c_1[2] == tm[2] - 1$  ? true

(ii)  $c_1[1, 3] \geq tm[1, 3]$  ? true.

∴ Accept message (∴ in causal order)

$P_1$  clock :  $(0, 1, 0)$



$e_{31} : (i=2, j=3) \quad C_3 [0,0,0] \xrightarrow{} tm [0,1,0]$   
 (i)  $C_3 [2] = tm [2] - 1$  ? true  
 (ii)  $C_3 [1,3] \geq tm [1,3]$  ? true  
 $\therefore$  Accept message  
 $P_3$  clock :  $(0,1,0)$

$e_{12}$  : Sender's rule.  $P_1$  clock :  $(1,1,0)$

$e_{22} : C_2 [0,1,0] \xrightarrow{} tm [1,1,0]$   
 $(i=1, j=2)$   
 (i)  $C_2 [1] = tm [1] - 1$  ? true  
 (ii)  $C_2 [2,3] \geq tm [2,3]$  ? true  
 $\therefore$  Accept message  
 $P_2$  clock :  $(1,1,0)$

$e_{32} : C_3 [0,1,0] \xrightarrow{} tm [1,1,0]$   
 $(i=1, j=3)$   
 (i)  $C_3 [1] = tm [1] - 1$  ? true  
 (ii)  $C_3 [2,3] \geq tm [2,3]$  ? true  
 $\therefore$  Accept message  
 $P_3$  clock :  $(1,1,0)$

$e_{23}$  : Sender's rule.  $P_2$  clock :  $(1,2,0)$

$e_{13} : C_1 [1,1,0] \xrightarrow{} tm [1,2,0]$   
 $(i=2, j=1)$   
 (i)  $C_1 [2] = tm [2] - 1$  ? true  
 (ii)  $C_1 [1,3] \geq tm [1,3]$  ? true  
 $\therefore$  Accept message  
 $P_1$  clock :  $(1,2,0)$

$e_{33} : C_3 [1, 1, 0] \text{ tm } [1, 2, 0]$   
 $(i = 2, j = 3)$

(i)  $C_3 [2] == \text{tm}[2] - 1 ? \text{true}$

(ii)  $C_3 [1, 3] \geq \text{tm}[1, 3] ? \text{true}$

$\therefore \text{Accept message}$

$P_3 \text{ clock} : (1, 2, 0)$

Case (ii)

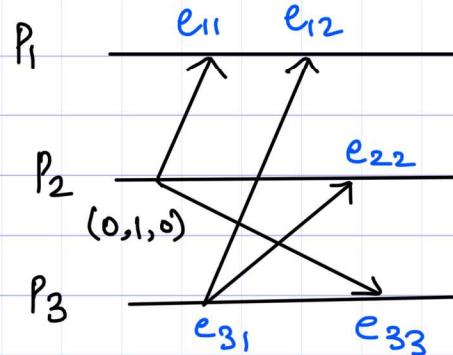
$e_{11} : C_1 (0, 0, 0) \text{ tm } (0, 1, 0)$   
 $(i = 2, j = 1)$

(i)  $C_1 [2] == \text{tm}[2] - 1 ?$   
 $\text{true}$

(ii)  $C_1 [1, 3] \geq \text{tm}[1, 3] ? \text{true}$

$\therefore \text{Accept message}$

$P_1 \text{ clock} : (0, 1, 0)$



$e_{31} : \text{Sender's rule . } P_3 \text{ clock} : (0, 0, 1)$

$e_{33} : C_3 (0, 0, 1) \text{ tm } (0, 1, 0)$   
 $(i = 2, j = 3)$

(i)  $C_3 [2] == \text{tm}[2] - 1 ? \text{true}$

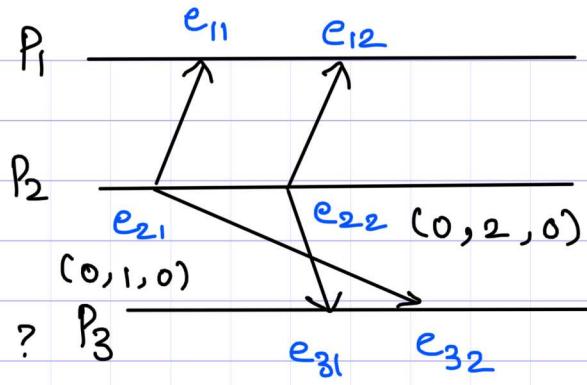
(ii)  $C_3 [1, 3] \geq \text{tm}[1, 3] ? \text{true}$

$\therefore \text{Accept message.}$

$P_3 \text{ clock} : (0, 1, 1)$

Similarly,  $e_{12}$  &  $e_{22}$  are accepted.

Case (iii) (not in causal order)



$e_{11} : C_1(0,0,0) \text{ tm}(0,1,0)$   
 $(i = 2, j = 1)$

(i)  $C_1[2] == \text{tm}[2] - 1$ ? true

(ii)  $C_1[1,3] \geq \text{tm}[1,3]$ ? true

∴ Accept message

$P_1 \text{ clock: } (0,1,0)$

→ 2nd msg delivered first.

$e_{31} : C_3(0,0,0) \text{ tm}(0,2,0)$   
 $(i = 2, j = 3)$

(i)  $C_3[2] == \text{tm}[2] - 1$ ? false

There is a gap!  $P_3$  missed a message from  $P_2$ .

∴ Do not accept the message, buffer it.

$P_3 \text{ buffer: } (0,2,0)_{e_{31}}$

$e_{32} : C_3(0,0,0) \text{ tm}(0,1,0)$   
 $(i = 2, j = 3)$

(i) & (ii) are true.

∴ Accept the message.

$P_3 \text{ clock: } (0,1,0)$

Retrieve message from Buffer & try delivery.

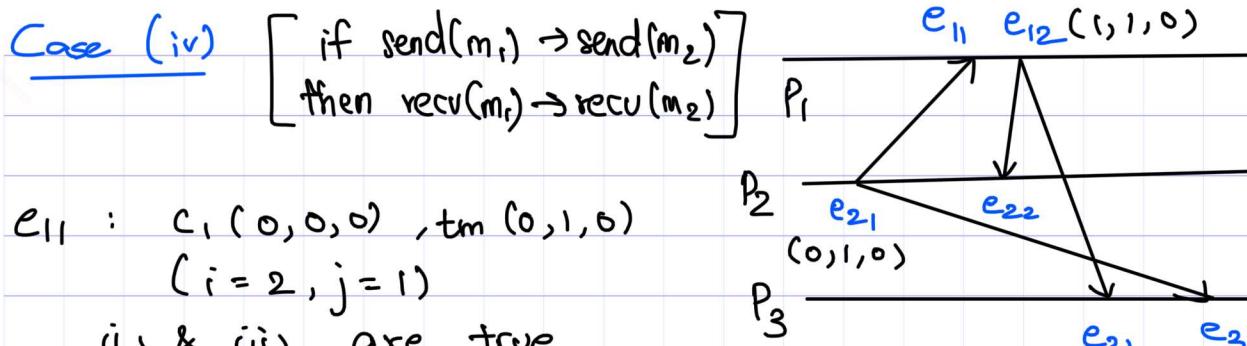
$e_{31}' : C_3(0,1,0) \text{ tm}(0,2,0)$

(i) & (ii) are true now.

∴ Accept message.

$P_3 \text{ clock: } (0,2,0)$ .

$e_{12}$  will be accepted.



$e_{22}$  : will be accepted similarly.  
 $C_2$  clock :  $(1,1,0)$

$e_{31}$  :  $C_3(0,0,0)$   $tm(1,1,0) \downarrow (i = 1, j = 3)$

$(i_1) C_3[1] == tm[1] - 1 ?$  true.

$\boxed{\begin{array}{l} \text{Conveys to } P_3 \text{ that } P_1 \text{ recv'd.} \\ \text{a msg. from } P_2 \text{ which} \\ P_3 \text{ hasn't yet recv'd.} \end{array}}$

$(ii_1) C_3[2,3] \geq tm[2,3] ?$  false.  
 Buffer  $e_{31}$  in  $P_3$ . Do not accept.

$e_{32}$  :  $C_3(0,0,0)$   $tm(0,1,0)$   
 $(i = 2, j = 3)$

$(i_1 \& ii_1)$  are true.  
 $\therefore$  Accept message.  
 $P_3$  clock :  $(0,1,0)$

Attempt redelivery.

$e'_{31}$  :  $C_3(0,1,0)$   $tm(1,1,0)$   
 $(i = 1, j = 3)$

$(i_1 \& ii_1)$  are true.  
 $\therefore$  Accept message.  
 $P_3$  clock :  $(1,1,0)$

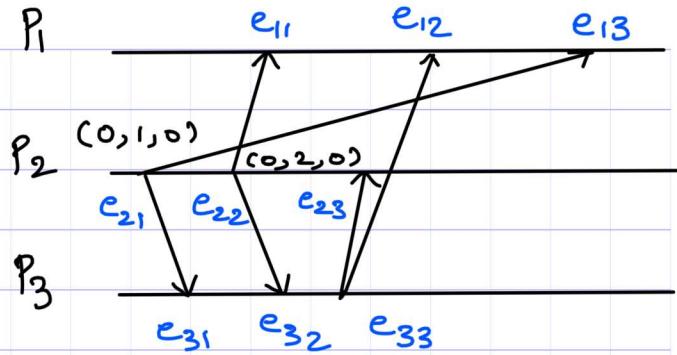
Case (v)

$e_{31} : C_3(0,0,0)$   
 $tm(0,1,0)$

( $i = 2, j = 3$ )  
 (i) & (ii) are true.

