

Distributed System

Dhanya Krishnan

Prof. Shahul Hameed



Unit - 1 Introduction

Problems in distributed systems?

- i) Absence of global memory
- ii) Absence of global clock

Need?

- ↳ To accomplish a common task
- ↳ Network solutions worry abt resource sharing, not load sharing

Cloud → example of DS — uses WSDL for communication

Addressing problem 1: Lack of global / shared memory

Deadlocks in DS

- ↳ State can change during message passing overhead. Deadlock detection becomes inconsistent

Characteristics of DS

- ↳ No common physical clock
 - ↳ No common shared memory
 - ↳ Geographical separation
 - ↳ autonomous & heterogeneous.
- ↳ processes are connected via a communication network.

Motivation for DS

- ↳ Focus on resource sharing
- ↳ Inherently distributed computation
- ↳ Access to remote resources
- ↳ Inc performance/cost ratio
- ↳ Reliability / Scalability / Modularity / Expandability.

Challenges

i) System Perspective

- ↳ Communication mechanism - RPC, ROI
- ↳ Design of distributed programs, thread / process management
- ↳ universal nomenclature to easily identify resource/process
- ↳ Synchronization
- ↳ Data storage & access
- ↳ Consistency & replication
- ↳ Fault-tolerance
- ↳ Security & Transparency
- ↳ API building
- ↳ Scalable & modular algo.

ii) Algorithm Perspective

- ↳ Maintaining global state & time
- ↳ Co-ordinating mechanisms
- ↳ Group communication & ordered msg delivery
- ↳ Monitor events & predicates
- ↳ Distributed program design
- ↳ Debugging programs
- ↳ Shared memory abstraction
- ↳ Reliable & fault tolerant
- ↳ Real-time scheduling

Applications

- ↳ Mobile Systems
- ↳ P2P Comm.
- ↳ Sensor Networks
- ↳ Distributed data mining
- ↳ Grid computing
- ↳ Pub/Sub, CDN

Clocks are impractical:

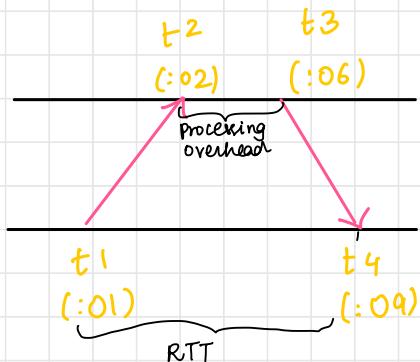
↳ Difficult to sync

↳ Timezones

Assumptions for clock syncing:

We need to assume common channel is error-free & is FIFO

e.g.:



true RTT =

$$\text{RTT} - \text{proc. overhead} = (t_4 - t_1) - (t_3 - t_2) \\ = 8 - 4 = 4$$

$$\text{true } \frac{\text{RTT}}{2} = \frac{4}{2} = 2$$

$$t_4 \rightarrow \text{rcvd time stamp} + \frac{\text{RTT}}{2} = 6 + 2 = 8 \leq t_4$$

Cutiaian's algorithm

$$T_{\text{new}} = T_{\text{curr.}} + \frac{\text{RTT}}{2}$$

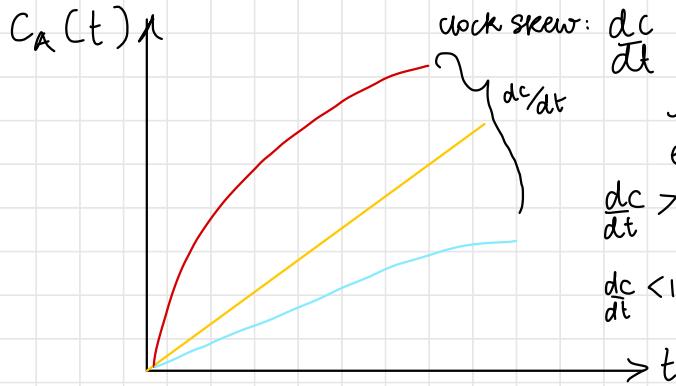
Clock drift → rate at which clock ticks

for e.g. quartz clock: 10^{-6} s ; diff of 1sec in 11.6 days

Clock skew → diff b/w 2 clocks at a pt in time.

Physical clock sync

$C_A(t) \cdot t \rightarrow$ for processor 'A', at time 't', clock = t
 it is perfectly synced (ideal case)



clock skew: $\frac{dC}{dt}$: change of clock value wrt time / some other clock

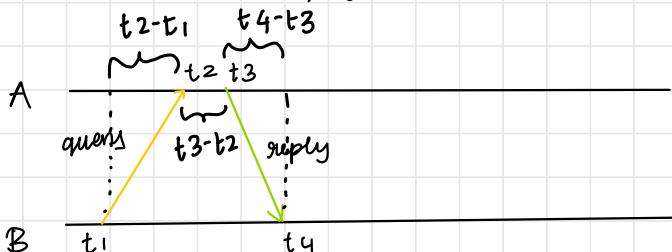
if eq 1, (synchronized)
 else, not in sync

$\frac{dC}{dt} > 1$ (clock runs faster)

$\frac{dC}{dt} < 1$ (clock slower)

$\frac{d^2C}{dt^2} \rightarrow$ clock drift

B wants to sync w/A



$t_2 - t_1 \rightarrow$ time delay (spent in channel / A is busy)
 / prop.

