

Chuyên đề hệ thống phân tán



KHOA CÔNG NGHỆ THÔNG TIN
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN

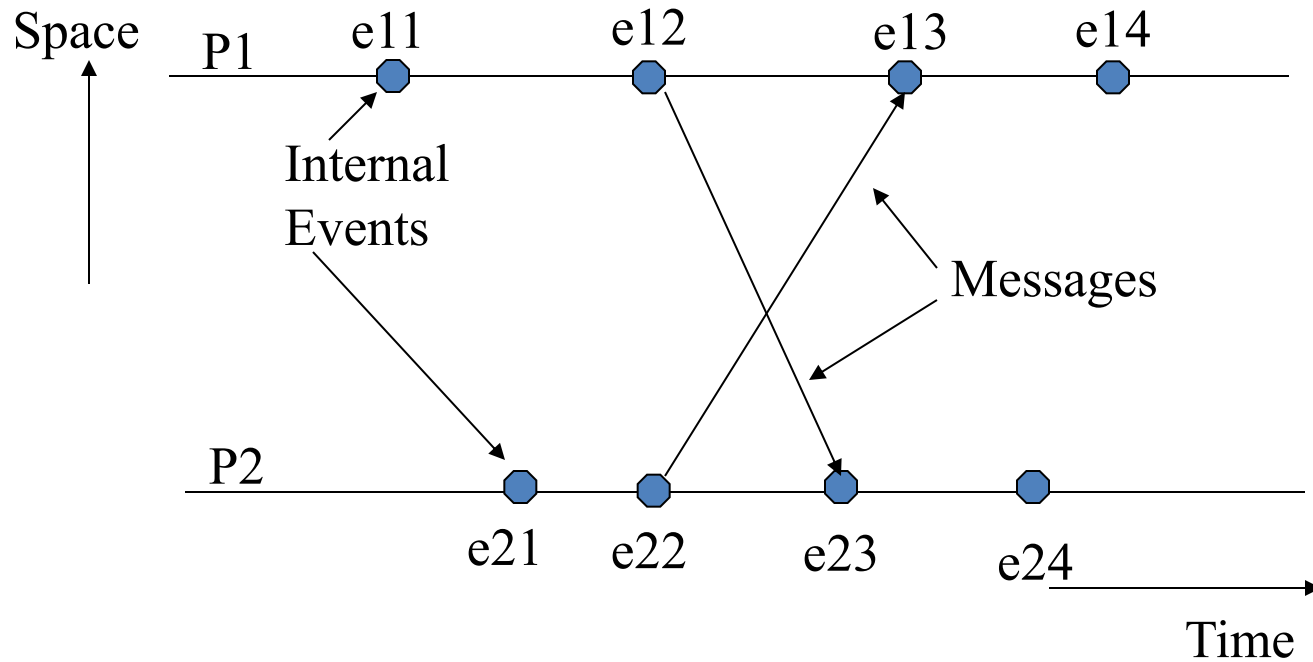
Theory aspects

- ☐ Logical Clocks – đồng hồ logic
- ☐ Causal Ordering – thứ tự nhân quả
- ☐ Global State Recording – Trạng thái toàn cục
- ☐ Termination Detection – Nhận dạng kết thúc

Lamport's Clock

- Happened before relation:
 - $a \rightarrow b$: Event a occurred before event b . Events in the same process.
 - $a \rightarrow b$: If a is the event of sending a message m in a process and b is the event of receipt of the same message m by another process.
 - $a \rightarrow b, b \rightarrow c$, then $a \rightarrow c$. “ \rightarrow ” is transitive.
- Causally Ordered Events
 - $a \rightarrow b$: Event a “causally” affects event b
- Concurrent Events
 - $a \parallel b$: if $a \nrightarrow b$ and $b \nrightarrow a$

