

Time



Global States



Time & Global State

- Introduction (Time in Distributed Sys)
- Physical Clocks**
- Why to synchronize multiple clocks in DS?
- Physical Clock Synchronization
- NTP (Network Time Protocol)
- Logical Clocks**
- Distributed Debugging**
- Global States**
- Logical Time & Global States**



Introduction (1)

- ❑ Time is to order events.
- ❑ DS = N processes & Communication by messages & No Shared memory
- ❑ Process P_i State $S_i =$ (its variable values + its objects in OS)
- ❑ Events in DS = **send(m)**, **receive(m)**, **internal event**.
- ❑ No Global Time (problem in DS)

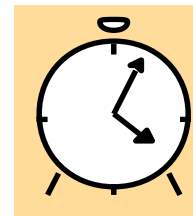
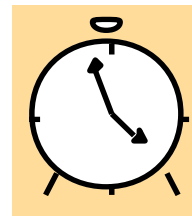
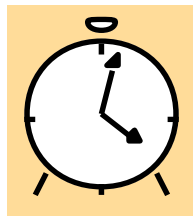
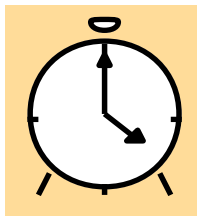


Introduction (2)

- One process (one physical clock) even multi-threaded →
 - single, total order for events : $e \longrightarrow_i e'$
- History $(p_i) = h_i = \langle e_i^0, e_i^1, e_i^2, \dots \rangle$
- No Global Time (problem in DS) = we can not accurately take the event times. So, we do not know the order of events.
- Two different observers may not agree on the interval between 2 events
- Two different observers may see the events in reversed order (Impossible if there is **causality** relation between the 2 events)

Physical Clocks

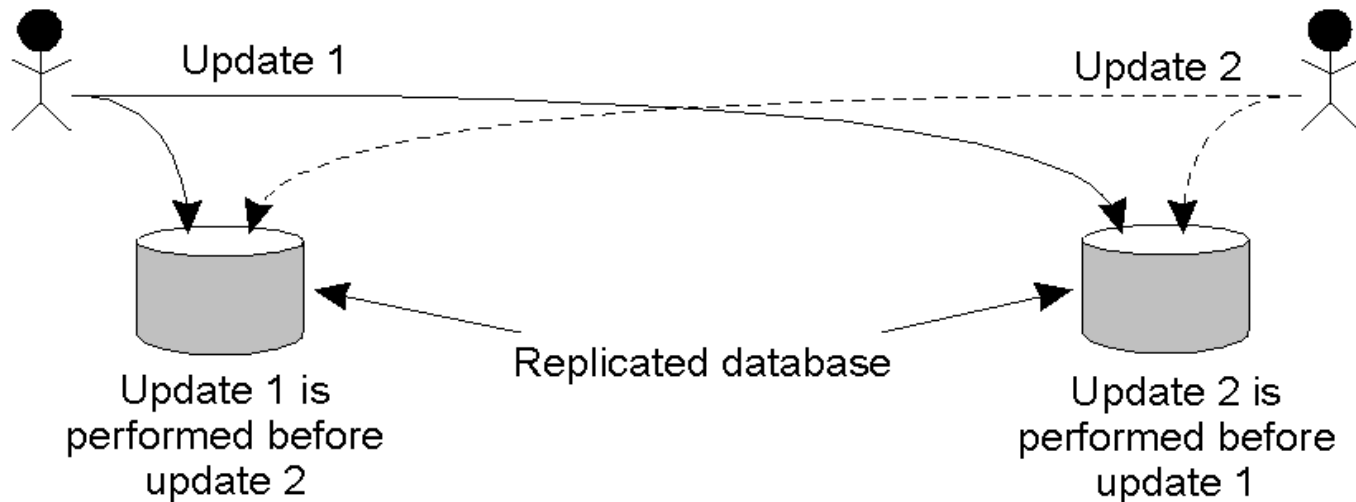
- In DS, each computer has its own clock.
 - Clock = A high-frequency **oscillator** (quartz crystal kept under tension) , **counter**, and **holding register**
- Clock **drift** = crystals run at different frequencies.
- Clock **skew** = the difference between two reads for two clocks at the same time.