$t_3 - t_2 \rightarrow$ time server takes to process

B has to realise delays in the system

$$(t_4 - t_1) - (t_3 - t_2) = \text{transmission delay}$$

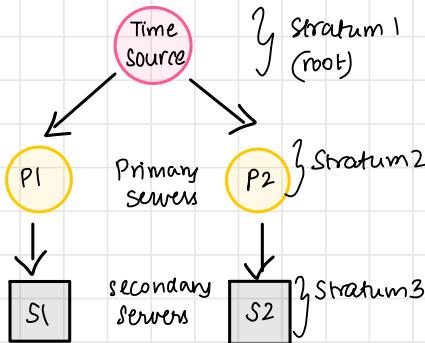
$$\frac{(t_4 - t_1) - (t_3 - t_2)}{2} = \text{avg trans. delay.}$$

Mechanisms to sync

- ↳ Cristian's Algorithm $(T_{new} = T_{server} + RTT/2)$
- ↳ Berkley's Algorithm
- ↳ Network Time Protocol

Network Time Protocol

↳ Clients across Internet sync to UTC despite msg delay.



multicast mode: time sent to all other nodes

procedure call mode: similar to Cristian's algo

Symmetric mode: for master servers : pair of servers exchange msgs.

Lamport's logical clock → uses a "happened before" relationship

↳ events can be processing/communicating

$$S = \{a, b\}$$

$a \rightarrow b$: binary relation

Strict partial order

i) irreflexive : an event can't happen before itself $a \not\rightarrow a$

ii) anti-symmetric : if $a \rightarrow b$ then $b \not\rightarrow a$

iii) Transitive : $a \rightarrow b \wedge b \rightarrow c$ then $a \rightarrow c$

each event has clk val (integer)

Rule1: $a \rightarrow b$, then $c_i(a) < c_j(b)$, both belong to same system/clock/processor

$$c_i \quad c_j$$

clock value

→ for processing

Rule2: $a \rightarrow b$, then $c_i = (c_i(m) + 1, c_j + 1)$

$$c_i(m) + 1, c_j + 1$$

↳ for communicating events

NOTE: if $a \rightarrow b$ $c(a) < c(b)$

but vice versa may not be valid

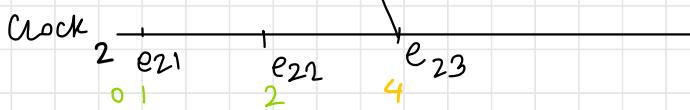
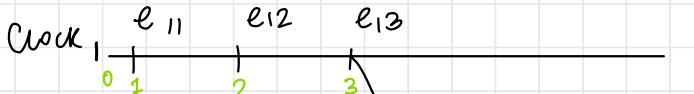
Y limitation of

Lamport's Logical clock

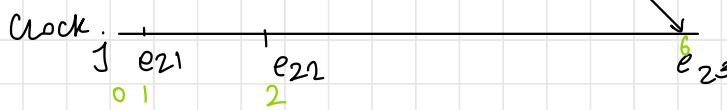
Clock not smart enough
to diff b/w processing & recr.
event

To update clock,

$$c_i(e) = c_i(e) + d \quad /d usually/$$



Using Rule 2



Limitation of Lamport's Logical Clock

We know that

$$\text{if } a \rightarrow b \Rightarrow c(a) \leq c(b)$$

but, if $c(a) \leq c(b) \not\Rightarrow a \rightarrow b$ (Inconsistent)

Vector Clocks

$$c_j[\text{event}] = c_j[\text{event}] + 1 \rightarrow \text{Processing event}$$

$$c_j[k] = \max \left(\underbrace{c_j[k] + 1}_{\text{clock of sender during msg}}, c_i[m] \right)$$

clock of rcv inc

For clock values t_a, t_b

$$\text{if } t_a = t_b \Rightarrow \nexists t_a[k] = t_b[k]$$

$$t_a \neq t_b \Rightarrow \exists t_a[k] \neq t_b[k]$$

Drawbacks of logical clocks.

↳ Overlooks concurrent events

↳ Limitation discussed - can't determine which event occurred in which process in what order.

↳ Needs more memory

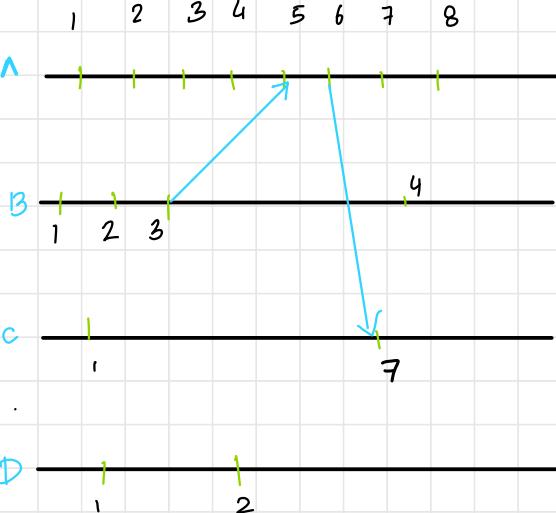
↳ No of processes may not be known in advance