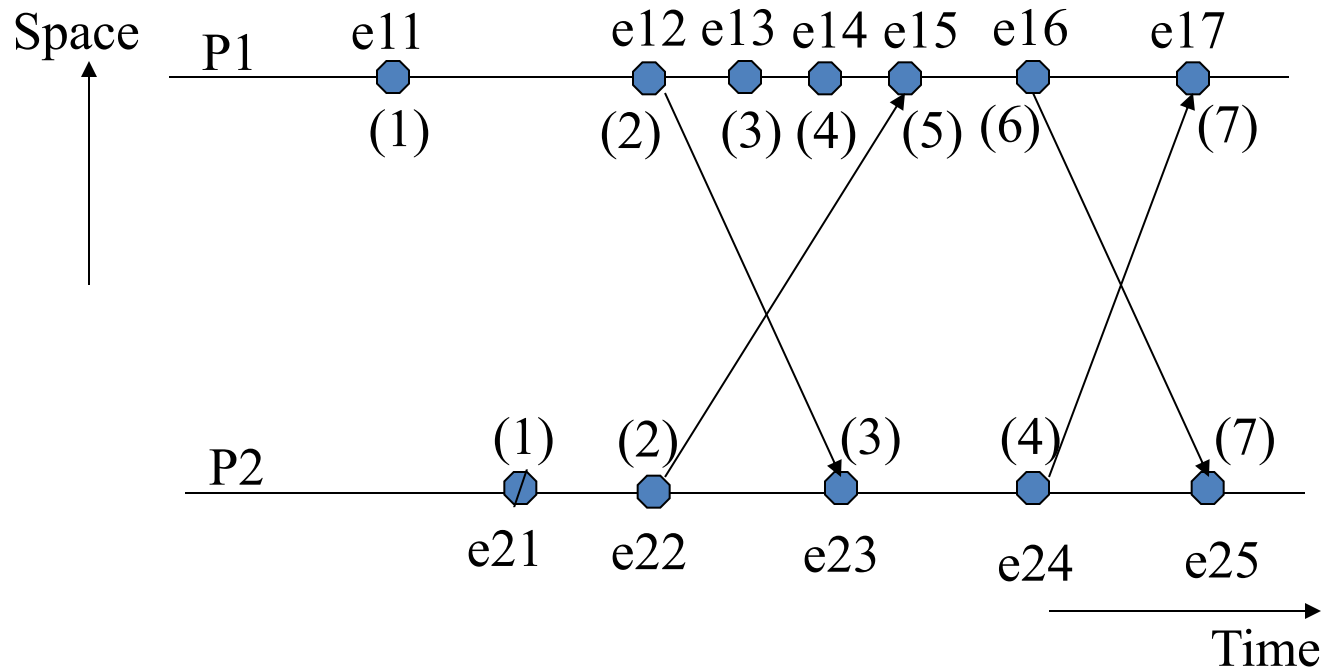
Lược đồ không-thời gian



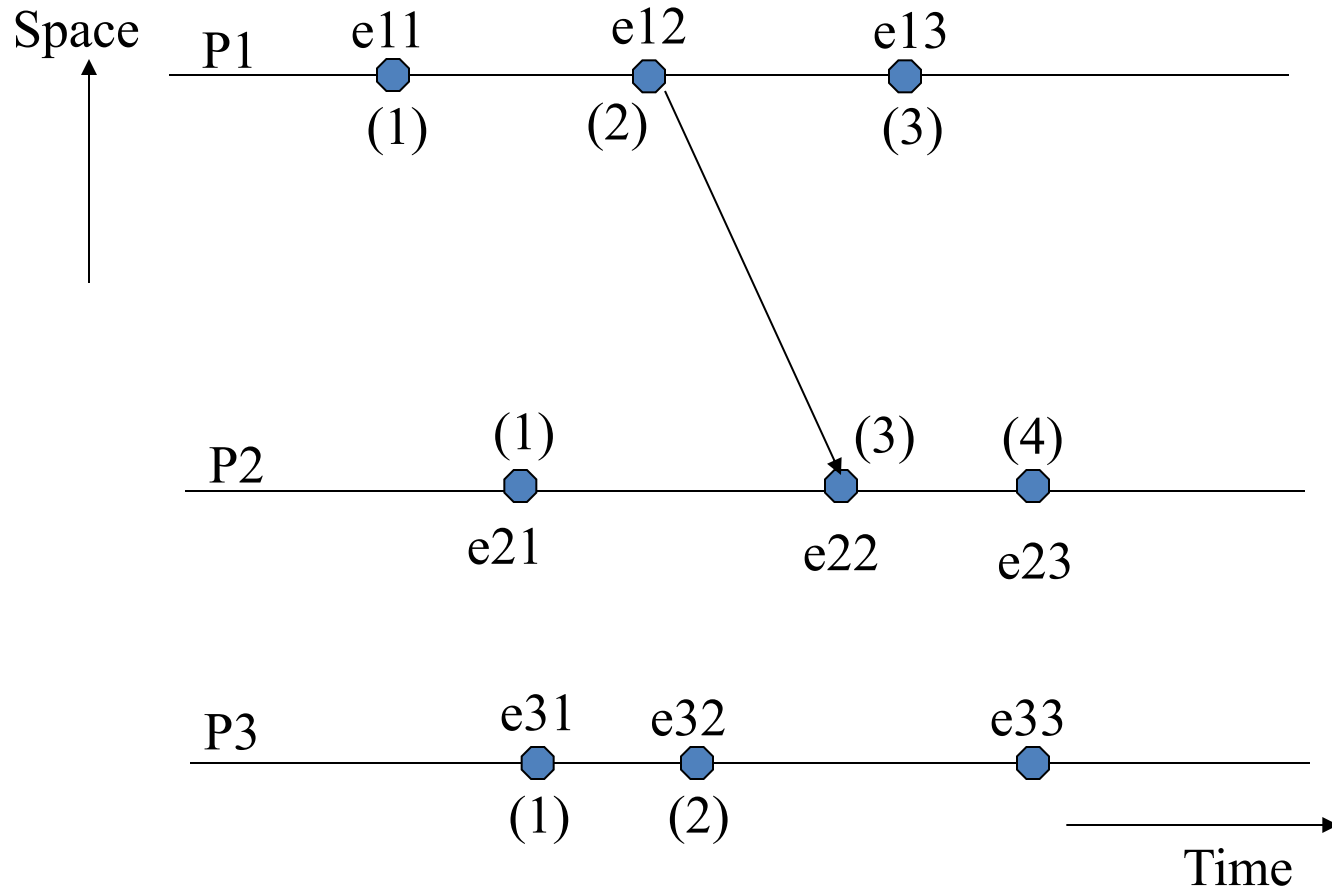
Logical Clocks – đồng hồ logic

- Qui luật:
 - C_i là giờ của Process P_i .
 - Nếu $a \rightarrow b$ trong process P_i , $C_i(a) < C_i(b)$
 - Nếu a : gửi msg m trong P_i ; b : nhận msg m trong P_j ; thì, $C_i(a) < C_j(b)$.
- Luật cập nhật đồng hồ:
 - R1: $C_i = C_i + d$ ($d > 0$); đồng hồ đc cập nhật giữa 2 sự kiện liên tiếp.
 - R2: $C_j = \max(C_j + d, t_m + d)$; ($d > 0$); Khi P_j nhận m với time stamp t_m (t_m gán bởi bên gửi P_i ; $t_m = C_i(a)$, sự kiện gửi m).
- Giá trị thường dùng cho d là 1

Lược đồ không-thời gian



Hạn chế của Lamport's Clock



$C(e11) < C(e32)$ nhưng không có quan hệ nhân quả.

Các thứ tự sự kiện dạng này không mô tả đc bằng Lamport's Clock.

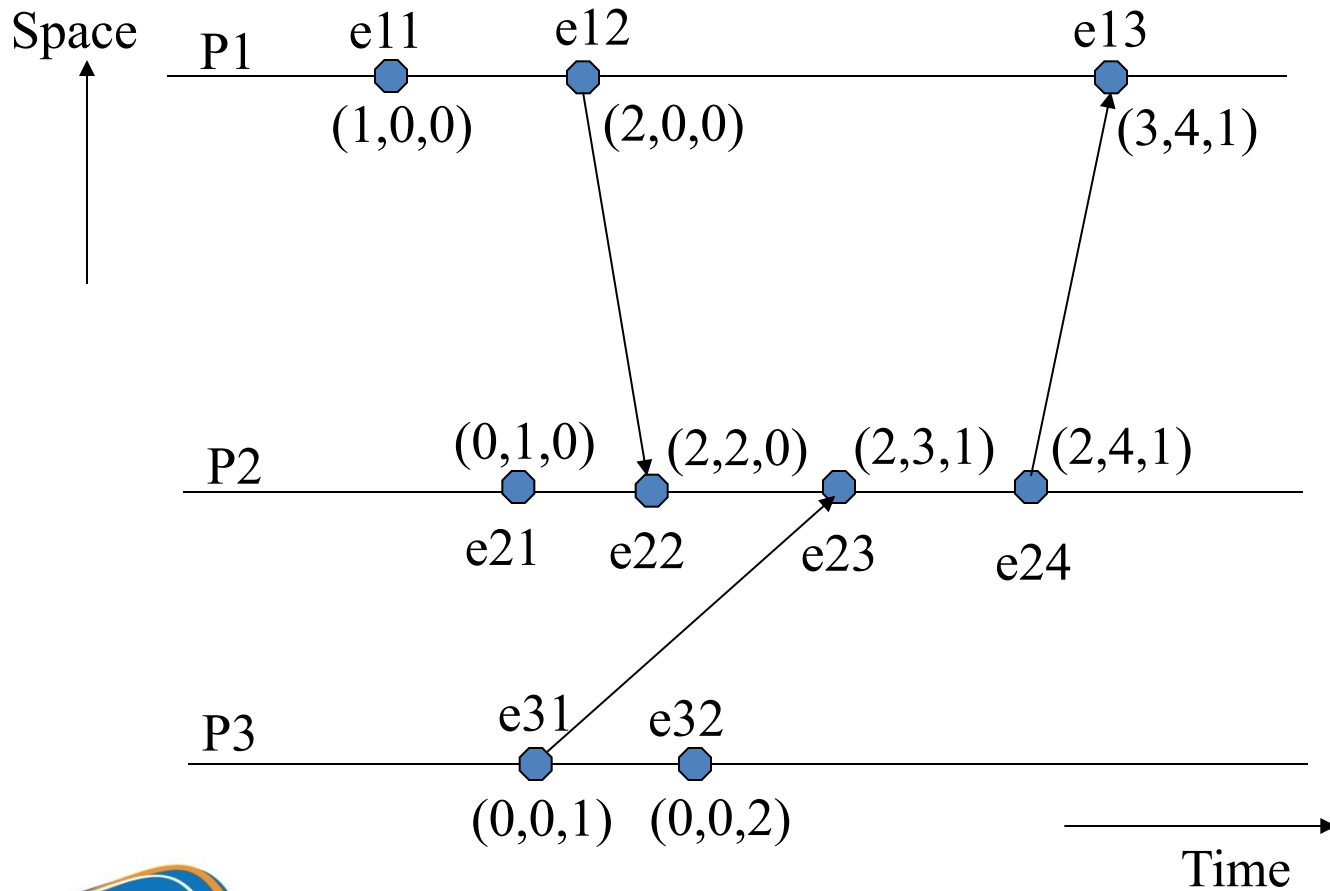
Vector Clocks –đồng hồ vector

- Lưu giữ mối qhe trao đổi giữa các process, để quá trình phục hồi (recovery).
- $C_i[1..n]$: là một “vector” clock lưu tại process P_i các thành phần của vector là giá trị clock (“assumed/best guess”) của các process khác.
- $C_i[j]$ ($j \neq i$) là giá trị dự đoán tốt nhất của P_i về P_j 's clock.
- Luật cập nhật Vector clock:
 - ▣ $C_i[i] = C_i[i] + d$, ($d > 0$); cho sự kiện liên tiếp của P_i
 - ▣ Với mọi k , $C_j[k] = \max(C_j[k], tm[k])$, khi m có time stamp tm gửi bởi P_i và được nhận tại P_j .