Network

Why to synchronize multiple clocks in DS?

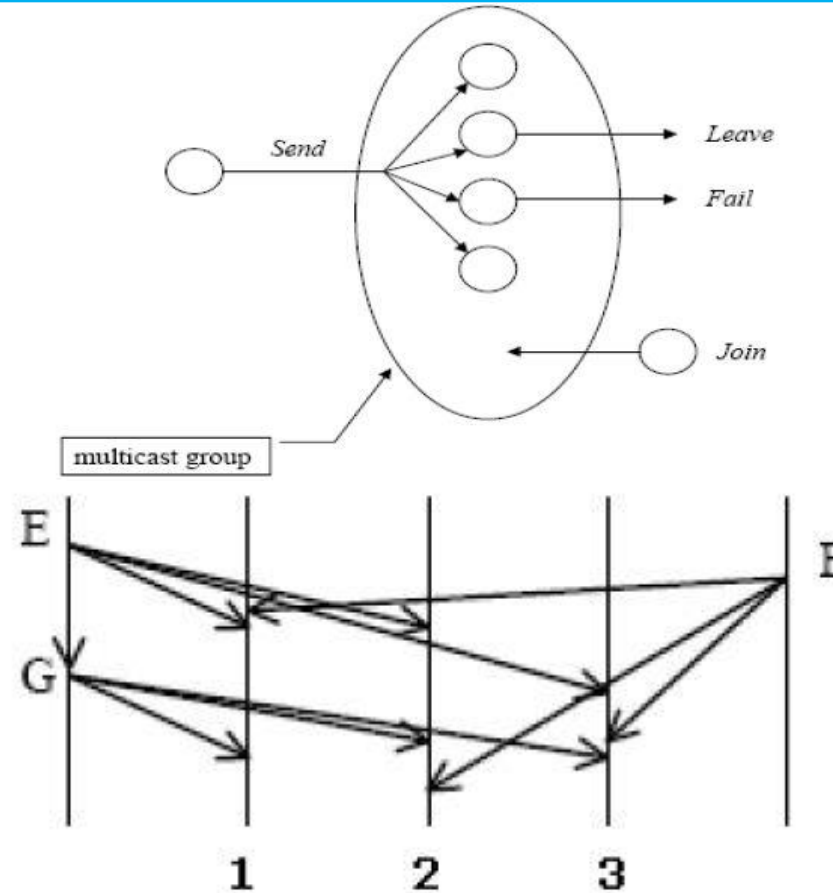
(Ex : Replicated Bank System)



- Balance = \$1000
- Update 1 = add(\$100)
- Update 2 = add(1%)
- balance 1 = 1110\$ && balance 2 = 1111\$

Why to synchronize multiple clocks in DS?

(Ex : Group Communication)

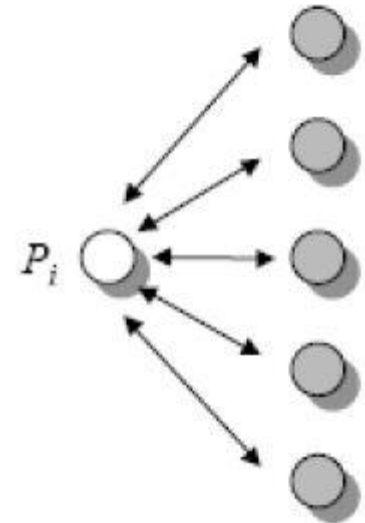


DS?