NOTE: Vector Clocks are strongly consistent

$$\text{if } t_a < t_b \Rightarrow a \rightarrow b$$

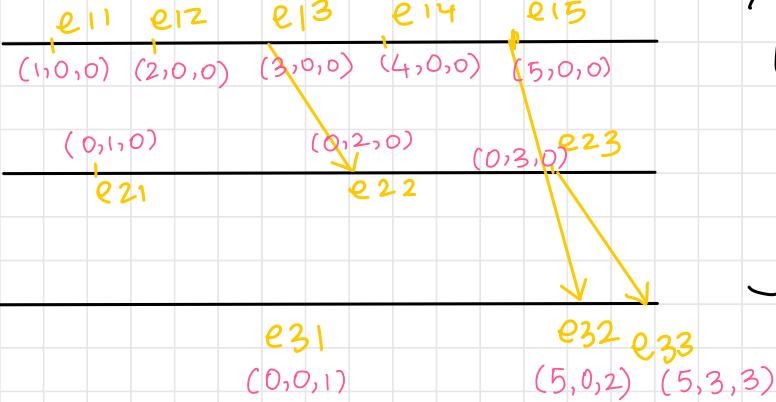
Solves Lamport clock problem

eg:



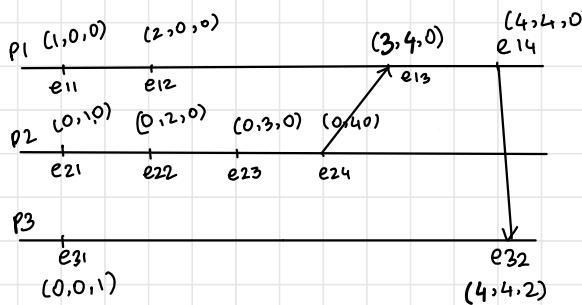
Lamport's clock

eg:



Vector
Clock

Vector Clock: vector of values rep each process (time-space diagram)

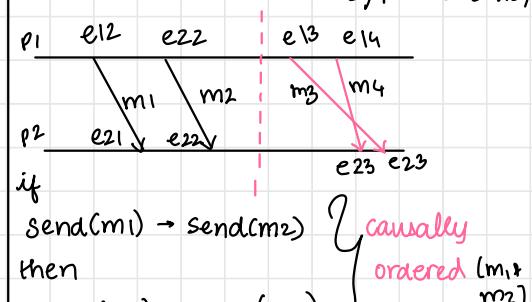


vector clock limitation

- ↳ large vectors when many processors
- ↳ suitable for LANs
ad hoc
- ↳ not for wireless (not fixed no. of processors) → need to tune

How to ensure ordering of messages

↳ Causal ordering (consider only msg pass events)

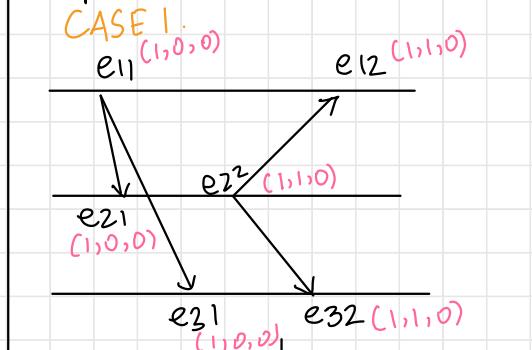


but m_3, m_4 are not

1) order important as meaning may be lost w/o it.

2) if shared resources exist, order it.

CASE 1.



sender rule (i)

i) inc i^{th} value

k update t

ii) broadcast

receiver's rule (j)

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

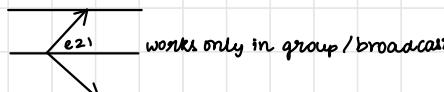
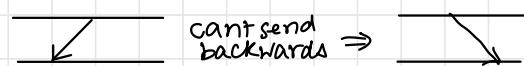
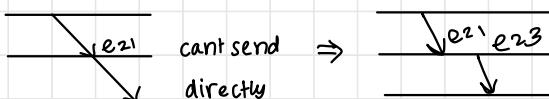
ii) $\forall k: c_j[k] \geq t_m[k]$

$k \neq i$

if holds, accept msg & update clock

$$\text{if } t_a \rightarrow t_b \times t_b \rightarrow t_a \Rightarrow t_a \sqcap t_b$$

Common mistakes to avoid



e₂₁ C₂(0,0,0)
 $t_m = (1,0,0)$

i=1, j=2

e₃₁ C₃(0,0,0)
 $t_m = (1,0,0)$

i=1, j=3

e₁₂ C₁(1,0,0)
 $t_m = (1,1,0)$

i=2, j=1

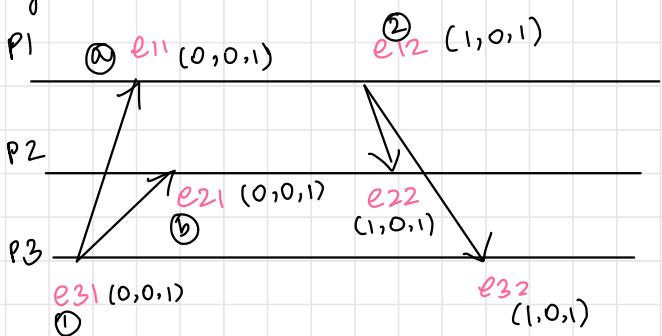
e₃₂ C₃(1,0,0)
 $t_m = (1,1,0)$

i=2, j=3

for receiver,
 rule 1 satisfied
 ✕ rule 2

NOTE: All Cases for Broadcasting

eg:



gender events

① e₃₁ C₃(0,0,1)