∴ Accept message.

$C_3 \text{ clock} : (0,1,0)$



$e_{32} : C_3(0,1,0) \quad tm(0,2,0)$   
 $(i = 2, j = 3)$

(i) & (ii) are true.

∴ Accept message.

$C_3 \text{ clock} : (0,2,0)$

$e_{33} : \text{Sender's rule. } C_3 \text{ clock} : (0,2,1)$

$e_{23} : \text{will be accepted.}$

$C_2 \text{ clock} : (0,2,1)$

$e_{11} : C_1(0,0,0) \quad tm(0,2,0) \quad (i=3, j=1)$

(i) false.

Buffer the message.

$P_1 \text{ buffer} : (0,2,0)_{e_{11}}$

$e_{12} : C_1(0,0,0) \quad tm(0,2,1) \quad (i=3, j=1)$

(i) true

(ii) false

Buffer the message

$P_1 \text{ buffer} : e_{11}, (0,2,1)_{e_{12}}$

$e_{13} : C_1(0, 0, 0) \text{ tm}(0, 1, 0) \quad (i=2, j=1)$

Rule (i) & (ii) are true.

∴ Accept message.

$C_1 \text{ clock} : (0, 1, 0).$

Attempt redelivery.

$e_{11}' : C_1(0, 1, 0) \text{ tm}(0, 2, 0)$

(i) & (ii) are true.

∴ Accept message.

$C_1 \text{ dock} : (0, 2, 0).$

$e_{12}' : C_1(0, 2, 0) \text{ tm}(0, 2, 1)$

(i) & (ii) are true.

∴ Accept message.

$C_1 \text{ clock} : (0, 2, 1)$

Case (vi)

$e_{31} : C_3(0, 0, 0)$   
 $\text{tm}(0, 3, 0)$   
 $(i=2, j=3)$

(i) false

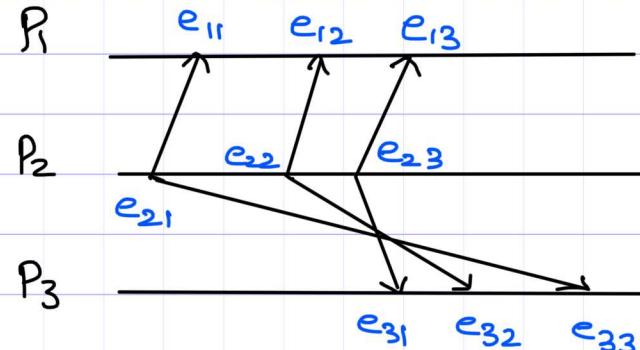
Buffer  $e_{31}$  in  $P_3$ .

$P_3 \text{ Buffer} : (0, 3, 0) e_{31}$

$e_{32} : C_3(0, 0, 0) \text{ tm}(0, 2, 0) \quad (i=2, j=3)$

(i) false.

$P_3 \text{ Buffer} : e_{31}, (0, 2, 0) e_{32}$



$e_{33} : C_3(0, 0, 0) \text{ tm } (0, 1, 0) \quad (i=2, j=3)$   
(i) & (ii) are true.

∴ Accept message.

$C_3 \text{ clock} : (0, 1, 0)$ .

Attempt redelivery.

$e_{31}' : C_3(0, 1, 0) \text{ tm } (0, 3, 0) \quad (i=2, j=3)$   
(i) false

Buffer  $e_{31}'$  again!

$P_3 \text{ Buffer} : e_{32}, (0, 3, 0)_{e_{31}}$ ,

$e_{32}' : C_3(0, 1, 0) \text{ tm } (0, 2, 0) \quad (i=2, j=3)$   
(i) & (ii) are true.

Accept message.

$C_3 \text{ clock} : (0, 2, 0)$

$e_{33}'' : C_3(0, 2, 0) \text{ tm } (0, 3, 0) \quad (i=2, j=3)$   
(i) & (ii) are true.

Accept message.

$C_3 \text{ clock} : (0, 3, 0)$ .

Case (vii)

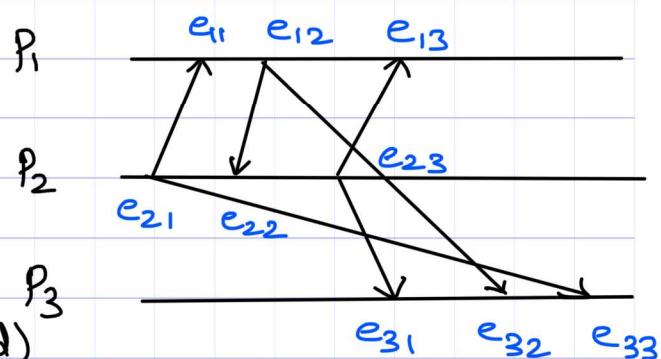
(Re-entrance,  
but not from  
the same  
sender)

$e_{12} : C_2 : (0, 1, 0) \text{ (send)}$

$e_{11} : \text{Accept. } C_1 : (0, 1, 0)$

$e_{12} : C_1 : (1, 1, 0) \text{ (send)}$

$e_{22} : \text{Accept. } C_2 : (1, 1, 0)$



$e_{23} : C_2 : (1, 2, 0)$  (send)

$e_{13} : \text{Accept. } C_1 : (1, 2, 0)$

$e_{31} : C_3 (0, 0, 0) \text{ tm}(1, 2, 0) (i=2, j=3)$   
(i) false.

$P_3 \text{ Buffer} : (1, 2, 0) e_{31}$

$e_{32} : C_3 (0, 0, 0) \text{ tm}(1, 1, 0) (i=1, j=3)$

(i) true

(ii) false

$P_3 \text{ Buffer} : e_{31}, (1, 1, 0) e_{32}$

$e_{33} : C_3 (0, 0, 0) \text{ tm}(0, 1, 0) (i=2, j=3)$

(i) & (ii) are true.

$\therefore \text{Accept. } P_3 \text{ clock} : (0, 1, 0)$

Attempt redelivery.

$e_{31}' : C_3 (0, 1, 0) \text{ tm}(1, 2, 0) (i=2, j=3)$

(i)  $C_3[2] == \text{tm}[2]-1$  ? true.

(ii)  $C_3[1, 3] \geq \text{tm}[1, 3]$  ? false.

$P_3 \text{ Buffer} : e_{32}, (1, 2, 0) e_{31}'$

$e_{32}' : C_3 (0, 1, 0) \text{ tm}(1, 1, 0) (i=1, j=3)$

(i)  $C_3[1] == \text{tm}[1]-1$  ? true.

(ii)  $C_3[2, 3] \geq \text{tm}[2, 3]$  ? true.

$\therefore \text{Accept message.}$

$C_3 \text{ clock} : (1, 1, 0).$

$e_{33}'' : C_3 (1, 1, 0) \text{ tm}(1, 2, 0) (i=2, j=3)$

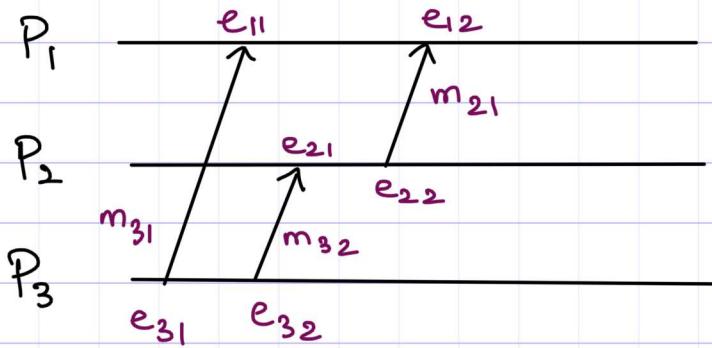
(i) & (ii) are true.

$\therefore \text{Accept the message.}$

$C_3 \text{ clock} : (1, 2, 0).$

If 2 send events causally affect each other, then the corresponding 2 receive events must also causally affect each other.

## Causal Ordering of Messages - Non Broadcast



$\text{send}(m_{31}) \rightarrow \text{send}(m_{32})$

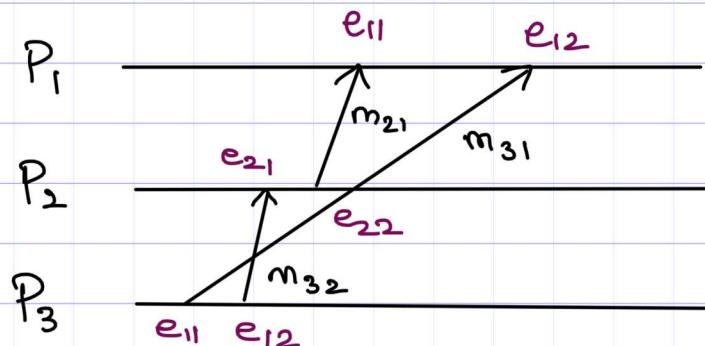
$\text{send}(m_{32}) \rightarrow \text{rec}(m_{32})$

$\text{rec}(m_{32}) \rightarrow \text{send}(m_{21})$

$\text{send}(m_{21}) \rightarrow \text{rec}(m_{21})$

Conclusion  $\Rightarrow \text{rec}(m_{31}) \rightarrow \text{rec}(m_{21})$

Delivery is causal.