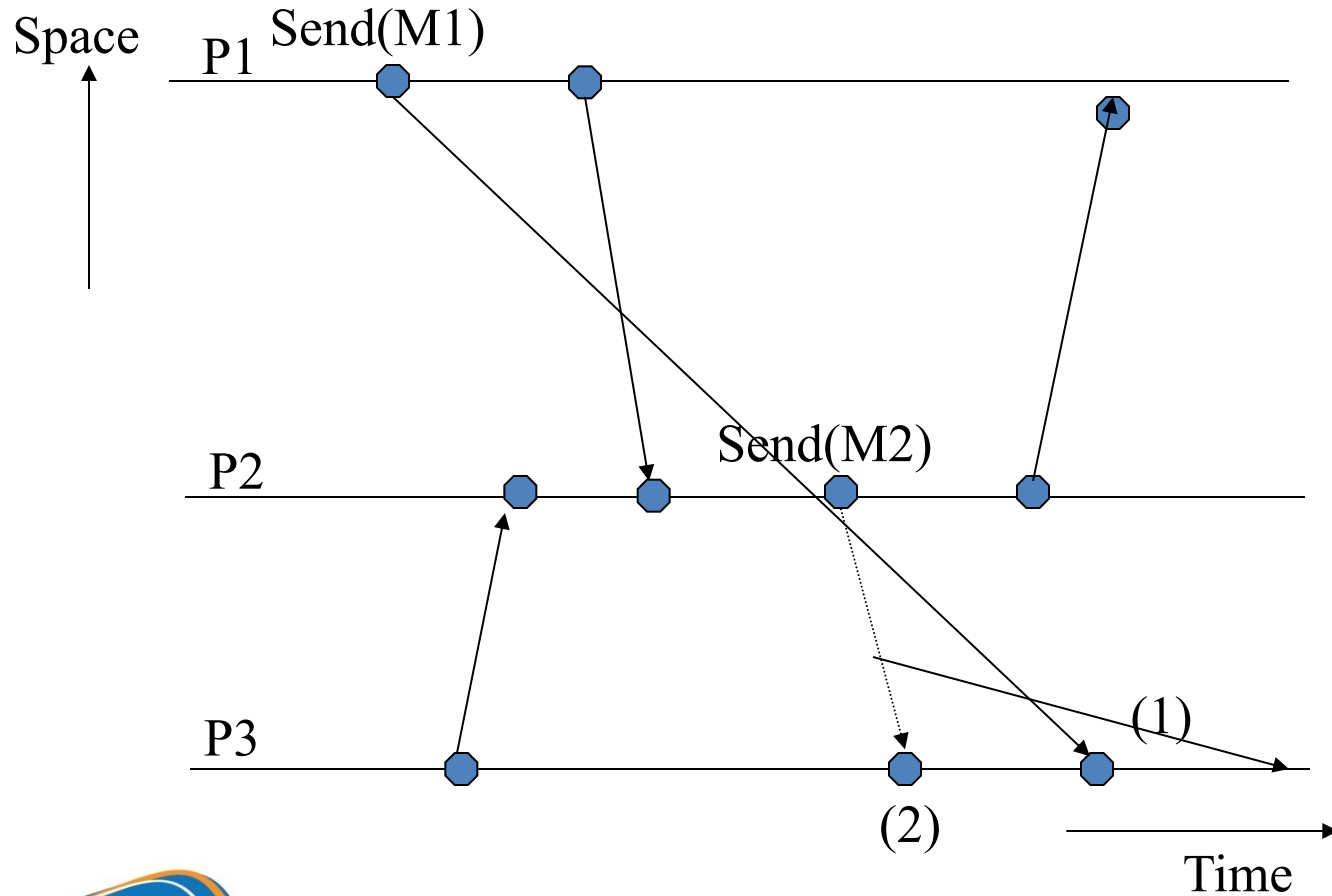
So sánh Vector Clocks

1. Equal: $t^a = t^b$ iff $\forall i, t^a[i] = t^b[i]$
2. Not Equal: $t^a \neq t^b$ iff $t^a[i] \neq t^b[i]$, for at least one i
3. Less than or equal: $t^a \leq t^b$ iff $t^a[i] \leq t^b[i]$, for all i
4. Less than : $t^a < t^b$ iff $t^a[i] \leq t^b[i]$ and $t^a[i] \neq t^b[i]$, for all i
5. Concurrent: $t^a \parallel t^b$ iff $t^a \not< t^b$ and $t^b \not< t^a$
6. Not less than or equal ...
7. Not less than ..

Vector Clock ...



Thứ tự nhân quả của Messages



Thứ tự Message ...

- ☐ Không nhất thiết phải duy trì một clock.
- ☐ Thứ tự của các message giữa các processes trong hệ phân tán.
- ☐ (v.d.) $\text{Send}(M1) \rightarrow \text{Send}(M2)$, M1 nên đc nhận trước M2.
- ☐ Điều này ko đc đảm bảo trong mạng máy tính. M1 có thể từ P1 đến P2 và M2 có thể từ P3 đến P4.
 indexStr
- ☐ Thứ tự Message:
 - ☐ Chuyển message cho ứng dụng khi mà message trước nó đã đc chuyển.
 - ☐ Ngược lại, lưu tạm trong buffer.

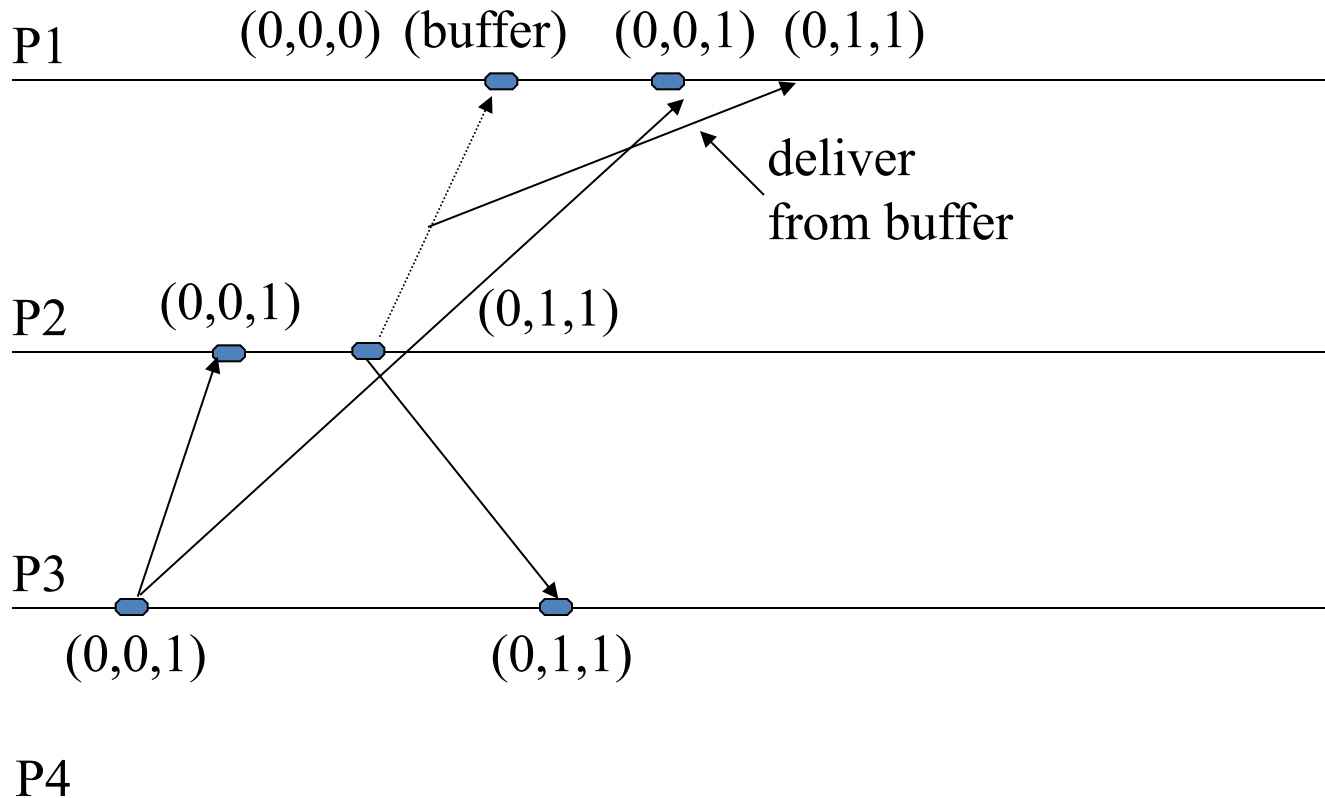
Thuật toán BSS

- ☐ BSS: Birman-Schiper-Stephenson Protocol
- ☐ Broadcast based: gửi msg tới tất cả processes.
- ☐ Chuyển message tới process khi mà message ngay trước nó đã đc chuyển.
- ☐ Ngược lại, buffer the message.
- ☐ Thực hiện bằng cách sử dụng vector gửi kèm message.