(Ex : Distributed Critical Section Control)

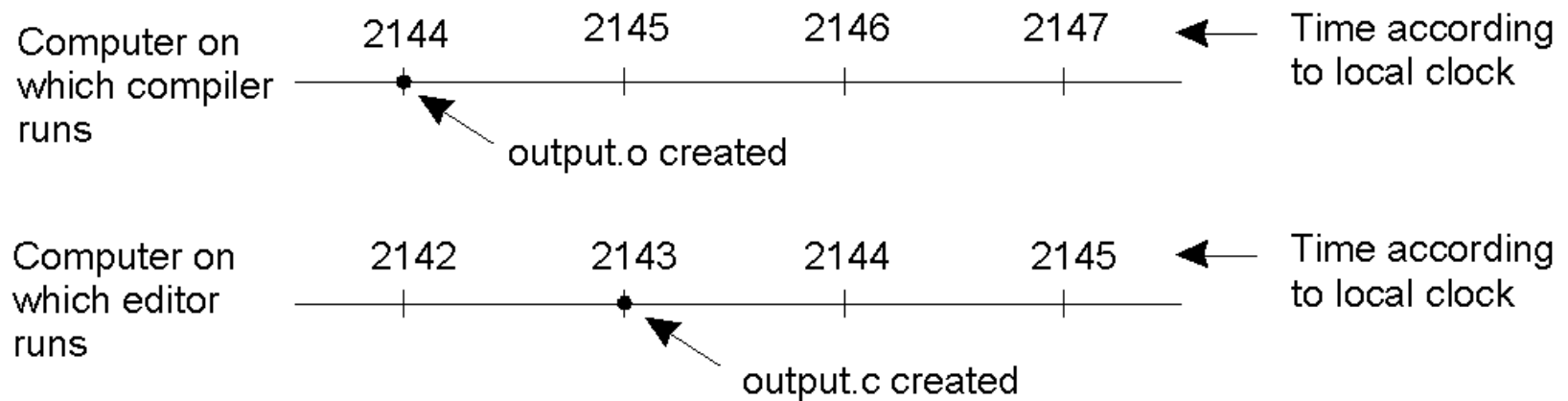


- ❑ **Centralized solution** = Process P_i sends to coordinator
 - If critical section is empty \rightarrow reply = ok
 - Else, request to queue && when queue is empty send reply
 - When process finishes, it sends Release() to coordinator.
 - **Pb** = coordinator dies
- ❑ **Decentralized solution** : Process P_i sends to all other p waits for N replies. \rightarrow Reply =
 - ❑ I don't care (Process doesn't want to enter)
 - ❑ I am before him
 - ❑ I am after him
- **Pb** = one process dies



Why to synchronize multiple clocks in DS?

(Ex : Make-file Tool)



- ❑ When each machine has its own clock, an event that occurred after another event may be assigned an earlier time.
- ❑ Accurate time is needed for some application domains.

Why to synchronize multiple clocks in DS?

(Some other examples)



- ❑ Ex 1 = Replicated Bank System
- ❑ Ex 2 = Group Comm.
- ❑ Ex 3 = Make-file Tool
- ❑ Ex 4 = Distributed Critical Section
- ❑ Ex 5 = Distributed Shared Memory (Replicated Data + Group Communication).
- ❑ Ex 6 = Transaction Serialization.
- ❑ Ex 7 = Distributed Debugging
- ❑ Ex 8 =