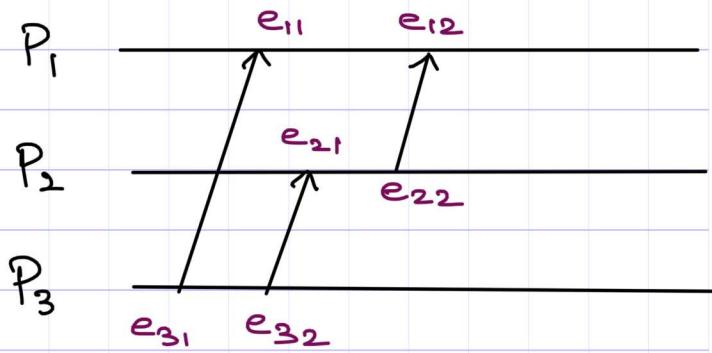
Non-causal delivery.

$m_{21}$  should be delivered after  $m_{31}$ .

## Schipper - Eggli - Sandoz Protocol

e<sub>31</sub>  
 $C(0, 0, 1)$   
 $V_3(x, x, x)$

Send  $t, V_3$ .



Update  $V_3$ .

$V_3((0, 0, 1), x, x)$

$t$  - timestamp  
 $V_m$  - message vector

e<sub>32</sub>

$C(0, 0, 2)$      $V_3((0, 0, 1), x, x)$

Send  $t, V_3$ .

Update  $V_3$ .

$V_3((0, 0, 1), (0, 0, 2), x)$

e<sub>11</sub>

$C(0, 0, 0)$      $V_1(x, x, x)$   
 $t(0, 0, 1)$      $V_m(x, x, x)$

→ Both null.

→ nothing has happened yet.

Accept the message.

$V_1(x, x, (0, 0, 1))$

copied here  
because message received  
from  $P_3$ .

Merge the vectors  
& copy the  
timestamp to the  
vector.

$C(0, 0, 1)$

e<sub>21</sub>

$$C(0, 0, 0) \quad V_2(x, \textcircled{x}, x)$$
$$t(0, 0, 2) \quad V_m((0, 0, 1), \textcircled{x}, x)$$

no conflict

$$V_2((0, 0, 1), x, (0, 0, 2)) \quad C(0, 0, 2)$$

e<sub>22</sub>

$$C(0, 1, 2) \quad V_2((0, 0, 1), x, (0, 0, 2))$$

Send  $C$ ,  $V_2$ .  
Update  $V_2$ .

$$V_2((0, 1, 2), x, (0, 0, 2))$$

e<sub>12</sub>

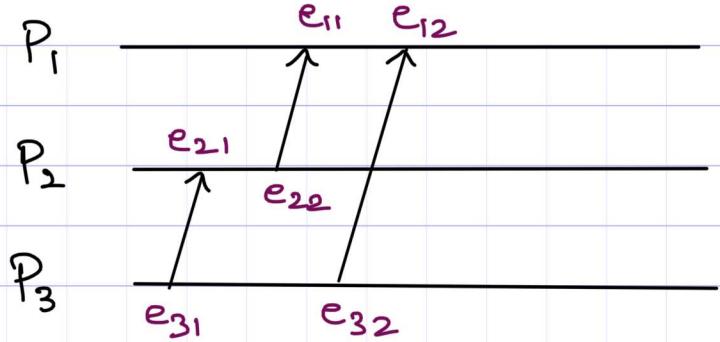
$$C(0, 0, 1) \quad V_1(x, x, (0, 0, 1))$$
$$t(0, 1, 2) \quad V_m((0, 0, 1), x, (0, 0, 2))$$

$e_{31}$

Send event

$$C_3(0, 0, 1)$$

$$V_3(x, x, x)$$



Send  $C_3, V_3$ .

Update  $V_3$ .

$$V_3(x, (0, 0, 1), x).$$

$e_{32}$

Send event

$$C_3(0, 0, 2)$$

$$V_3(x, (0, 0, 1), x)$$

Send  $C_3, V_3$ .

Update  $V_3$ .

$$V_3((0, 0, 2), (0, 0, 1), x).$$

$e_{21}$

$$C_2(0, 0, 0)$$

$$V_2(0, 0, 0)$$

Receive

$$C_m(0, 0, 1) \in V_m(x, x, x). \text{ from sender.} \rightarrow \text{First msg.}$$

Accept the message.

Update  $V_2$ .

$$V_2(x, x, (0, 0, 1))$$

Merge Clock Values.

$$C_2(0, 0, 1)$$

e<sub>22</sub>

Send event

$C_2(0, 1, 1)$

$V_2(x, x, (0, 0, 1))$

Send  $C_2 \in V_2$ .

Update  $V_2$ .

$V_2((0, 1, 1), x, (0, 0, 1))$

e<sub>11</sub>

$C_1(0, 0, 0)$        $V_1(x, x, x)$

→ no conflict.

Receive

$C_m(0, 1, 1)$        $V_m(x, x, (0, 0, 1))$

Accept the message.

Update  $C_1$ .       $C_1(0, 1, 1)$

Update  $V_1$ .       $V_1(x, \underbrace{(0, 1, 1)}, \underbrace{(0, 0, 1)})$

↓  
timestamp copied to  
sender's location.

+ merged  $V_m \in V_1$

Now  $P_1$  is aware  
that  $P_2$  has sent a neg. to  $P_2$ .

e<sub>12</sub>

C<sub>1</sub> (0, 1, 1)

V<sub>1</sub> (x, (0, 1, 1), (0, 0, 1))

Receive

C<sub>m</sub> (0, 0, 2)

V<sub>m</sub> (x, (0, 0, 1), x)

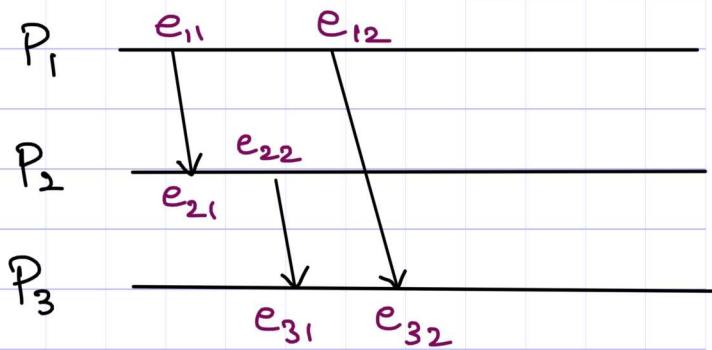
no conflict

Accept the message.

Update/merge C<sub>1</sub> & V<sub>1</sub>.

V<sub>1</sub> (x, (0, 1, 1), (0, 0, 2))

C<sub>1</sub> (0, 1, 2)



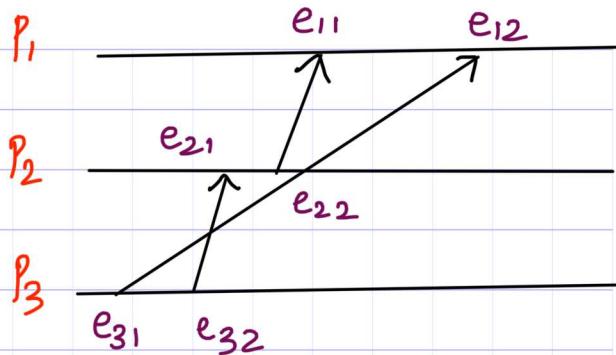
Tony.

## Example

e<sub>31</sub>

C<sub>3</sub> [ 0 0 0 ]

V<sub>3</sub> [ x x x ]



e<sub>32</sub>

C<sub>3</sub> [ 0 0 1 ]

V<sub>3</sub> [ (0 0 1) x x ]

send C<sub>3</sub> [ 0 0 1 ] ; V<sub>3</sub> [ (0 0 1) x x ] to e<sub>21</sub>.

update C<sub>3</sub> [ 0 0 2 ] ; V<sub>3</sub> [ (0 0 1) (0 0 2) x ]

e<sub>21</sub>

C<sub>2</sub> [ 0 0 0 ] V<sub>2</sub> [ x x x ]

recv C<sub>0</sub> [ 0 0 2 ] ; V<sub>m</sub> [ (0 0 1) x x ]

update C<sub>2</sub> [ 0 0 2 ] ; V<sub>2</sub> [ (0 0 1) (0 0 2) x ]

e<sub>22</sub>

C<sub>2</sub> [ 0 0 2 ] ; V<sub>2</sub> [ (0 0 1) (0 0 2) x ]

send C<sub>2</sub> [ 0 1 2 ] ; V<sub>2</sub> [ (0 0 1) (0 0 2) x ]

update C<sub>2</sub> [ 0 1 2 ] ; V<sub>2</sub> [ (0 1 2) (0 0 2) x ]

e<sub>11</sub>

C<sub>1</sub> [ 0 0 0 ] V<sub>1</sub> [ x x x ]

recv C<sub>1</sub> [ 0 1 2 ] V<sub>m</sub> [ (0 0 1) (0 0 2) x ]

send C<sub>1</sub> [ 0 1 2 ] V<sub>m</sub> [ (0 0 1) (0 0 2) x ]

to buffer.

e<sub>12</sub>

C<sub>1</sub> [0 0 0 0]

V<sub>1</sub> [x x x]

recv. [0 0 1]

V<sub>m</sub> [x x x]

update C<sub>1</sub> [0 0 1] V<sub>1</sub> [(0 0 1) x x]

since clock C<sub>1</sub> is updated, check buffer.