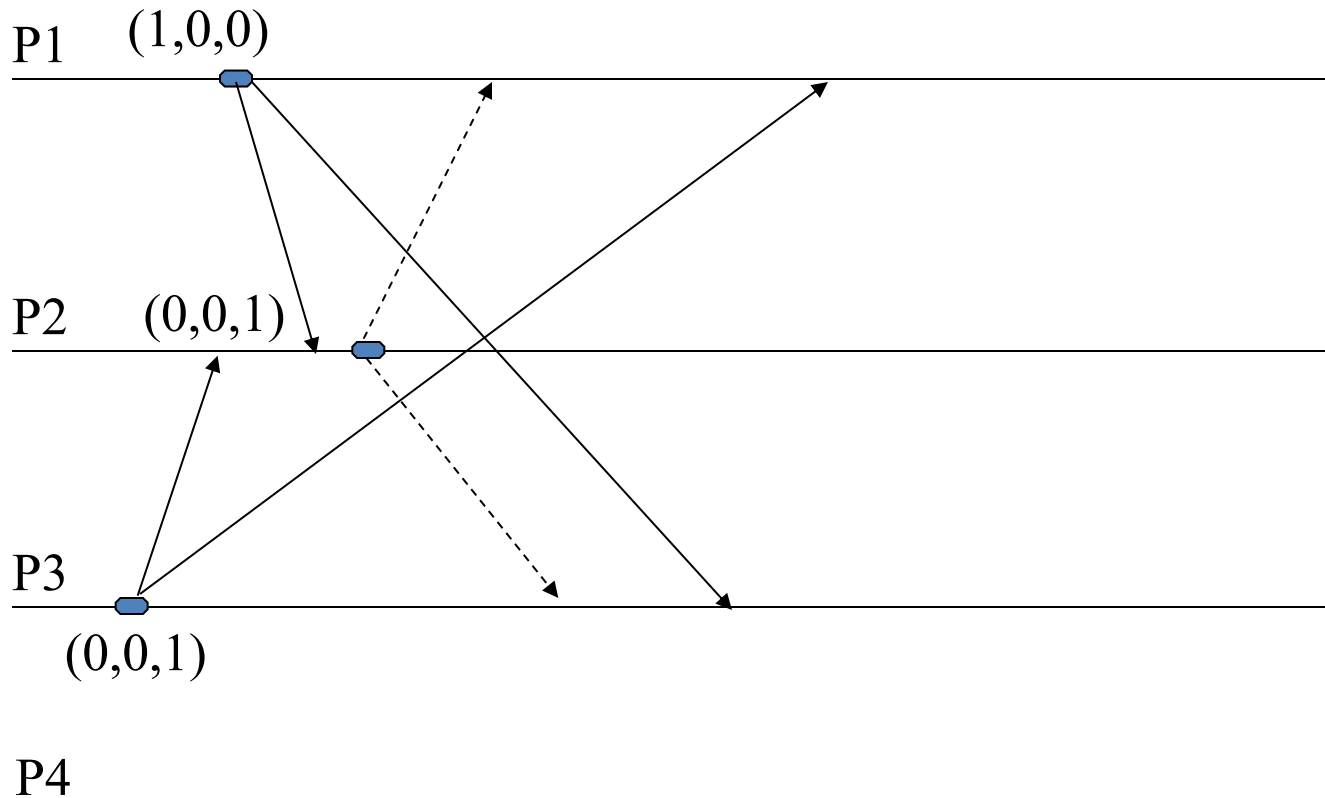
Thuật toán BSS ...

1. Process P_i tăng vector time $VT_{pi}[i]$, time stamps, và broadcasts message m . $VT_{pi}[i] - 1$ là số message trước m .
2. $P_j \neq P_i$ nhận m . m đc chuyển tới nếu:
 - a. $VT_{pj}[i] == VT_m[i] - 1$
 - b. $VT_{pj}[k] \geq VT_m[k]$ với mọi k in $\{1, 2, \dots, n\} - \{i\}$, n là tổng số processes.
 - c. Messages cùng lúc đc sắp xếp theo thời điểm nhận.
3. Khi m đc chuyển giao tại P_j , VT_{pj} đc cập nhật theo Rule 2 của vector clocks.
 - 2(a) : P_j nhận tất cả P_i 's messages trước m .
 - 2(b): P_j nhận tất cả messages mà P_i đã nhận đc trước khi gửi m

Thuật toán BSS ...



Thuật toán BSS ...



Thuật toán SES

- SES: Schiper-Eggli-Sandoz Algorithm. Không cần broadcast message.
- Mỗi process lưu 1 vector V_P kích thước $N - 1$, N số lượng processes.
- Mỗi phần tử của V_P chứa (P', t) : P' là id của process đích và t là vector timestamp.
- tm : thời điểm (logic) gửi m
- tP_i : thời điểm (logic) hiện tại tại p_i
- Ban đầu, V_P rỗng.

Thuật toán SES

□ Gửi Message:

- Gửi message M, time stamped t_m , cùng với V_{P1} đến $P2$.
- Thêm $(P2, t_m)$ vào V_{P1} . Ghi chồng lên $(P2, t)$, nếu có.
- $(P2, t_m)$ không đc gửi. Các message có $(P2, t_m)$ trong V_{P1} chỉ đc chuyển đến $P2$ khi mà $t_m < t_{P2}$.