Receiver events

i=3, j=1

② e₁₁ C₁(0,0,0) at rcvr, C_i[i] = $\frac{t[i]-1}{m}$
 $t_m = (0,0,1)$

if diff > 1

P1 e₁₁(1,0) e₁₂(2,0) e₂₁ i=1, j=2

m_1

P2 (10,6) e₂₁

e₁₁ C₁(0,0,0)

$t_m = (0,0,1)$

C₂(0,0,0)

$t_m = (2,0,0)$

$t'_1 (0,0,1)$

m_1 has been missed

Rule 1 explanation

accept → msg delivered

receive → msg at buffer

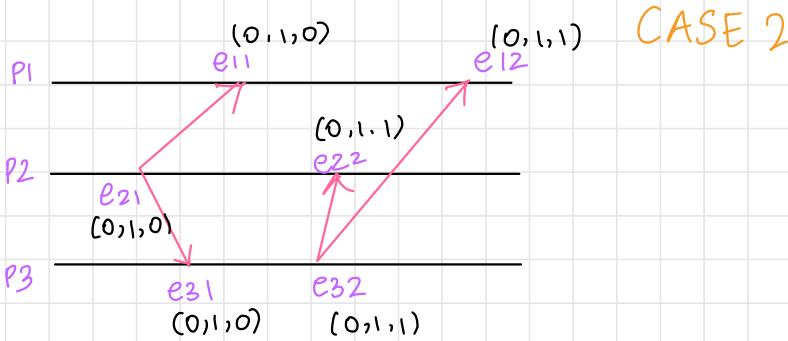
⑥ 11^w for e_{21} , $t_2 \rightarrow (0,0,1)$

⑦ $e_{22} C_2(0,0,1) \rightarrow t_2(1,0,1)$
 $i=1, j=2 \rightarrow$ implies p₂ has not yet rcv any msg from p₁
 $t_m = (1,0,1) \rightarrow t_2(1,0,1)$

⑧ 11^y for e_{32} , $t_3 \rightarrow (1,0,1)$

timestamp in Rule 2 implies that already p₁ knows a msg was rcvd from p₃

equal/
shared
knowledge



Send events

i) $e_{21} C_2(0,1,0)$

receive events

a) $e_{31} C_3(0,0,0) i=2$
 $t_m(0,1,0) j=3$
 $\Rightarrow t_3(0,1,0)$

b) 11^w $e_{11} C_1(0,0,0) i=2$
 $t_m(0,1,0) j=1$
 $\Rightarrow t_1(0,1,0)$

c) $e_{22} C_2(0,1,0) i=3$
 $t_m(0,1,1) j=2$
 $\Rightarrow t_2(0,1,1)$

d) $e_{12} C_1(0,1,0) i=3$
 $t_m(0,1,1) j=1$
 $\Rightarrow t_1(0,1,1)$

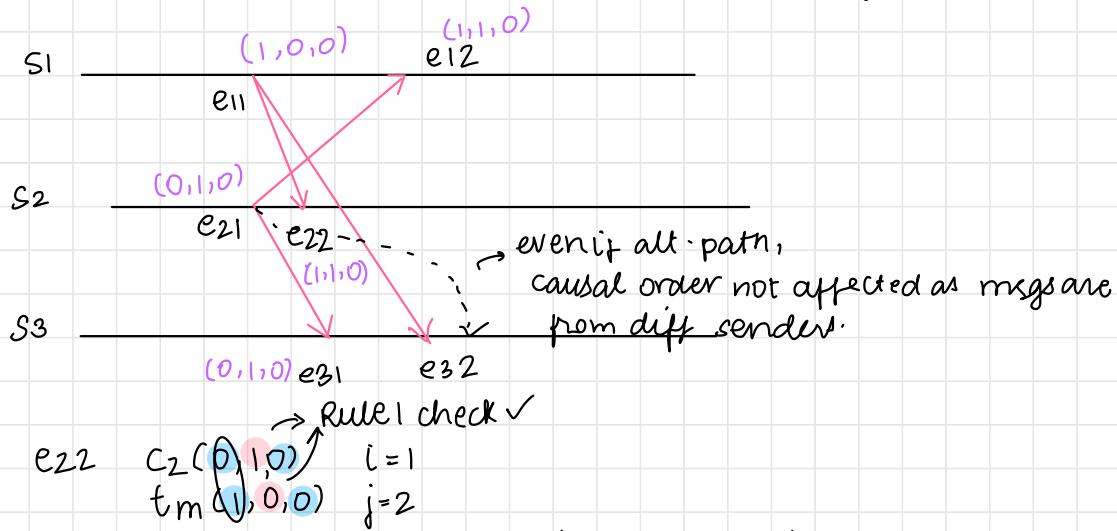
ii) $e_{32} C_3(0,1,1)$

NOTE: Rule 2 checks for non-sender clock value consistency b/w sender & receiver.