e<sub>11'</sub>

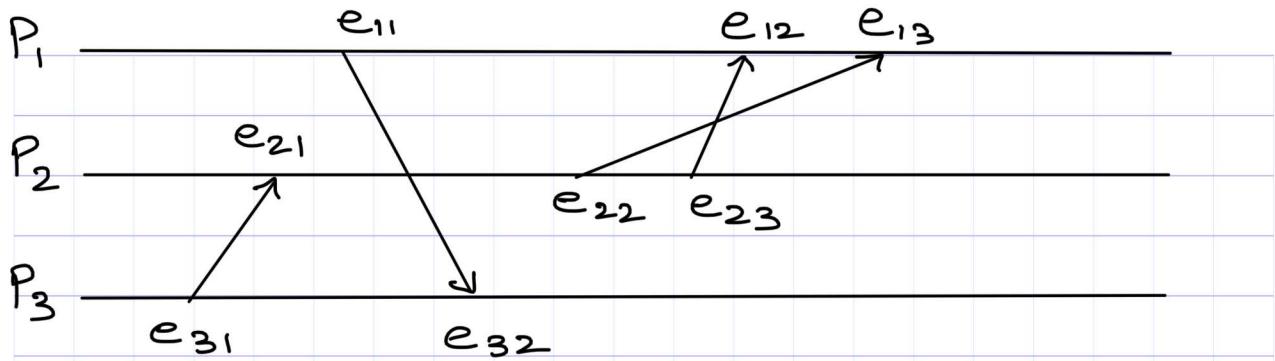
C<sub>1</sub> [0 0 1]

V<sub>1</sub> [(0 0 1) x x]

recv. [0 1 2]

V<sub>m</sub> [(0 0 1) (0 0 2) x]

update C<sub>1</sub> [0 1 2] V<sub>1</sub> [(0 0 1) (0 0 2) x]



e<sub>31</sub> (send)

$$C_3 [0, 0, 0]$$

$$V_3 [x, x, x]$$

$$C_3 [0, 0, 1]$$

$$V_3 [x, x, x]$$

Send to e<sub>21</sub>

$$V_3 [x, [0, 0, 1], x]$$

e<sub>21</sub> (receive)

$$C_2 [0, 0, 0]$$

$$V_2 [x, \cancel{x}, x]$$

$$\rightarrow C_m [0, 0, 1]$$

$$V_m [x, \cancel{x}, x]$$

unset

Accept delivery.

$$C_2 [0, 0, 1] \quad C_2 [x, [0, 0, 1], x]$$

e<sub>11</sub> (send)

$$C_1 [0, 0, 0] \quad V_1 [x, x, x]$$

$$C_1 [1, 0, 0] \quad V_1 [x, x, x]$$

Send to e<sub>32</sub>

$$V_1 [x, x, [1, 0, 0]]$$

$e_{32}$  (receive)

$$\begin{array}{ll} C_3 [0, 0, 1] & V_3 [x, [0, 0, 1], \textcircled{x}] \\ C_m [1, 0, 0] & V_m [x, x, \textcircled{x}] \\ & \text{unset} \end{array}$$

Accept delivery

$$\begin{array}{l} C_3 [1, 0, 1] \\ V_3 [x, [0, 0, 1], [1, 0, 0]] \end{array}$$

 $e_{22}$  (send)

$$\begin{array}{ll} C_2 [0, 0, 1] & V_2 [x, [0, 0, 1], x] \\ C_2 [0, 1, 1] & V_2 [x, [0, 0, 1], x] \\ & \text{send to } e_{13} \end{array}$$

$$V_2 [[0, 1, 1], [0, 0, 1], x]$$

 $e_{23}$  (send)

$$\begin{array}{ll} C_2 [0, 2, 1] & V_2 [[0, 1, 1], [0, 0, 1], x] \\ & \text{send to } e_{12} \end{array}$$

$$V_2 [[0, 2, 1], [0, 0, 1], x]$$

 $e_{12}$  (receive)

$$\begin{array}{ll} C_1 [1, 0, 0] & V_1 [\textcircled{x}, x, [1, 0, 0]] \\ C_m [0, 2, 1] & V_m [[0, 1, 1], [0, 0, 1], x] \end{array}$$

buffer it.

 $e_{13}$  (receive)

$$\begin{array}{ll} C_1 [1, 0, 0] & V_1 [\textcircled{x}, x, [1, 0, 0]] \\ C_m [0, 1, 1] & V_m [\textcircled{x}, [0, 0, 1], x] \\ & \text{accept delivery} \end{array}$$

$$C_1 [1, 1, 1] \quad V_1 [[0, 1, 1], x, [1, 0, 0]]$$

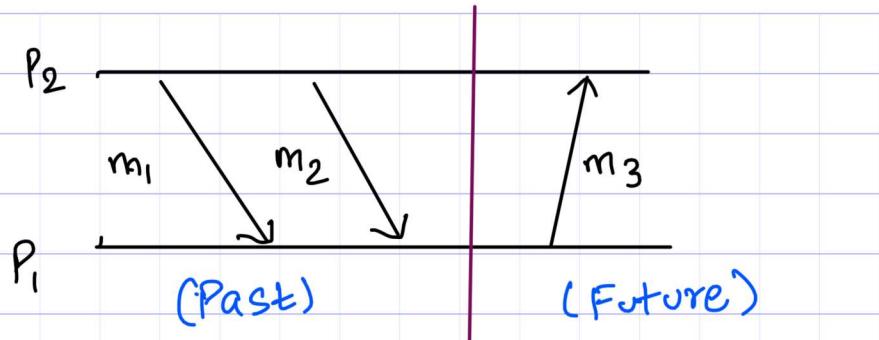
e<sub>12</sub>' (redelivery)

C, [1, 1, 1] V, [ [0, 1, 1], X , [1, 0, 0] ]  
Cm [0, 2, 1] Vm [ [0, 1, 1], X , , X ]

accept delivery

C, [1, 2, 1] V, [ [0, 2, 1], X , [1, 0, 0] ]

## Global States



Global state =  $\bigcup_i^n$  Local states;

Here,

$$GS = LS_1 \cup LS_2.$$

[Cut is a graphical line splitting the time-space graph into 2 parts]

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

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

$$GS = \left\{ \begin{array}{l} \text{send}(m_1), \text{send}(m_2), \\ \text{receive}(m_1), \text{receive}(m_2) \end{array} \right\}$$

strongly  
consistent

## Strongly consistent global state

1) If send event in GS, corresponding receive event will be present

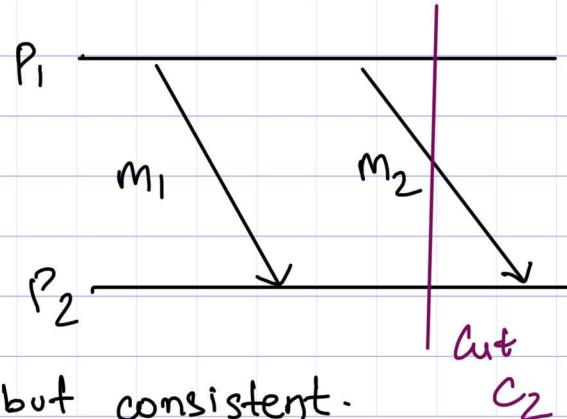
2) & no(send) = no(receive)

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

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

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

not strongly consistent, but consistent.



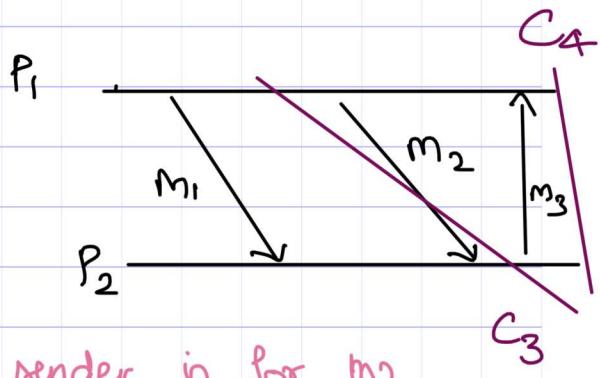
$m_2$  is in transit because only  $\text{send}(m_2)$  is present.

### Consistent Global State.

If receive event in GS, corresponding send events are present.

Transit would be there in consistent, but not in strongly consistent.

$$\begin{aligned} LS_1 &= \{\text{send}(m_1)\} \\ LS_2 &= \{\text{rec}(m_1), \text{rec}(m_2)\} \\ GS_1 &= LS_1 \cup LS_2 \\ \underline{\text{Inconsistent global state.}} \end{aligned}$$



$P_2$  cannot know who the sender is for  $m_2$  from the global state.

$$GS_2 = \{\text{send}(m_2), \text{send}(m_3), \text{rec}(m_3)\}$$

$GS_1 \rightarrow GS_2$  (causal ordering)

But causality fails here because in  $GS_1$  we have a  $\text{recv}(m_2)$  before the  $\text{send}(m_2)$  in  $GS_2$ .

But  $GS_1 \rightarrow GS_2 \Rightarrow$  causal ordering is inconsistent.

Maintaining strongly consistent states  
is practically impossible because  
of performance lag.

(Need to freeze, snapshot, resume operation  
each time event occurs)

## Consistent Global State Recording Protocol

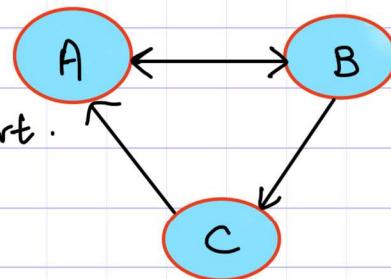
Any node can start the algorithm.

(Assume here A starts)

Send:

- (i) integrate msg. on channel & start.
- (ii) Sends a marker message M thro' all outgoing channels.