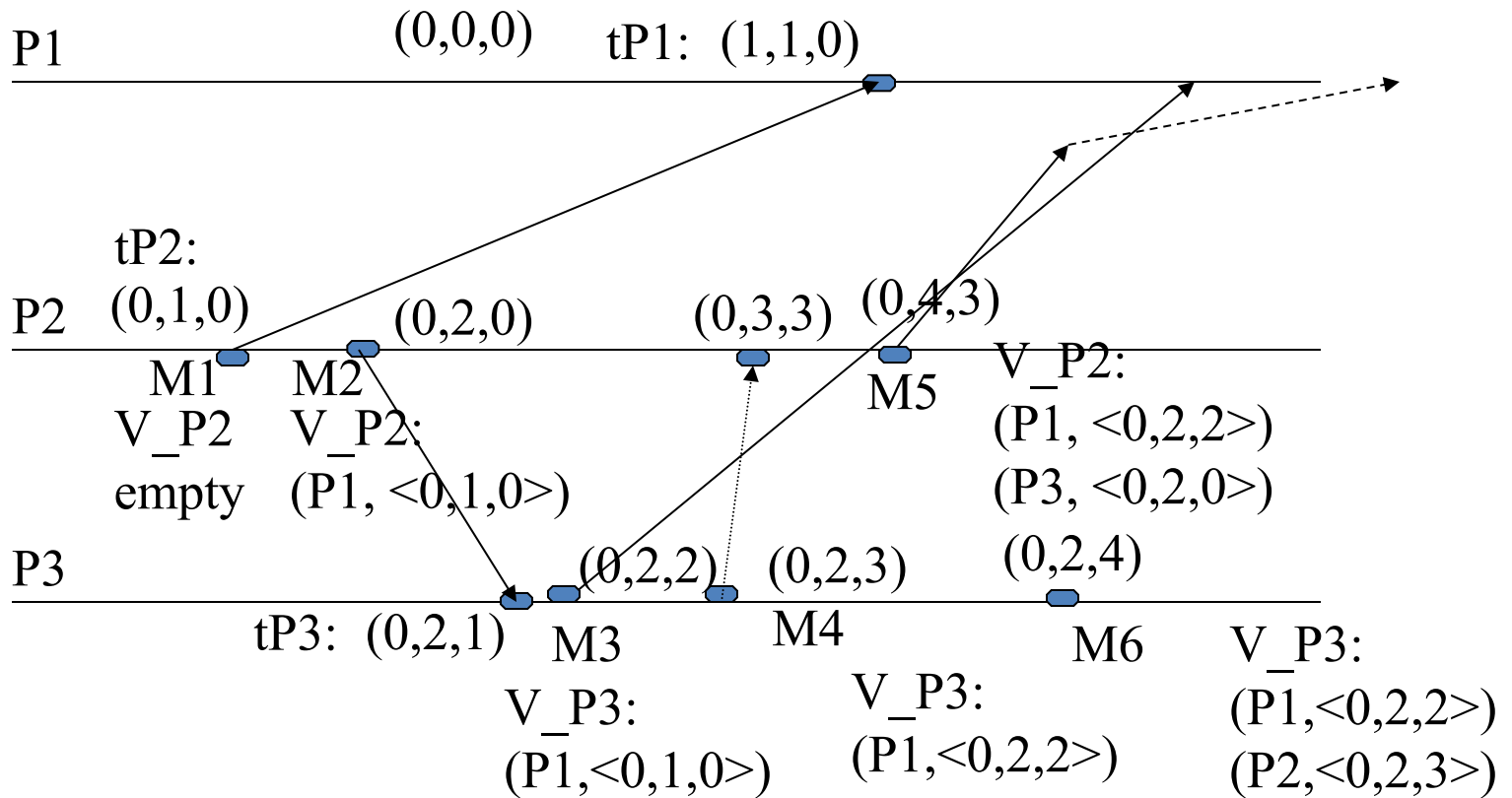
□ Chuyển giao message

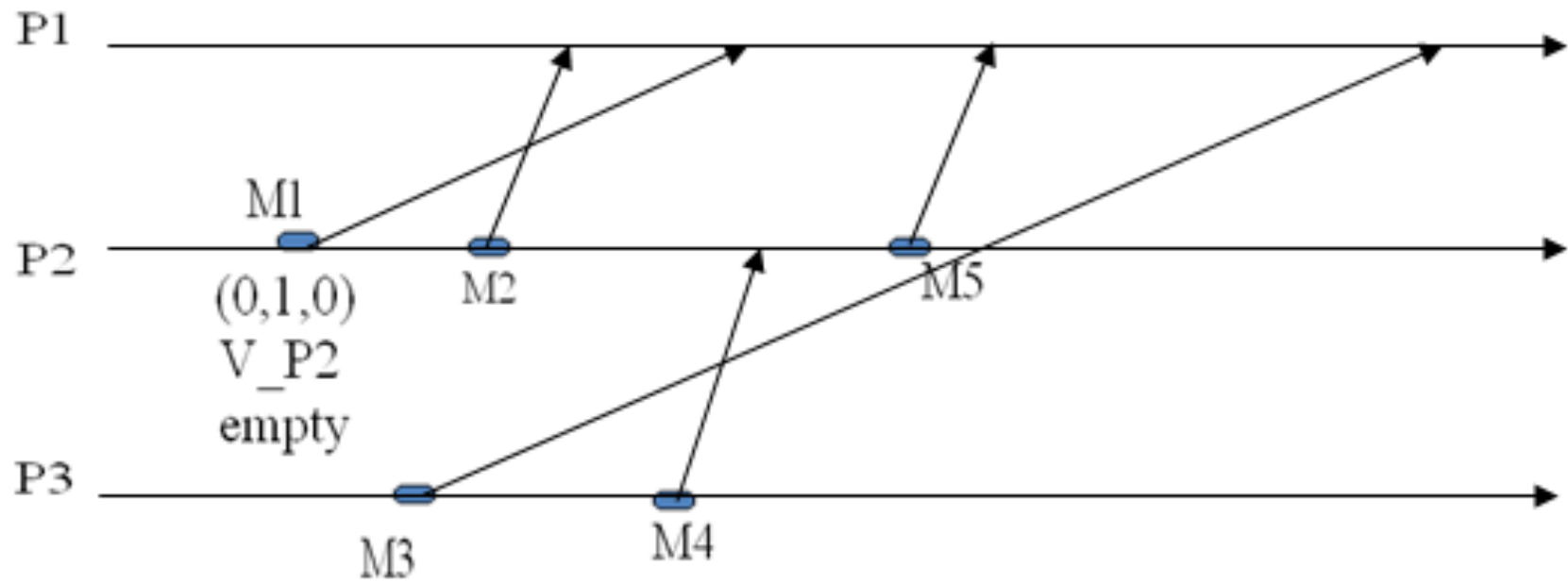
- If V_M (in the message) không chứa $(P2, t)$, thì có thể chuyển msg này.
- /* $(P2, t)$ exists */ If $t_m > t_{P2}$, buffer the message. (Don't deliver).
- else ($t_m \leq t_{P2}$) deliver it

Thuật toán SES ...

- Điều kiện $t \geq tP2$ nói lên điều gì?
 - t is message vector time stamp.
 - $t > tP2 \rightarrow$ For all j , $t[j] > tP2[j]$
 - Có tồn tại sự kiện trong process khác mà $P2$'s chưa cập nhật. Vì vậy $P2$ quyết định buffer the message.
- Khi $t < tP2$, message đc chuyển & $tP2$ được cập nhật với thông tin trong V_{P2} (sau phép trộn).

Ví dụ SES Buffering





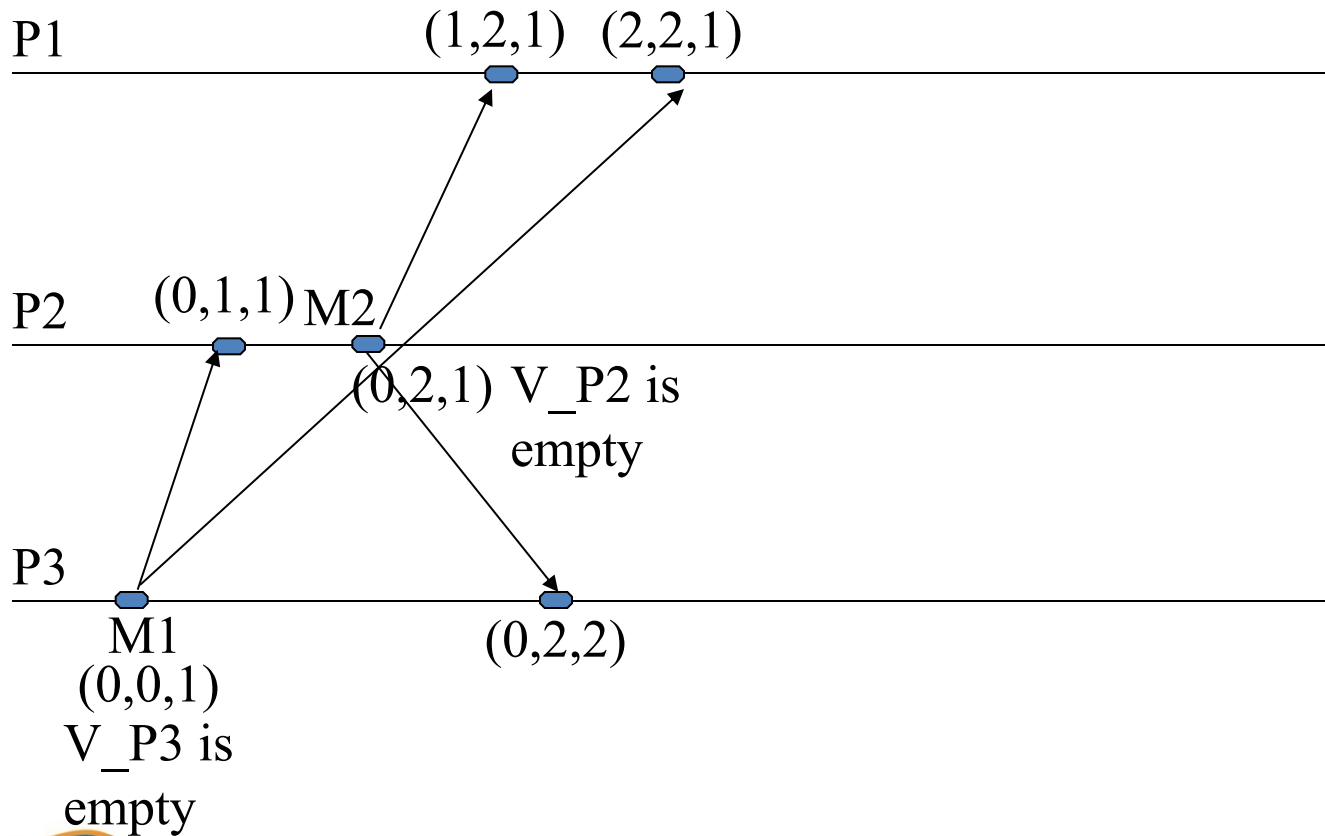
Ví dụ SES Buffering...

- M1 from P2 to P1: $M1 + Tm (= \langle 0, 1, 0 \rangle)$ + Empty V_{P2}
- M2 from P2 to P3: $M2 + Tm (\langle 0, 2, 0 \rangle)$ + $(P1, \langle 0, 1, 0 \rangle)$
- M3 from P3 to P1: $M3 + \langle 0, 2, 2 \rangle$ + $(P1, \langle 0, 1, 0 \rangle)$
- M3 gets buffered because:
 - ▣ $Tp1$ is $\langle 0, 0, 0 \rangle$, t in $(P1, t)$ is $\langle 0, 1, 0 \rangle$ & so $Tp1 < t$
- When M1 is received by P1:
 - ▣ $Tp1$ becomes $\langle 1, 1, 0 \rangle$, by rules 1 and 2 of vector clock.
- After updating $Tp1$, P1 checks buffered M3.
 - ▣ Now, $Tp1 > t$ [in $(P1, \langle 0, 1, 0 \rangle)$].
So M3 is delivered.

Thuật toán SES ...

- Khi chuyển giao message:
 - Merge V_M (in message) with V_P2 as follows.
 - If (P,t) is not there in V_P2 , merge.
 - If (P,t) is present in V_P2 , t is updated with $\max(t[i] \text{ in } V_m, t[i] \text{ in } V_P2)$. {Component-wise maximum}.
 - Message không chuyển cho tới khi t trong V_M nhỏ hơn t trong V_P2
 - Cập nhật site $P2$'s local, logical clock.
 - Kiểm tra các message đang được buffered, logical clock update.