But why is there a greater than symbol?

$$c_j[K] \geq t_m[K] \quad K \neq \text{sender}$$

Because i) rcvr could've initiated a send ii) concurrent msg could be sent



Rule 2 inference: sender has not received msg initiated by receiver yet

to update clock, merge the values & take greatest)

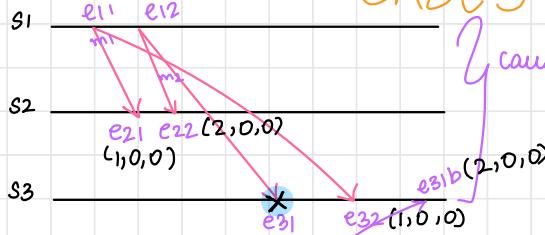
$$e_{22} \quad t_2(1,1,0)$$

III for $e_{12} \quad (1,1,0)$

$$\begin{array}{ll} e_{31} \quad c_3(0,0,0) & i=2 \\ t_m(0,1,0) & j=3 \\ t_3(0,1,0) \end{array}$$

$$\begin{array}{ll} e_{32} \quad c_3(0,1,0) & i=1 \\ t_m(1,0,0) & j=3 \\ t_3(1,1,0) \end{array}$$

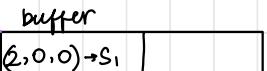
CASE 3



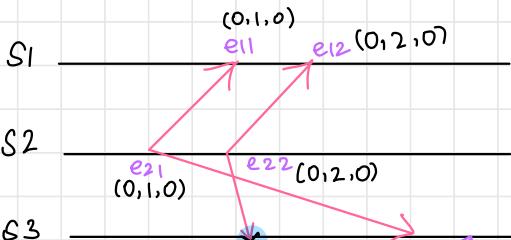
e_{21} $c_2(0,0,0)$ $i=1$
 $t_{m}(1,0,0)$ $j=2$
 $t_2(1,0,0)$

Π^{w} e_{22} $c_2(1,0,0)$
 $t_{m}(2,0,0)$
 $t_2(2,0,0)$

e_{31} $e_{3}(0,0,0)$ rule violated as diff > 1
 $t_{m}(2,0,0) \Rightarrow$ decline message & move to buffer
do not update clock now



e_{32} $c_3(0,0,0)$ $i=1$
 $t_{m}(1,0,0)$ $j=3$
 \Rightarrow accept & update clock
retrieve msg from buffer
 e_{31b} $c_3(1,0,0)$ $i=1$
 $t_{m}(2,0,0)$ $j=3$
 \Rightarrow accept & update



e_{11} $c_1(0,0,0)$
 $t_{m}(0,1,0)$
 $\Rightarrow c_1(0,1,0)$
 $i=2$ & $j=1$

e_{31} $c_3(0,0,0)$ $i=2$
 $t_{m}(0,2,0)$ $j=3$
buffer

$(0,2,0), S_2$	
----------------	--

Stored in buffer
& clock not updated

e_{12} $c_1(0,1,0)$ $i=2$
 $t_{m}(0,2,0)$ $j=1$
 $c_1(0,2,0)$

e_{32} $c_3(0,0,0)$
 $t_{m}(0,1,0)$
 $c_3(0,1,0)$
 $i=2$ & $j=3$

e_{31b} $c_3(0,1,0)$ $i=2$
 $t_{m}(0,2,0)$ $j=3$
 $c_3(0,2,0)$

NOTE: Check buffer on
every receive event &
obtain in FIFO order

CASE 3 extended

S_1 $(1, 0, 0)$ $(2, 0, 0)$
 e_{11} e_{12} $e_{13} (3, 0, 0)$

S_2 e_{21} e_{22} e_{23}

S_3 causal e_{31} e_{32} e_{33} $(1, 0, 0)$ $(2, 0, 0)$ $(3, 0, 0)$

non-causal

e_{31} $i=1, j=3$

$C_3 (0, 0, 0)$

$t_m (2, 0, 0)$

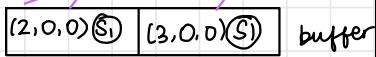
no clock update

e_{32} $i=1, j=3$

$C_3 (0, 0, 0)$

$t_m (3, 0, 0)$

no clock update

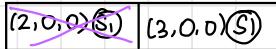


e_{33} $C_3 (0, 0, 0)$ $i=1$
 $t_m (1, 0, 0)$ $j=3$
 $\Rightarrow C_3 (1, 0, 0)$ $y \text{ accepted}$