Karpot & Chandy



Receive:

- (i) Record States.
- (ii) follow send rule.

(Starting the algorithm is done by a leader election method)

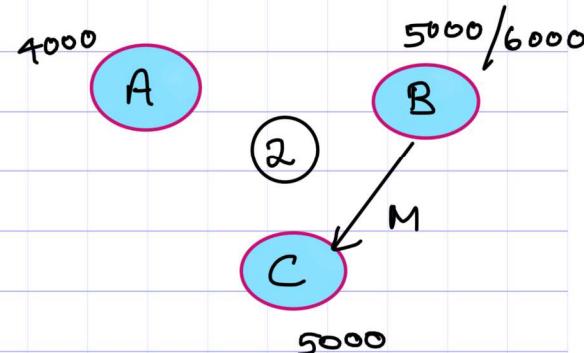
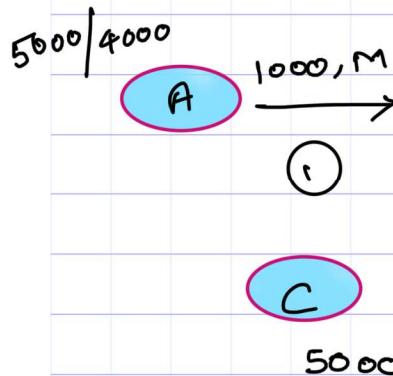
Termination:

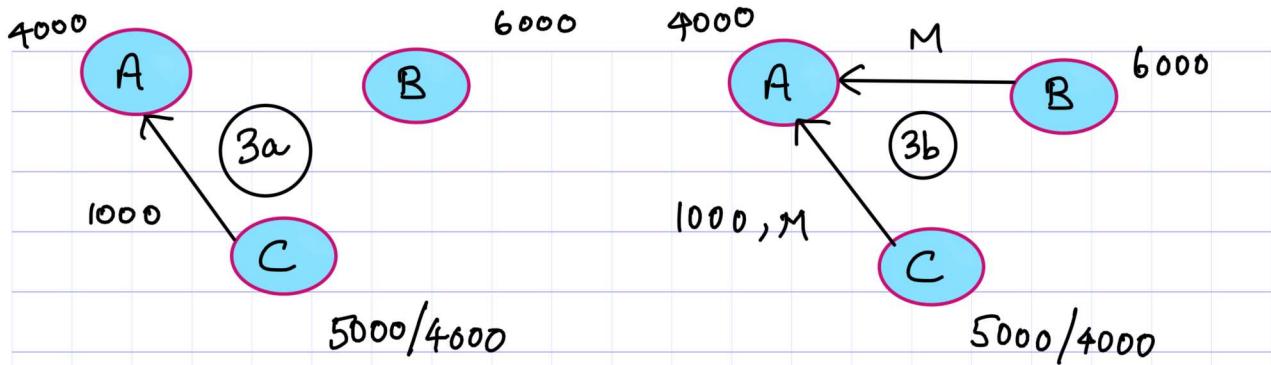
Once all incoming nodes receive the marker message M.

Assumptions:

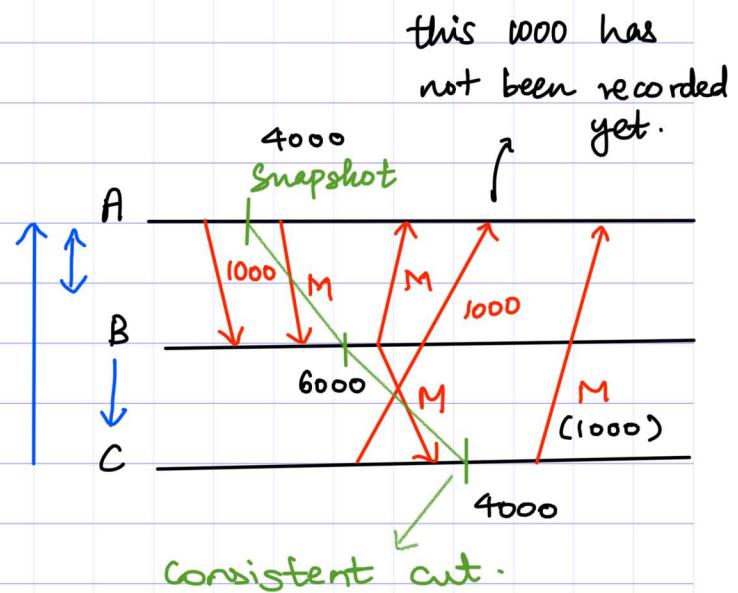
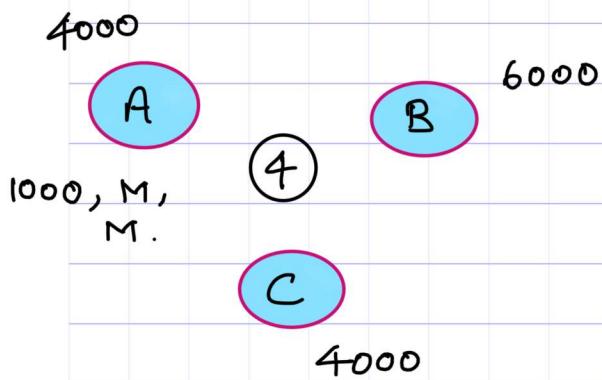
Reliable channel.

Ordering is FIFO (inherently causal order).

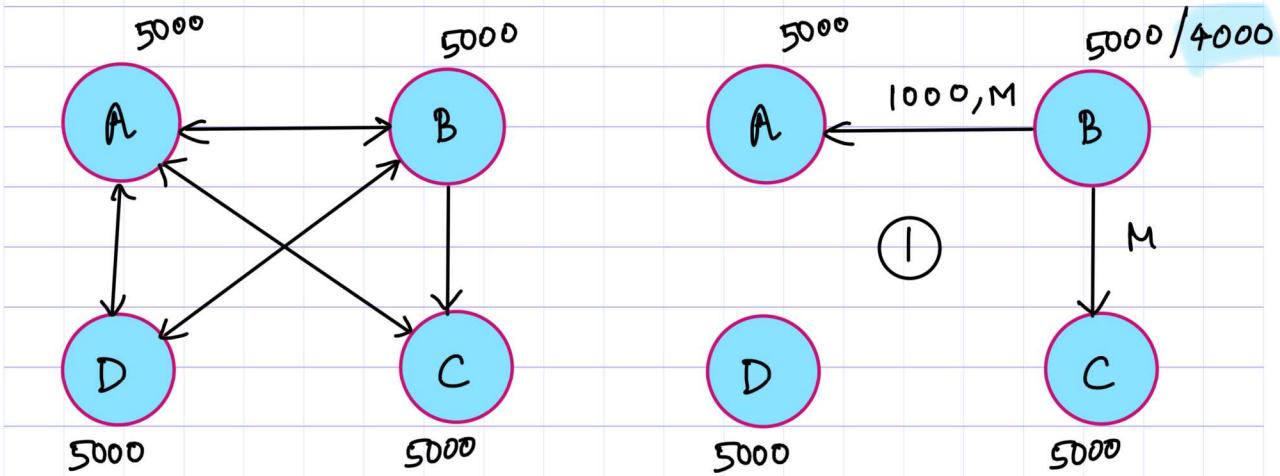




3a & 3b are concurrent events  
(happens parallelly)



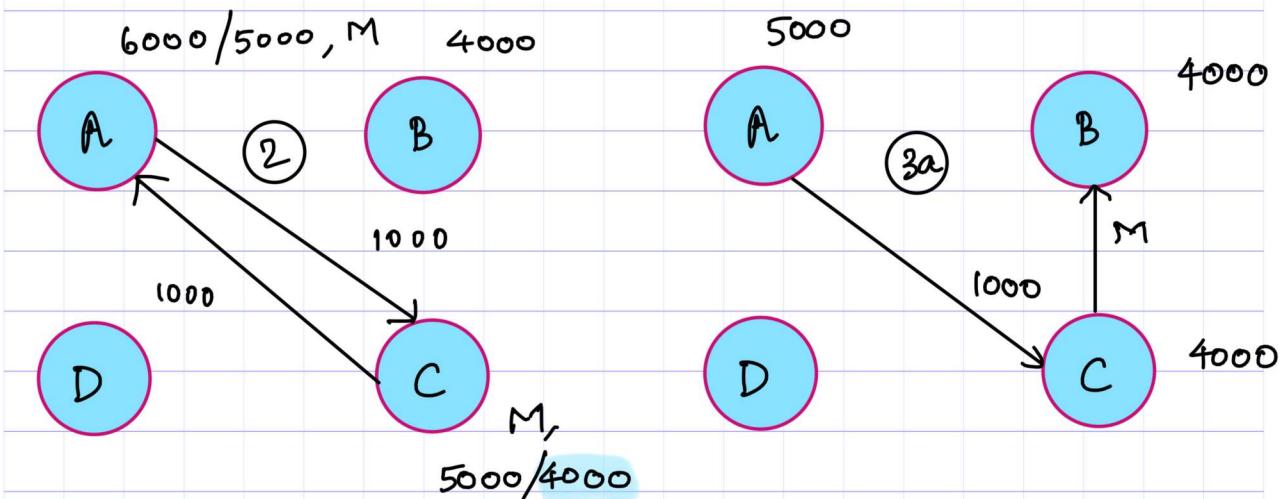
Example:



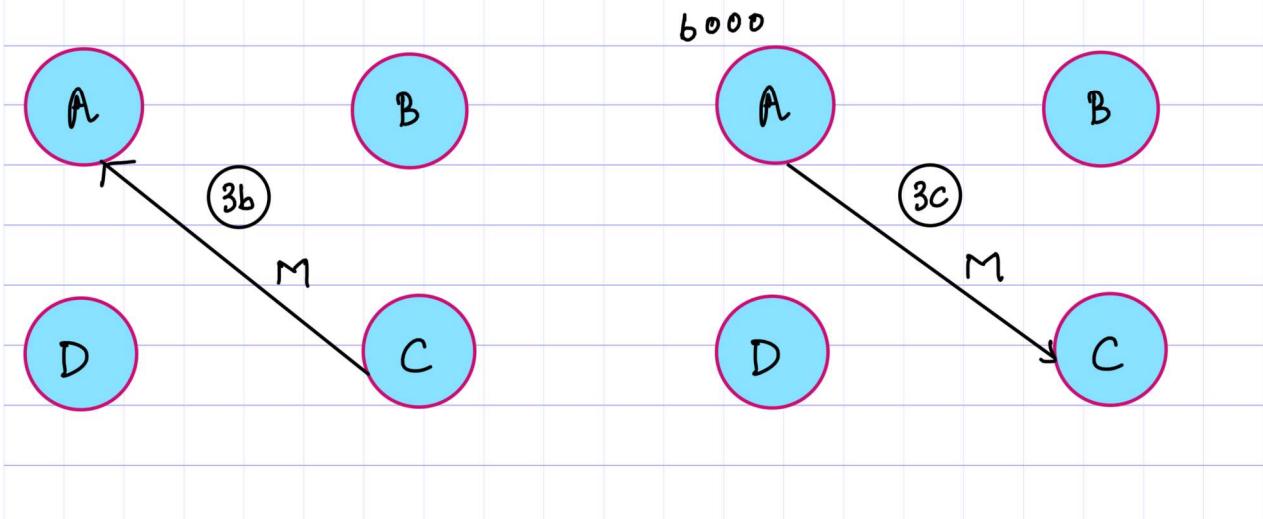
Leader was chosen

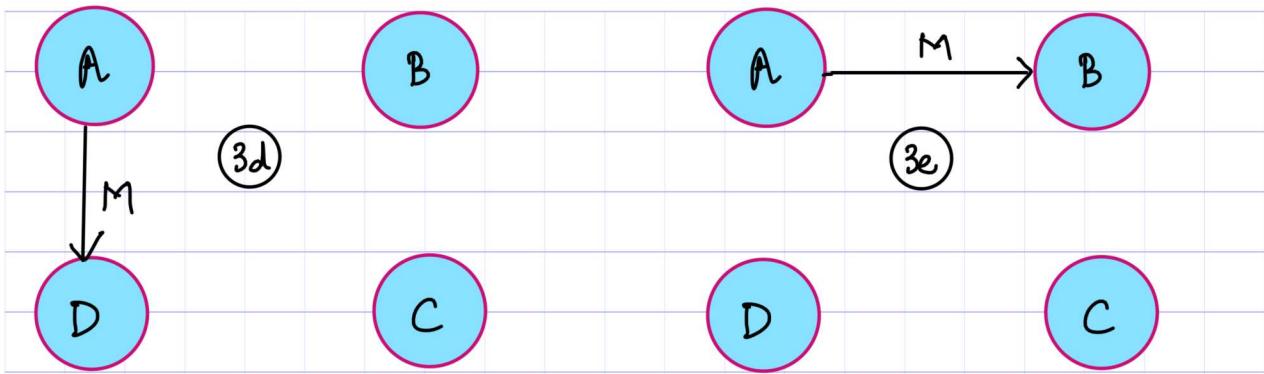
to be B.

B starts the algorithm

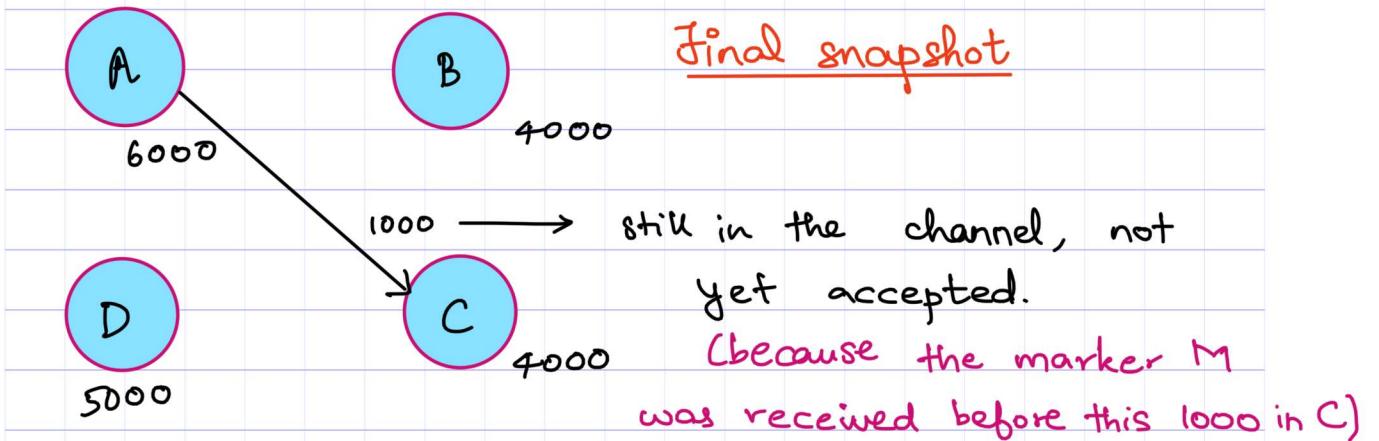


(C sends 1000 by the time it receives M from B)

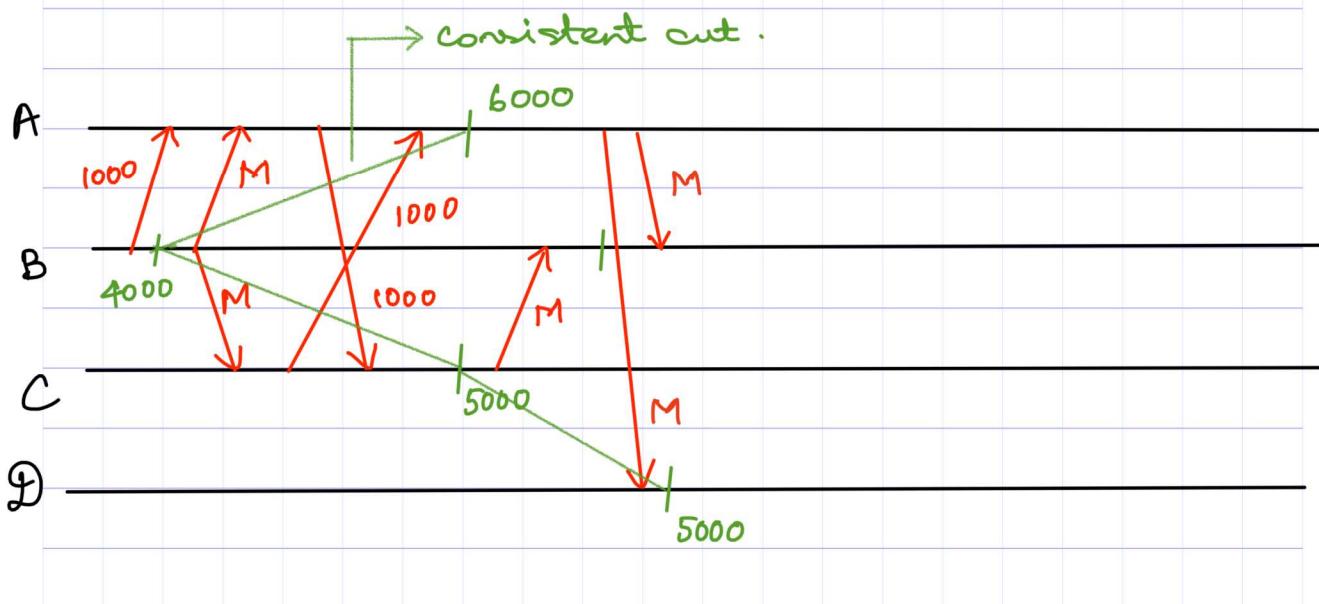




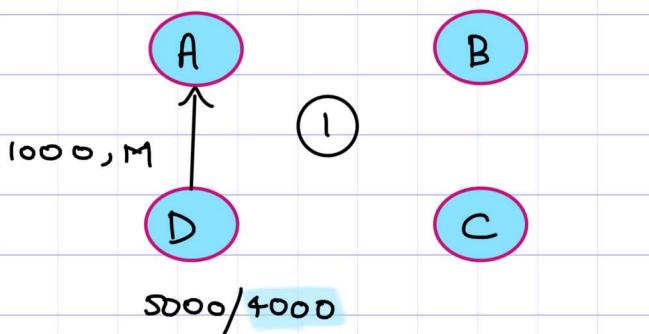
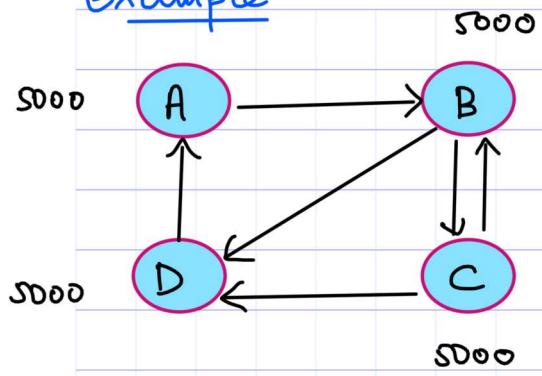
Sender receives marker M from all the incoming nodes, state is consistent.