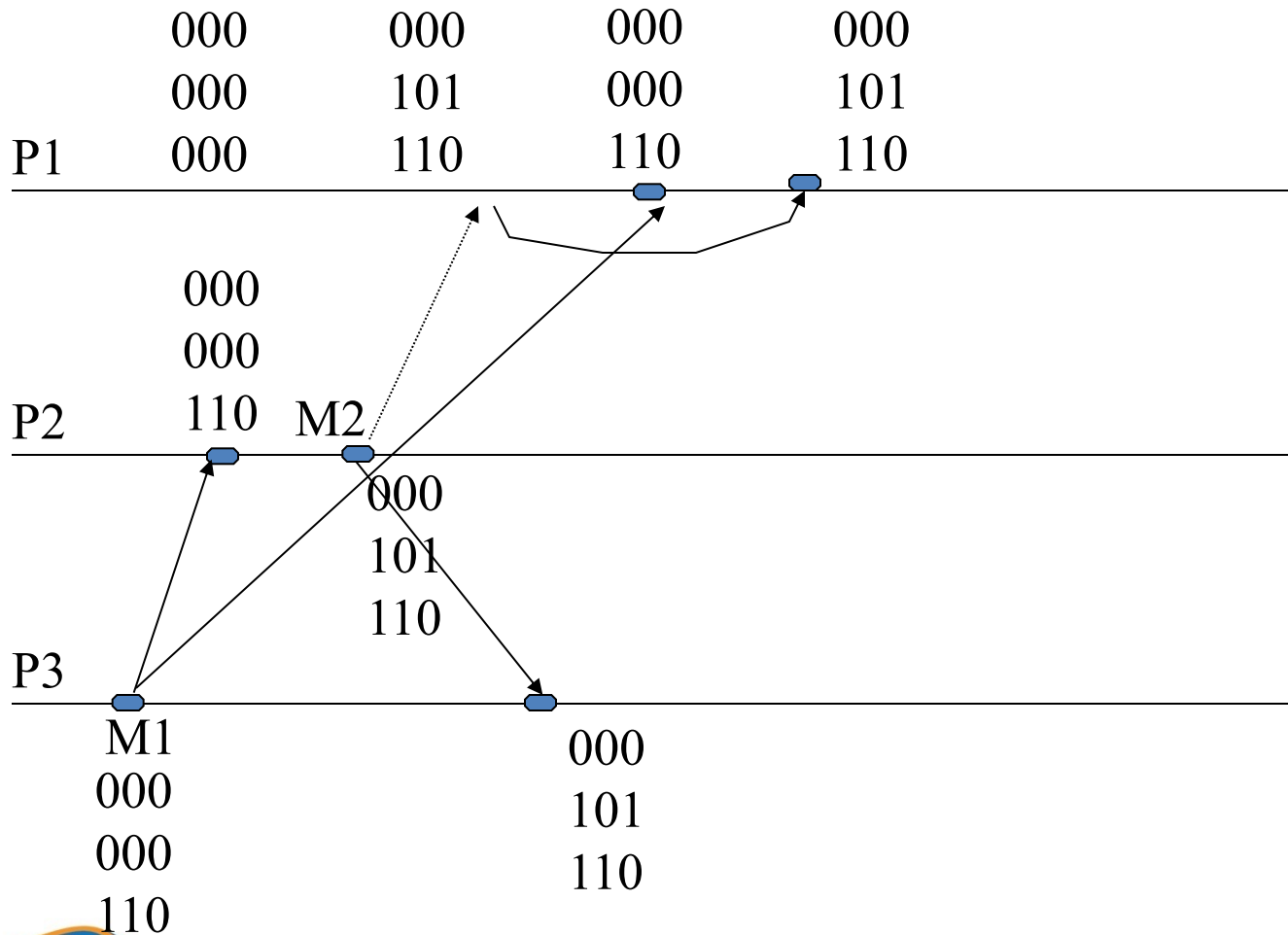
SES Algorithm ...



Multicasts

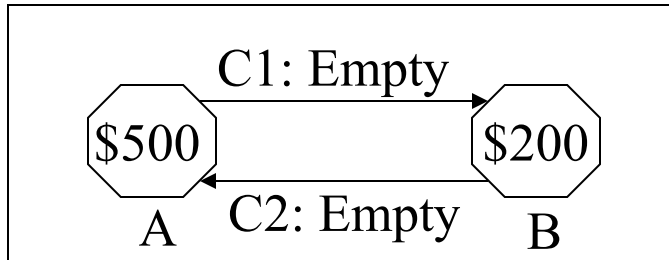
- Mỗi node lưu matran M ($n \times n$), n số lượng processes.
- Node i multicasts to j and k : tăng $M_{i,j}$ và $M_{i,k}$. M đc gửi cùng với message.
- Khi node j nhận m từ i , nó có thể đc chuyển nếu và chỉ nếu:
 - $M_{j,i} = M_{m,i} - 1$
 - $M_{j,k} \geq M_{m,k}$ for all $k \neq i$.
- Còn lại buffer the message
- Khi chuyển giao message: $M_j[x,y] = \max(M_j[x,y], M_m[x,y])$

Ví dụ Multicasts

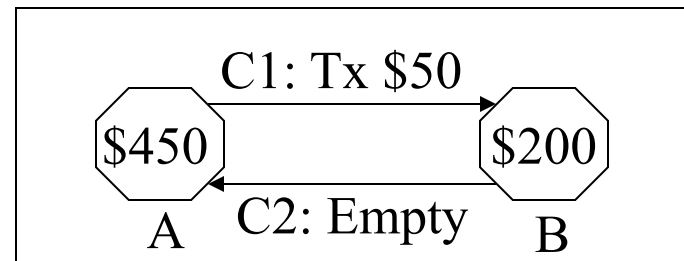


Global State – trạng thái toàn cục

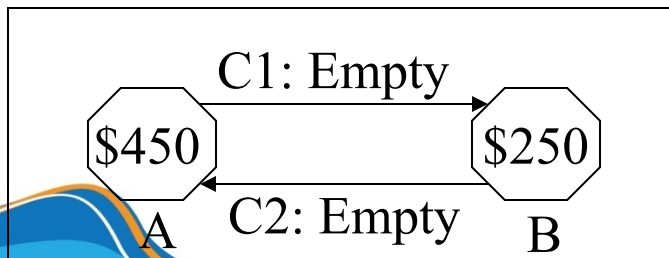
Global State 1



Global State 2



Global State 3



Lưu Global State...

- (v.d.,) Global state of A đc lưu ở (1) không phải ở (2).
 - ▣ Trạng thái của B, C1, và C2 đc lưu trong (2)
 - ▣ Khoản tiền \$50 sẽ xuất hiện trong global state
 - ▣ Lý do: A's state được lưu trước khi gửi message và C1's state *lưu sau khi gửi message*.
- Inconsistent global state if $n < n'$, trong đó
 - ▣ n số message đc gửi bởi A và được lưu trong A's state
 - ▣ n' số message gửi bởi A đang trong channel, và được lưu trong channel's state.
- Consistent global state: $n = n'$

Lưu Global State...

- Tương tự, để nhất quán $m = n'$
 - ▣ n' : no. of messages trong channel gửi đến B
 - ▣ m : no. of messages mà B được nhận từ channel này.
- Nên $n' < m$, không tồn tại hệ thống mà số lượng msg gửi ít hơn số lượng msg nhận
- Vậy, $n' \geq m$
- Consistent global state phải thỏa điều kiện trên.
- Consistent global state:
 - ▣ Channel state: thứ tự message đc gửi trước khi lưu sender's state, trừ đi các message đã đc nhận và đã lưu tại receiver's state.
 - ▣ Chỉ có transit messages được lưu trong channel state.

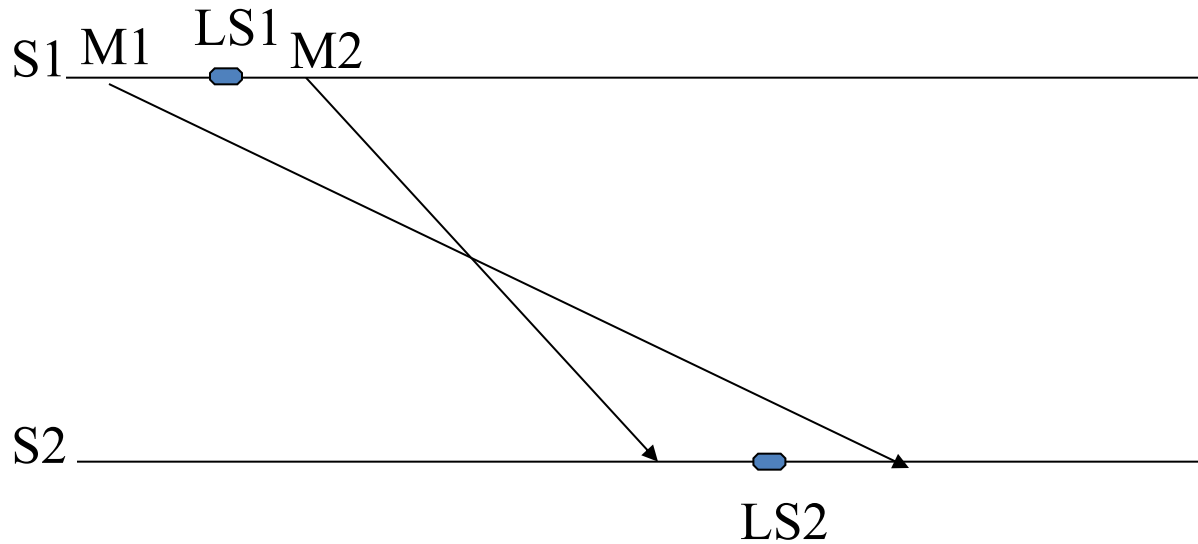
LUU Global State

- ☐ $\text{Send}(M_{ij})$: message M sent from S_i to S_j
- ☐ $\text{rec}(M_{ij})$: message M received by S_j , from S_i
- ☐ $\text{time}(x)$: Time of event x
- ☐ LS_i : local state at S_i
- ☐ $\text{send}(M_{ij})$ is in LS_i iff (if and only if) $\text{time}(\text{send}(M_{ij})) < \text{time}(\text{LS}_i)$
- ☐ $\text{rec}(M_{ij})$ is in LS_j iff $\text{time}(\text{rec}(M_{ij})) < \text{time}(\text{LS}_j)$
- ☐ $\text{transit}(\text{LS}_i, \text{LS}_j)$: set of messages sent/recorded at LS_i and NOT received/recorded at LS_j

Local Global State ...