e_{31b} $C_3 (1, 0, 0)$ $i=1$

$t_m (2, 0, 0)$ $j=3$

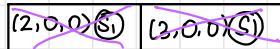
$\Rightarrow C_3 (2, 0, 0)$



e_{32b} $C_3 (2, 0, 0)$ $i=1$

$t_m (3, 0, 0)$ $j=3$

$\Rightarrow C_3 (3, 0, 0)$



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

S_1 e_{11} e_{12}

S_2 e_{21} $(0,1,0)$

e_{22} $(1,1,0)$

S_3 e_{31}

e_{32} $(0,1,0)$

e_{31b} $(1,1,0)$

$t_m(1,1,0)$ $i=1, j=3$



CASE - 4

① e_{11} $c(0,0,0)$

$t_m(0,1,0)$

$(i=2, j=1)$

$c(0,1,0)$

② e_{22} $c(0,1,0)$

$t_m(1,1,0)$

$(i=1, j=2)$

③

e_{31} $c(0,0,0)$ Rule 2 violated

$t_m(1,1,0)$ $i=1, j=3$

④ e_{32} $c(0,0,0)$ $i=2$

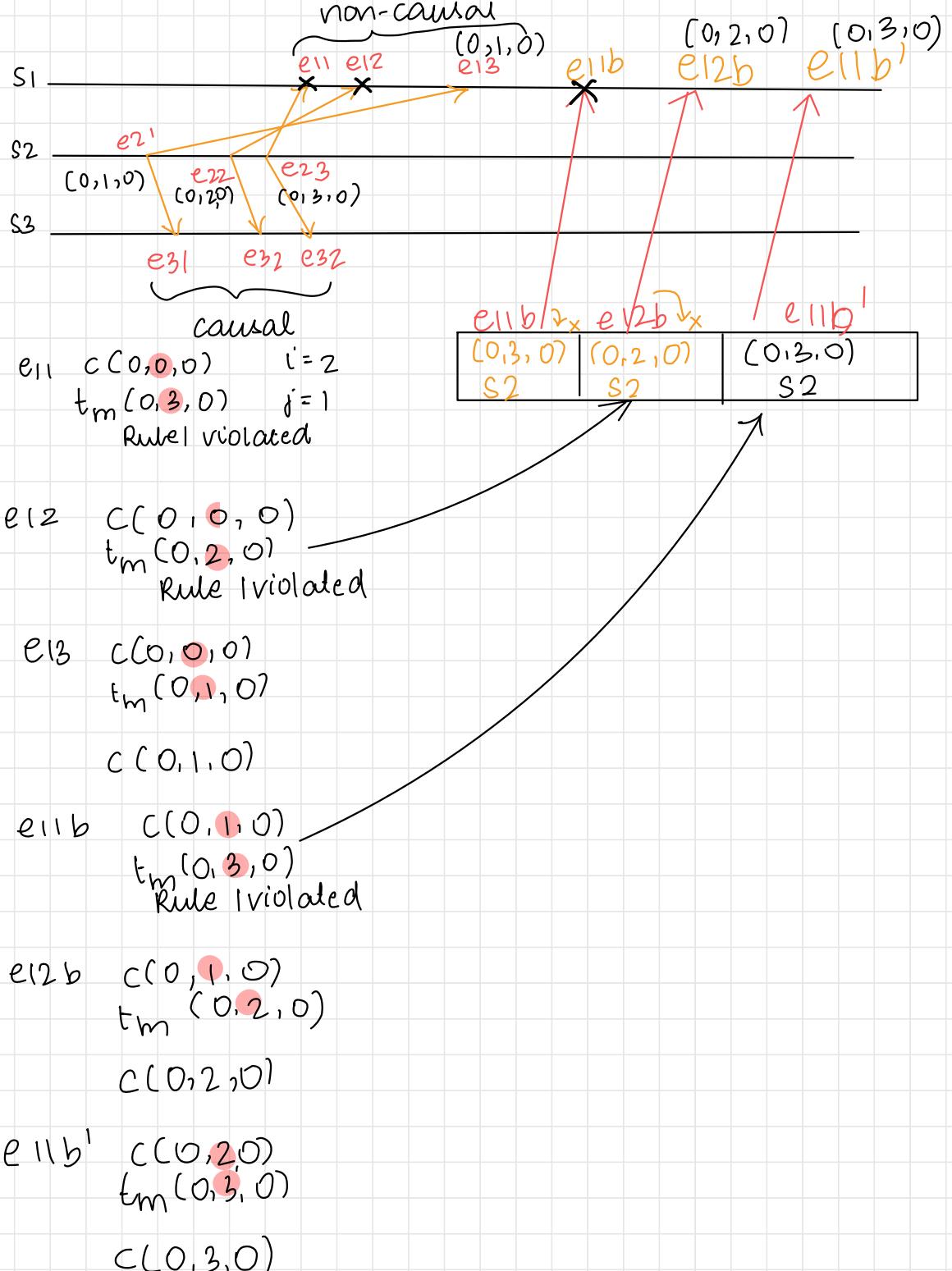
$t_m(0,1,0)$ $j=3$

$c(0,1,0)$

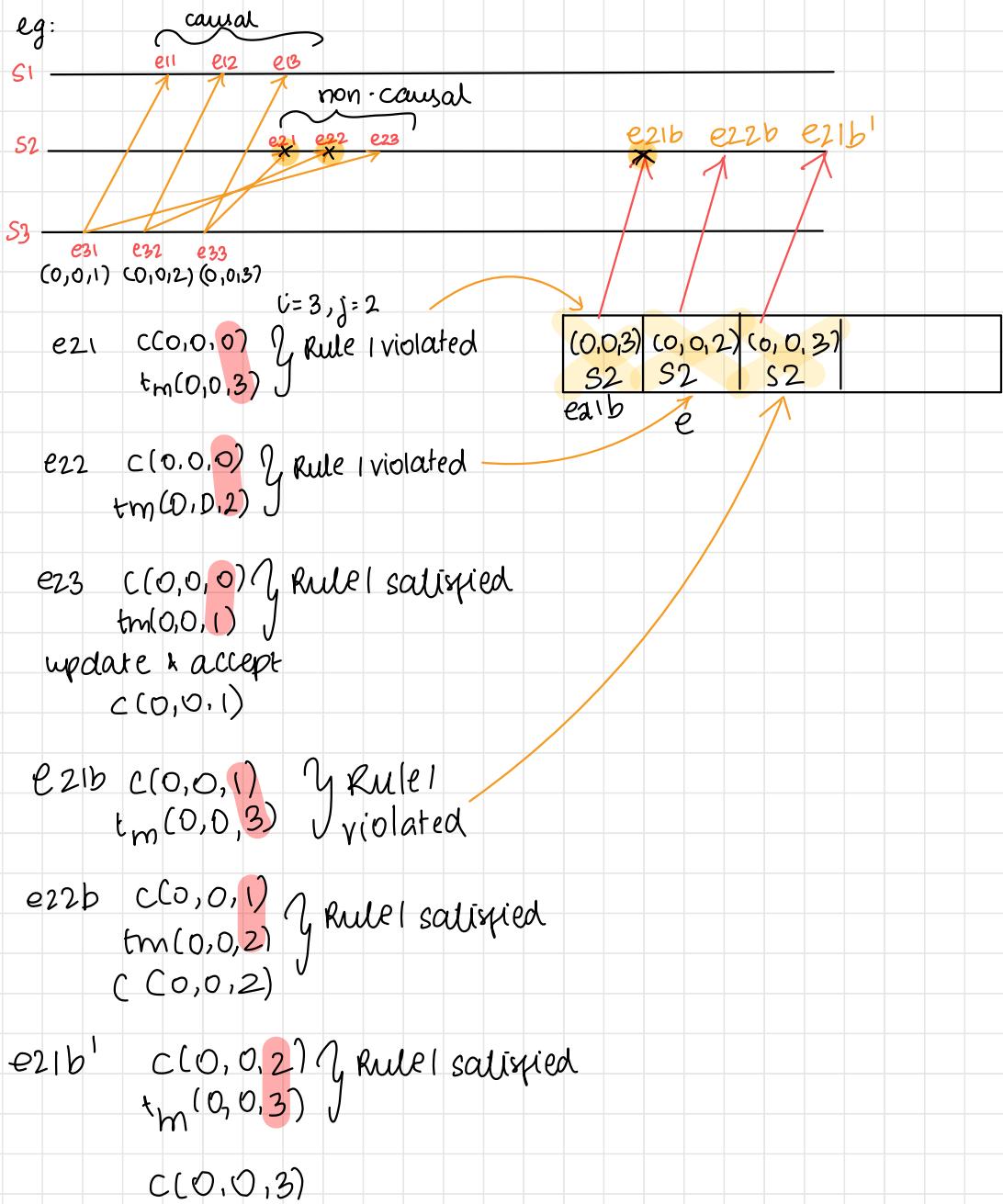