This 1000 can be a part of the next global snapshot.



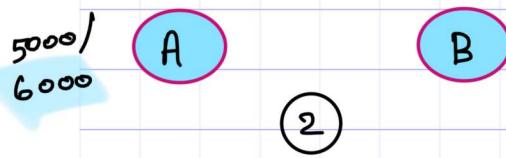
## Example



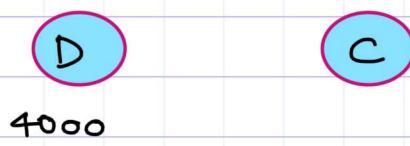
D is chosen as leader.

1000, Marker is sent in order by D.

Since D sends M, D records state as well.



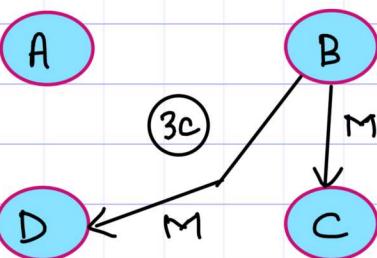
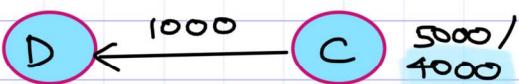
Marker recvd. by A,  
state recorded.



B recvs. M, B snapshots @ 5000,  
then recvs. 1000 from A.



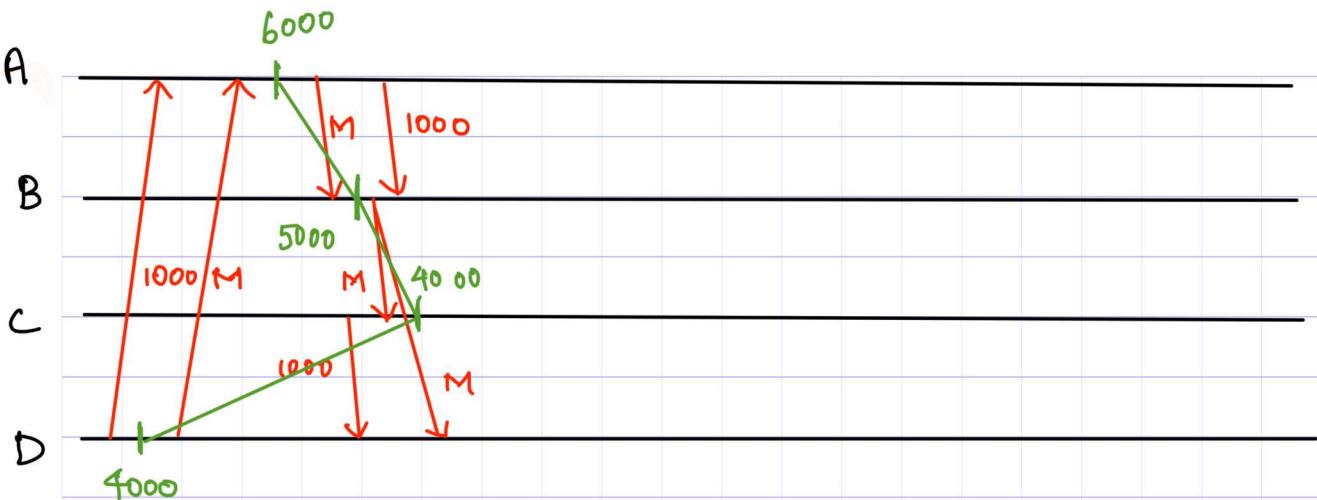
C concurrently sends  
1000 to D before  
getting M from B.



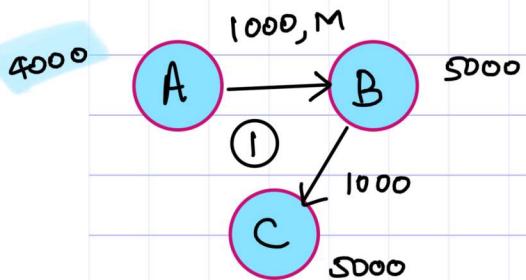
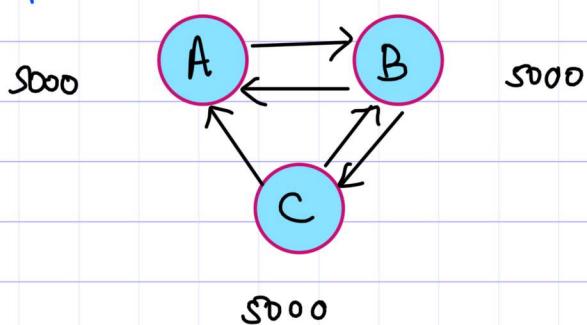
D receives 1000 followed  
by M from C. But  
D doesn't snapshot because  
it is the initiator.

4000/5000

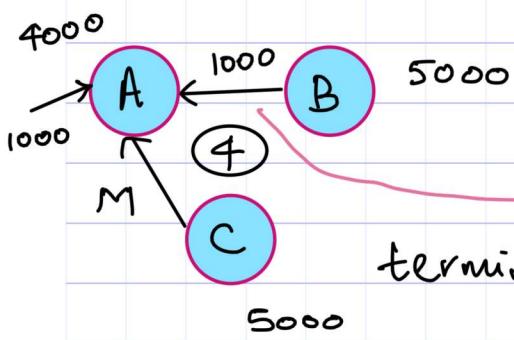
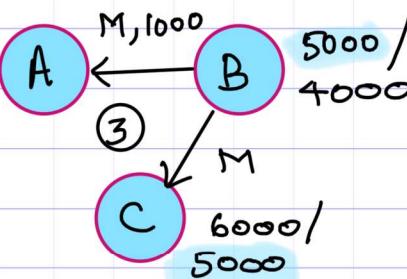
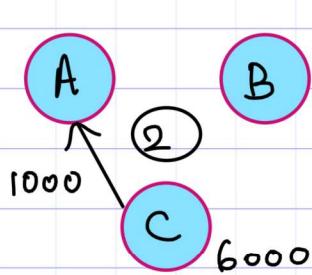
1000 from C → D is waiting in channel.



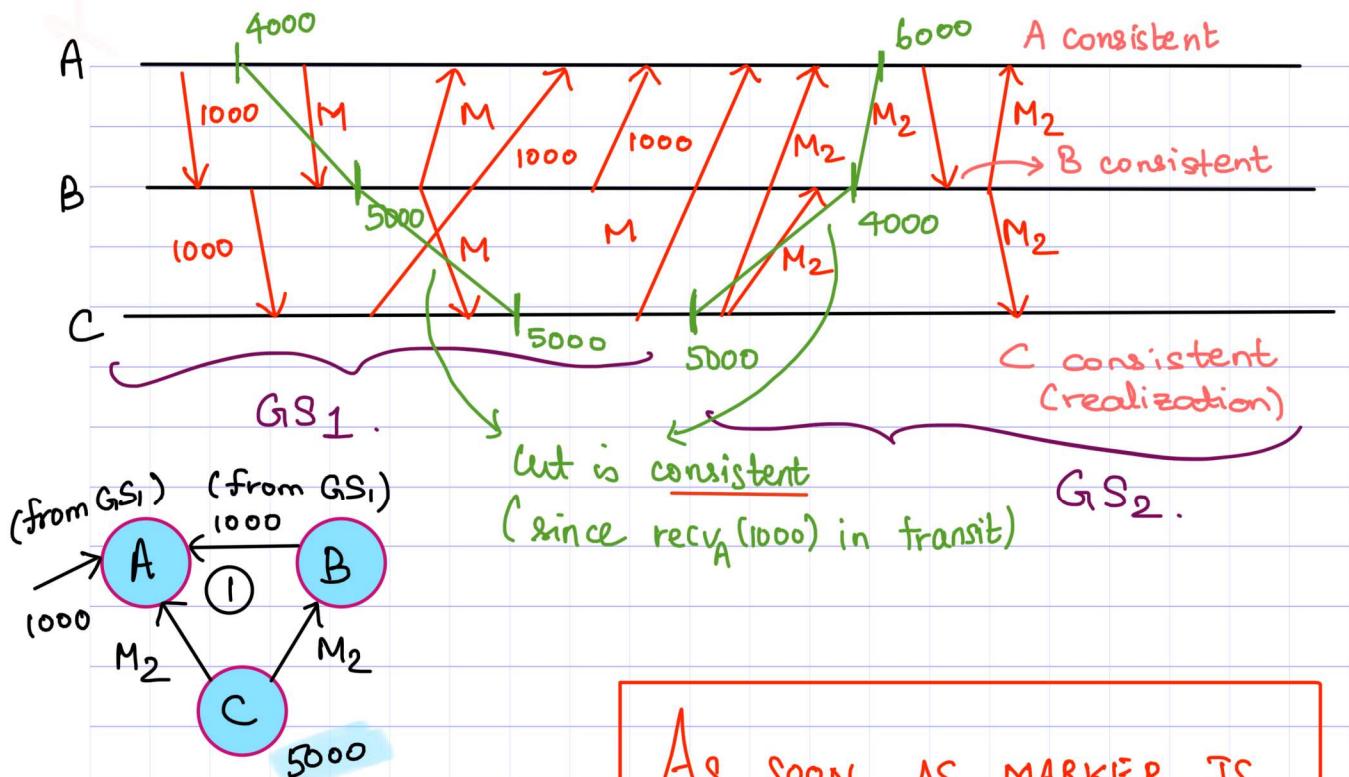
Example :



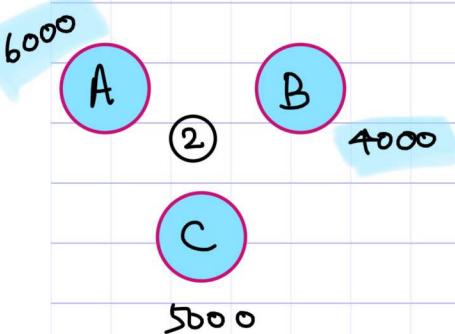
A sends 1000, M to B.  
B concurrently sends 1000 to C.



both send & receive are not  
terminated. recorded in snapshot  
yet because B already  
sent the marker.