- $\text{inconsistent}(\text{LS}_i, \text{LS}_j)$: set of messages NOT sent/recorded at LS_i and received/recorded at LS_j
- Global State, $\text{GS}: \{\text{LS}_1, \text{LS}_2, \dots, \text{LS}_n\}$
- Consistent Global State, $\text{GS} = \{\text{LS}_1, \dots, \text{LS}_n\}$ AND for all i, j in n , $\text{inconsistent}(\text{LS}_i, \text{LS}_j)$ is null.
- Transitive global state, $\text{GS} = \{\text{LS}_1, \dots, \text{LS}_n\}$ AND for all i, j in n , $\text{transit}(\text{LS}_i, \text{LS}_j)$ is null.

Lưu Global State ..

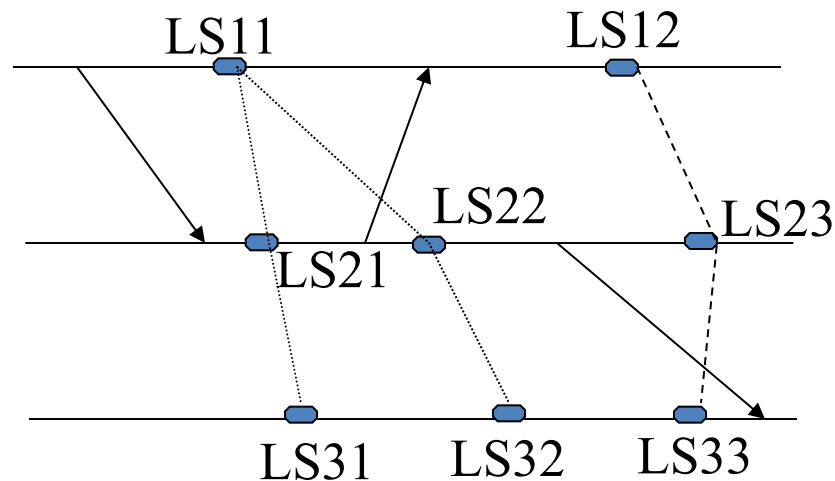


M1: transit

M2: inconsistent

Lưu Global State...

- Strongly consistent global state: consistent và transitless, i.e., tất cả các sự kiện gửi và sk nhận đã được lưu trong tất cả LSi.

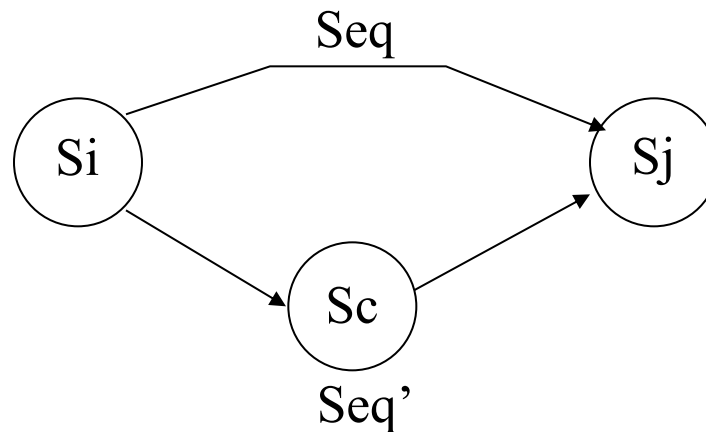


Thuật toán Chandy-Lamport

- Thuật toán phân tán để lưu consistent global state. Kênh truyền giả thiết là FIFO.
- Sử dụng 1 *marker* để khởi tạo thuật toán. Marker sort of dummy message, with no effect on the functions of processes.
- Gửi Marker tại P:
 - ▣ P lưu trạng thái của nó.
 - ▣ Với mỗi outgoing channel C, P gửi marker qua C trước khi P gửi các message khác qua C.
- Nhận Marker tại Q:
 - ▣ Nếu Q CHƯA lưu trạng thái: (a). Lưu trạng thái của C là rỗng. (b) SEND marker (sử dụng luật bên trên).
 - ▣ Ngược lại (Q đã lưu trạng thái): lưu trạng thái của C là chuỗi message đã nhận từ C, trong khoảng thời gian Q's state được lưu và trước khi Q nhận được marker.
- FIFO channel + markers làm thỏa tính nhất quán.

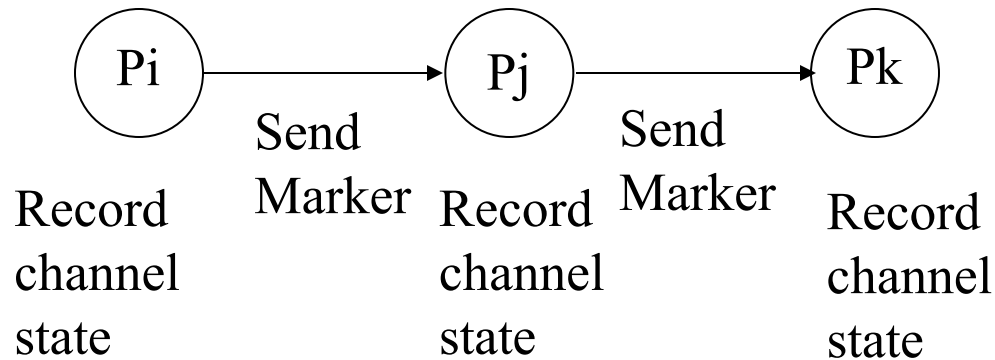
Thuật toán Chandy-Lamport

- Các process có thể khởi tạo marker, giá trị marker duy nhất: $\langle \text{process id, sequence number} \rangle$.
- Các processes có thể khởi tạo quá trình lưu trạng thái hệ thống cùng lúc.
- One possible way to collect global state: all processes send the recorded state information to the initiator of marker. Initiator process can sum up the global state.



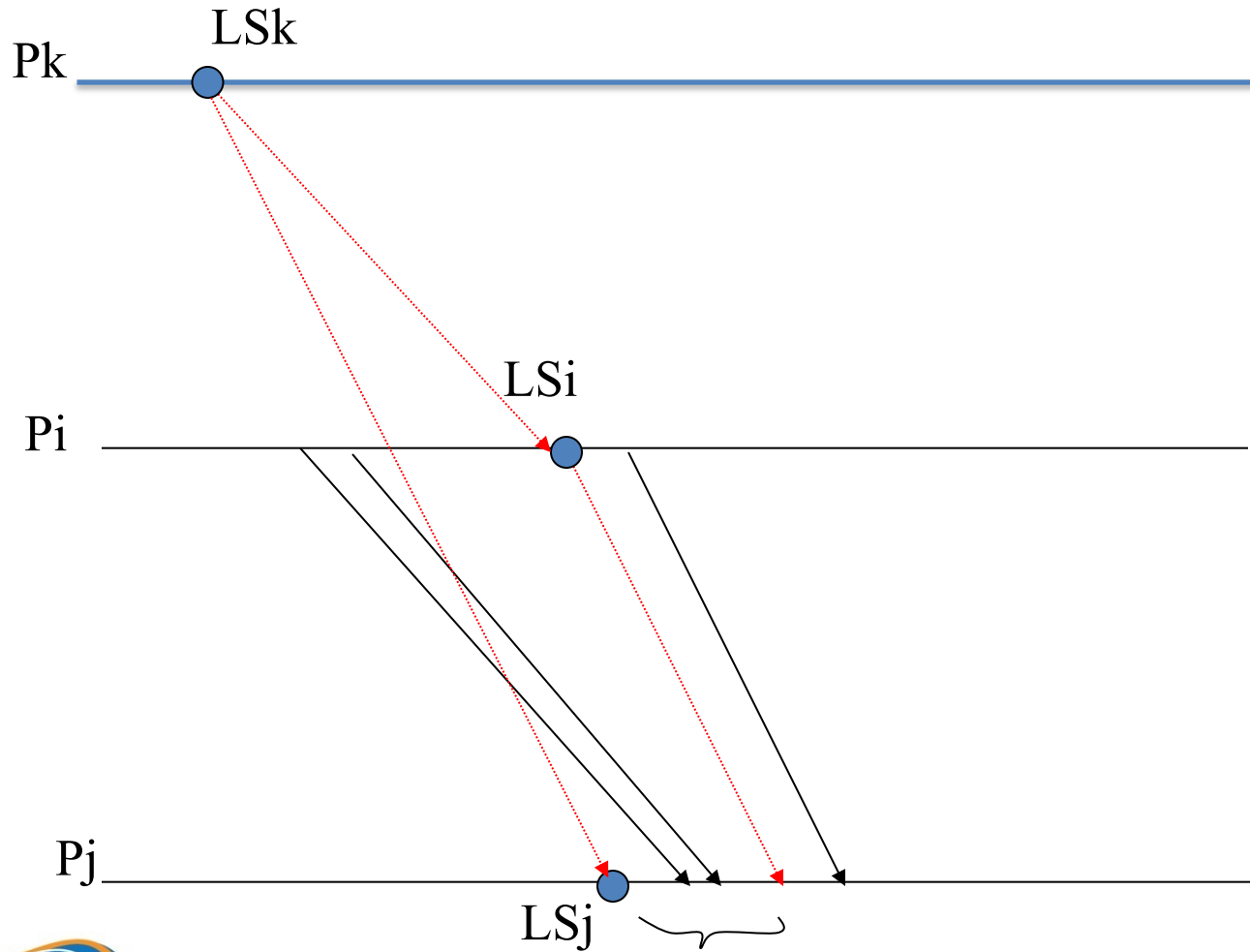
Thuật toán Chandy-Lamport ...