⑤ e_{31b} $c(0,1,0)$

$t_m(1,0)$

$c(1,1,0)$



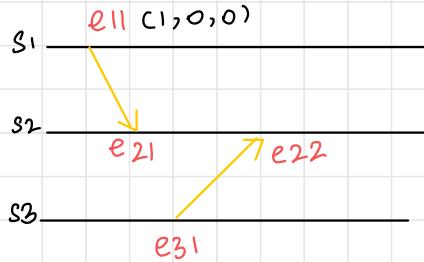
eg:



Causal Ordering for non-broadcasting

Send rule

- i) inc clock
- ii) send timestamp t_m & vector v_m
- iii) update v_i



rcv rule

- has other cases
i) check if $v_j[j]$ is \underline{x} & $v_m[j]$ is x
rcv msg, update clock & vector

$S_1 \rightarrow e_{11}$ → this is the msg to S_2
 $\{C_1(1, 0, 0)\}$ → initially empty

vector → Process id & timestamp

receiver id → sender's time

v_i after updation (after msg sent)

$v_i(x, 100, x)$

changed rcvr entry

$S_2 \rightarrow e_{21}$ updated clk value entered in
 $C_2(0, 0, 0)$ rcr position for msg send

$v_2(x, x, x)$

rcv msg: $t_m(1, 0, 0)$ $v_i(x, x, x)$

updation:

$C_2(1, 0, 0)$

$v_2(x, 100, x)$ → update in self position for msg rcv

$S_3 \rightarrow e_{31}$

$\{C_3(0, 0, 1)\}$ → sent to S_2
 $\{V_3(x, x, x)\}$

update V_3

$V_3(x, 001, x)$

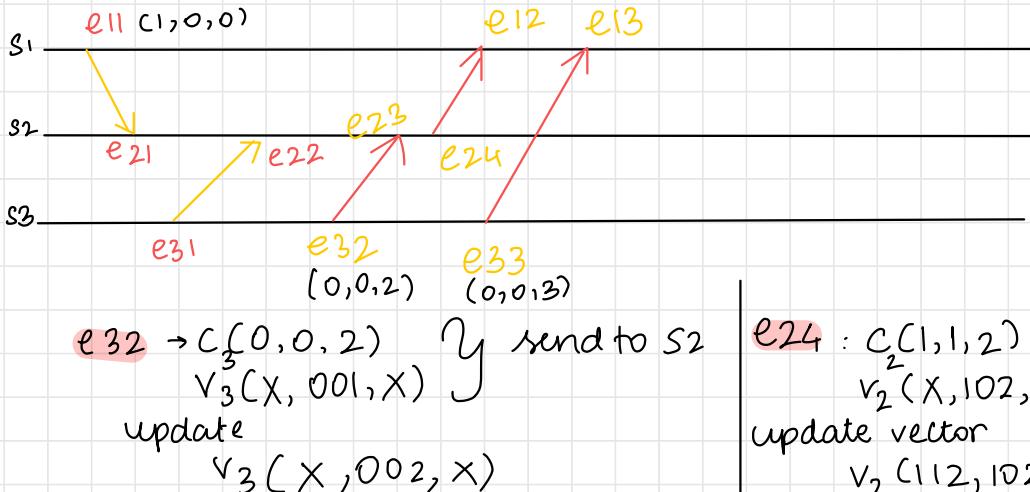
$S_2 \rightarrow e_{22}$

$C_2(1, 0, 0)$

$V_2(x, 100, x)$

RCV msg: $tm(0,0,1)$ $V_m(X, X, X)$
 updatation
 $C_2(1,0,1)$
 $V_2(X, 101, X)$

Extension of above



$e_{33} \rightarrow C_3(0,0,3)$ γ send to S_1
 $V_3(X, 002, X)$
 update
 $V_3(003, 002, X)$

$e_{23} \rightarrow C_2(1,0,1)$
 $V_2(X, 101, X)$

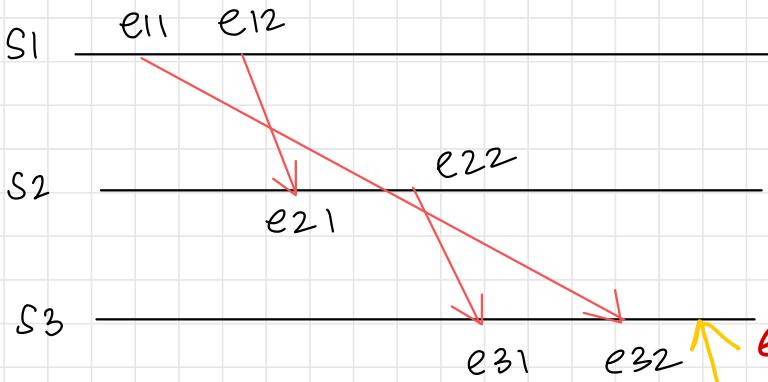
RCV msg : $C_m(0,0,2)$
 $V_m(X, 001, X)$

update : $C_2(1,0,2)$
 $V_2(X, 102, X)$

$e_{24} : C(1,1,2)$ γ send to S_1
 $V_2(X, 102, X)$
 update vector
 $V_2(112, 102, X)$

$e_{12} : C_1(1,0,0)$
 $V_1(X, 100, X)$
 RCV msg
 $Cm(1,1,2)$
 $Vm(X, 102, X)$
 update : $V_1(X, 102, X)$
 $C_1(1,1,2)$
 $V_1(112, 102, X)$

$e_{13} : C_1(1,1,2)$
 $V_1(112, 102, X)$
 RCV msg
 $Cm(0,0,3)$
 $Vm(X, 002, X)$
 update : $C_1(1,1,3)$
 $V_1(113, 102, X)$



e11 $C_1(1,0,0)$ send to S3
 $V_1(X, X, X)$
update
 $V_1(X, X, 100)$

e12 $C_1(2,0,0)$ send to S2
 $V_1(X, X, 100)$
update
 $V_1(X, 200, 100)$

e21 $C_2(0,0,0)$
 $V_2(X, X, X)$
rcv msg: $C_m(2,0,0)$
 $V_m(X, X, 100)$
update: $C_2(2,0,0)$
 $V_2(X, 200, 100)$

e22 $C_2(2,1,0)$ send to S3
 $V_2(X, 200, 100)$
update
 $V_2(X, 200, 210)$

e31 $C_3(0,0,0)$
 $V_3(X, X, X)$ violation
rcv msg: $C_m(2,1,0)$
 $V_m(X, 200, 100)$

$C_m(2,1,0)$	
$V_m(X, 200, 100)$	

e32: $C_3(0,0,0)$
 $V_3(X, X, X)$
rcv msg: $C_m(1,0,0)$
 $V_m(X, X, X)$
update $C_3(1,0,0)$
 $V_3(X, X, 100)$

e31b: $C_3(1,0,0)$
 $V_3(X, X, 100)$
rcv msg: $C_m(2,1,0)$
 $V_m(X, 200, 100)$
update: $C_3(2,1,0)$
 $V_3(X, 200, 210)$



e_{33} $c_3(2,1,1)$ γ send to
 $v_3(x, 200, 210) \checkmark$ S2
update
 $v_3(x, 211, 210)$

e_{34} $c_3(2,1,2)$ γ send to
 $v_3(x, 211, 210) \checkmark$ S2
update
 $v_3(x, 212, 210)$

e_{23} $c_2(2,1,0)$
 $v_2(x, 200, 210)$
rcv msg: $c_m(2,1,2)$
 $v_m(x, 211, 210)$

$c_m(2,1,2)$	
$v_m(x, 211, 210)$	

e_{24} $c_2(2,1,0)$
 $v_2(x, 200, 210)$
rcv msg: $c_m(2,1,1)$
 $v_m(x, 200, 210)$
update: $c_2(2,1,1)$
 $v_2(x, 211, 210)$

e_{23b} $c_2(2,1,1)$
 $v_2(x, 211, 210)$
rcv msg: $c_m(2,1,2)$
 $v_m(x, 211, 210)$
update: $c_2(2,1,2)$
 $v_m(x, 212, 210)$