C initiates Marker.

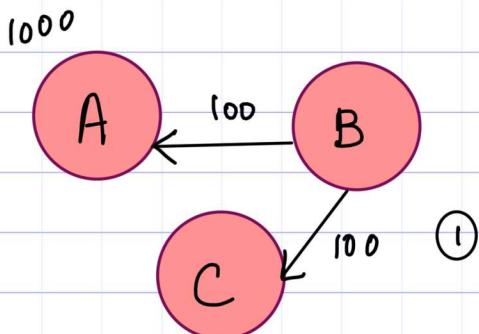
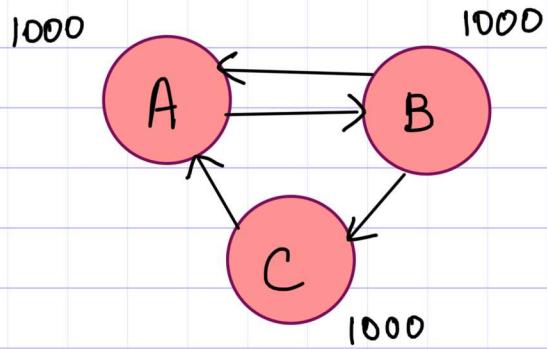


AS SOON AS MARKER IS  
RECEIVED TAKE A  
SNAPSHOT!

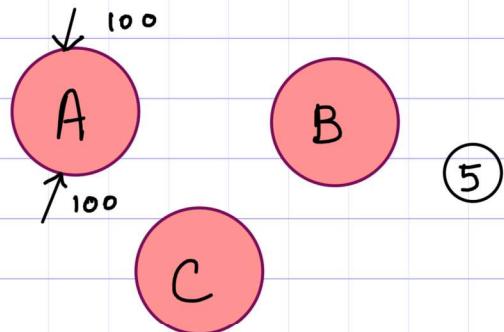
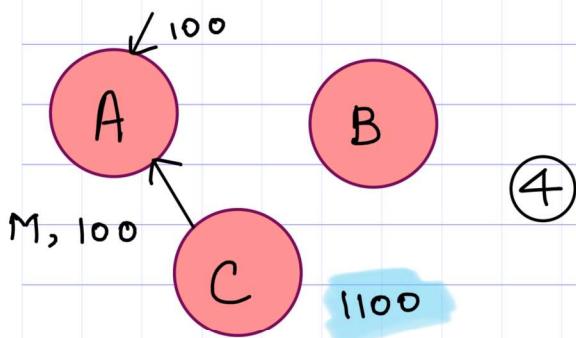
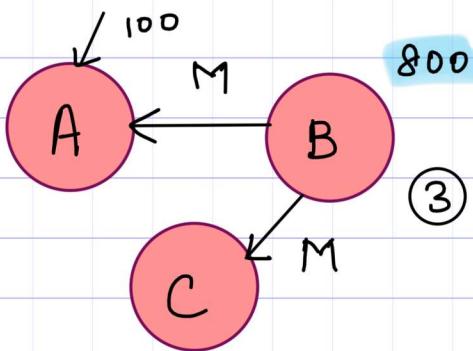
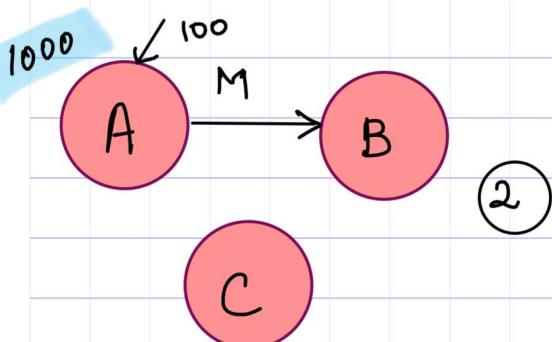
Both cuts are not  
strongly consistent because  
the receive & send of 1000  
from  $C \rightarrow A$  are recorded in  
2 separate snapshots.

Send  $B$  (1000) to  $A$  and receive  $A$  (1000) from  $B$   
will be recorded as part of  $AS_2$ .

## Global State Recording Algorithm

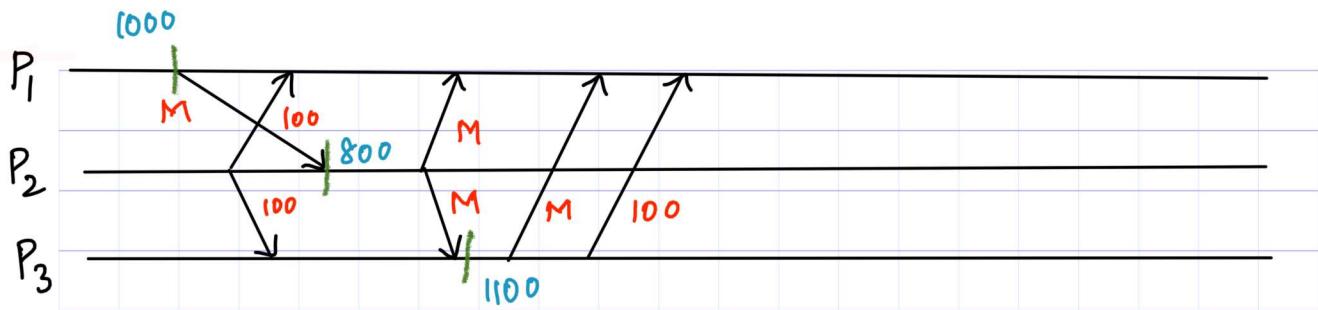


Transactions  $B \rightarrow A$   
 $B \rightarrow C$ .



Recording complete.

State is consistent.

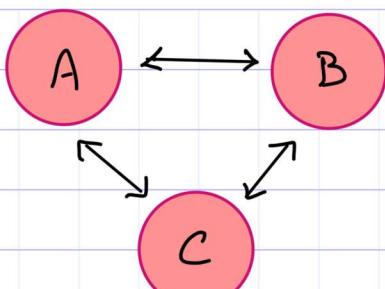


If  $P_1$  initiates the algorithm again,

$P_1$ 's state is recorded as 1200

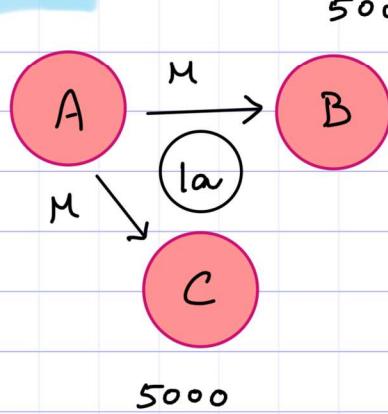
(1000 + 100 from  $P_2$  + 100 from  $P_3$ ).





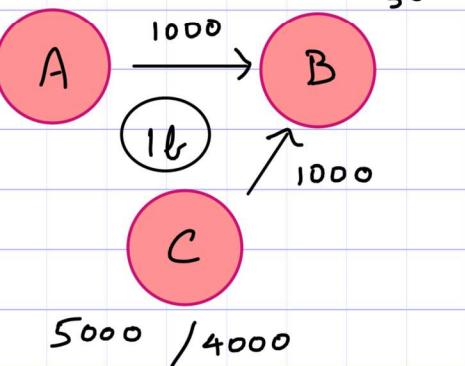
A initiates algorithm.

5000



5000

5000

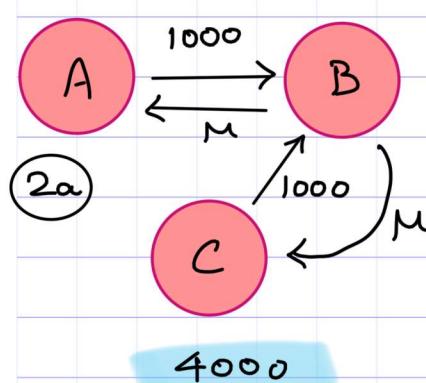


5000

1000

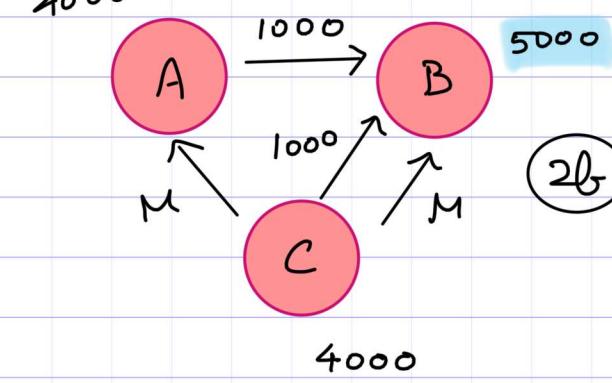
5000 / 4000

4000



4000

1000



5000

2b