□ Example:



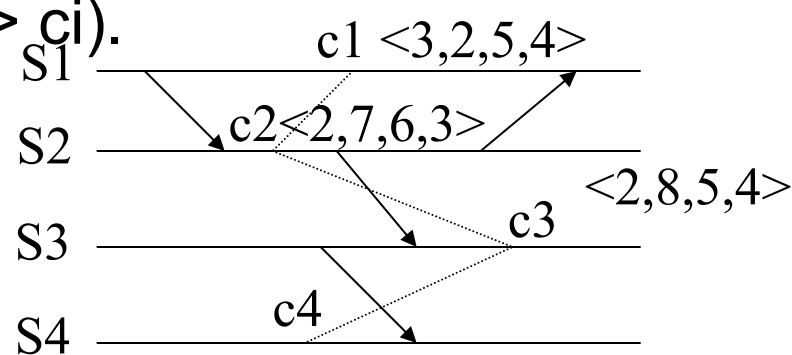
Channel state example: M1 sent to P_x at t_1 , M2 sent to P_y at t_2 ,

Chandy-Lamport Algorithm ...

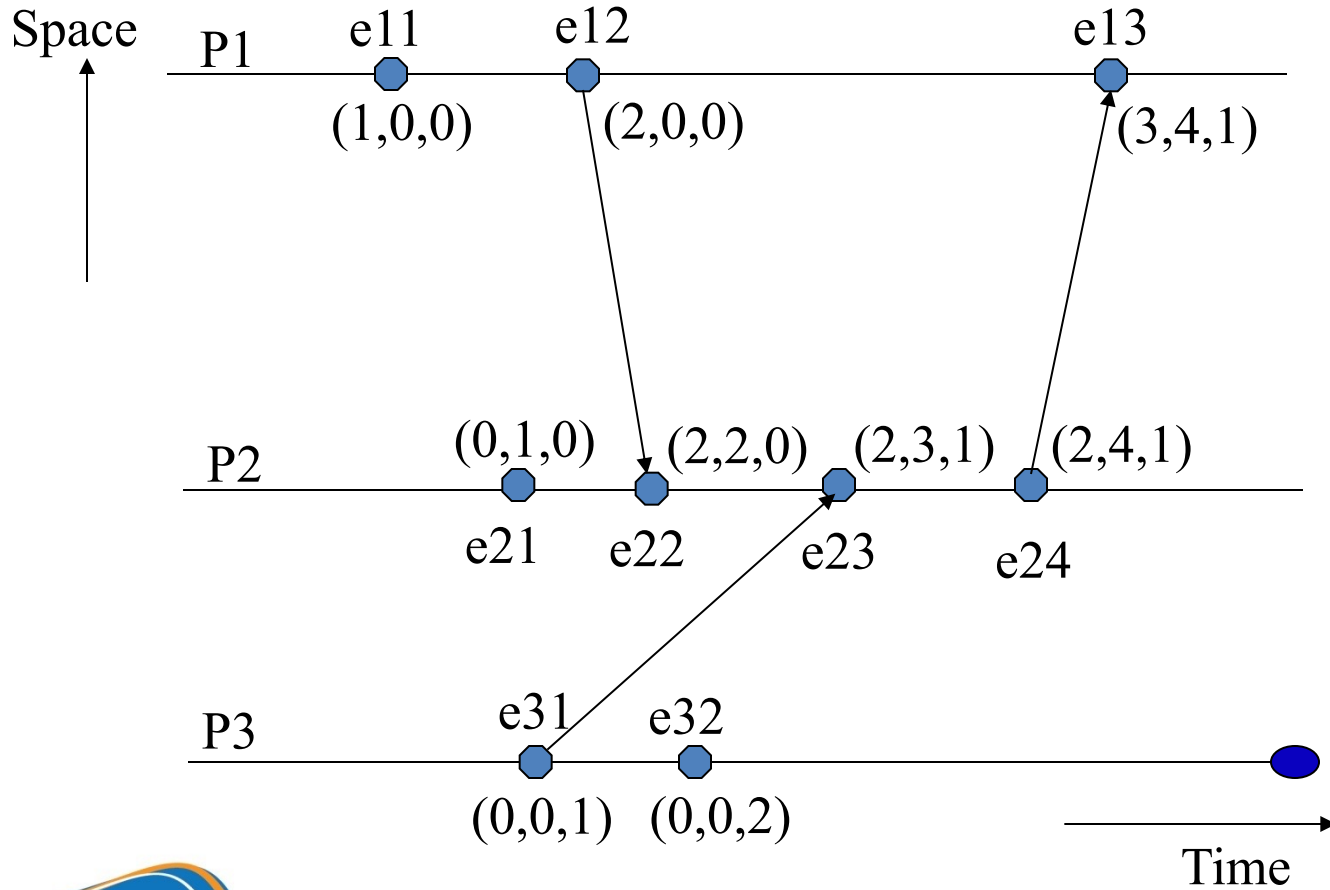


Cuts – lát cắt

- Cuts: biểu diễn bằng hình của một global state.
- Cut $C = \{c1, c2, \dots, cn\}$; c_i : sự kiện cut tại S_i .
- Consistent Cut: nếu mỗi message đc nhận tại S_i xảy ra trước cut event, thì phải đc gửi trước cut event tại bên Gửi.
- Định lý: Một cut là một consistent cut iff không tồn tại hai cut events có quan hệ nhân quả, i.e., $!(c_i \rightarrow c_j)$ và $!(c_j \rightarrow c_i)$.



Vector Clock ...



Thời điểm của một Cut

- $C = \{c1, c2, \dots, cn\}$ with vector time stamp VT_{ci} .
Vector time of the cut, $VT_c = \sup(VT_{c1}, VT_{c2}, \dots, VT_{cn})$.
- \sup is a component-wise maximum, i.e., $VT_{ci} = \max(VT_{c1}[i], VT_{c2}[i], \dots, VT_{cn}[i])$.
- Định lý: a cut is consistent iff $VT_c = (VT_{c1}[1], VT_{c2}[2], \dots, VT_{cn}[n])$.

Termination Detection

- ☐ Termination: completion of the sequence of algorithm. (e.g., leader election, deadlock detection, deadlock resolution).
- ☐ Use a *controlling agent* or a *monitor process*.
- ☐ Initially, all processes are idle. Weight of controlling agent is 1 (0 for others).
- ☐ Start of computation: message from controller to a process. Weight: split into half (0.5 each).
- ☐ Repeat this: any time a process send a computation message to another process, split the weights between the two processes (e.g., 0.25 each for the third time).
- ☐ End of computation: process sends its weight to the controller. Add this weight to that of controller's. (Sending process's weight becomes 0).
- ☐ *Rule*: Sum of W always 1.
- ☐ *Termination*: When weight of controller becomes 1 again.

Huang's Algorithm

- B(DW): computation message, DW is the weight.
- C(DW): control/end of computation message;
- Rule 1: Before sending B, compute $W1, W2$ (such that $W1 + W2$ is W of the process). Send B($W2$) to P_i , $W = W1$.
- Rule 2: Receiving B(DW) $\rightarrow W = W + DW$, process becomes active.
- Rule 3: Active to Idle \rightarrow send C(DW), $W = 0$.
- Rule 4: Receiving C(DW) by controlling agent $\rightarrow W = W + DW$, If $W == 1$, computation has terminated.

Huang's Algorithm