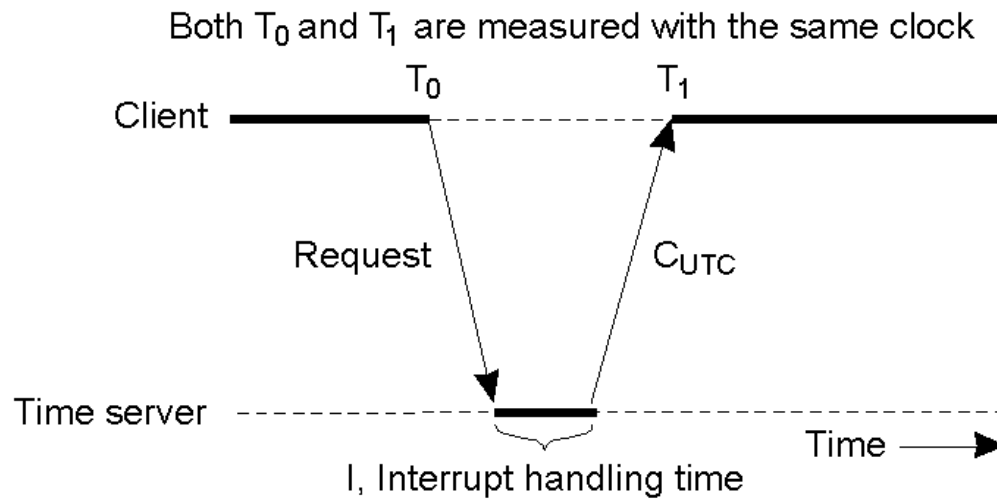


Physical Clock Synchronization

- Centralized / Decentralized algorithms
 - Centralized → Cristian & Berekely Algorithms
 - Decentralized → NTP
- External Synchronization
 - → $| S(t) - C_i(t) | < D$
 - Internet External Time Resource : UTC (Coordinated Universal Time)
 - By GPS → 1μ sec
 - By Radio stations → 0.1-10 m sec
 - By Telephone lines → several m sec
- Internal Synchronization
 - → $| C_i(t) - C_j(t) | < D$



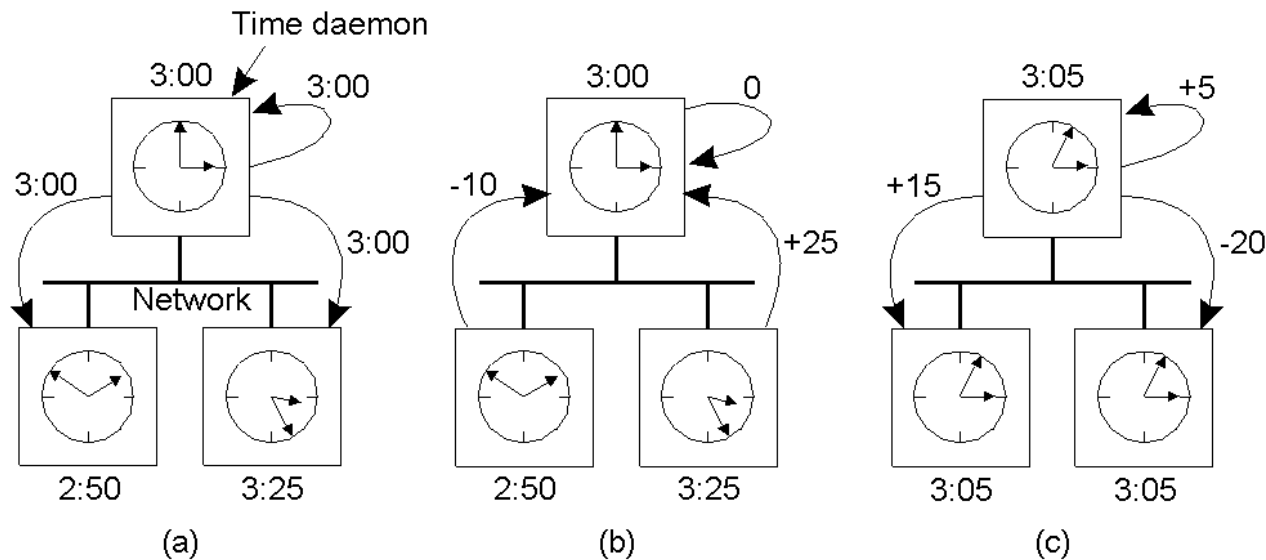
Cristian's method



- time server is passive.
- Periodically, client time = $C_{UTC} + (T_1 - T_0 - I) / 2$
- P_b : sender's clock is fast $\rightarrow C_{UTC} <$ sender's clock (but time must not run backward and changes must be introduced gradually)

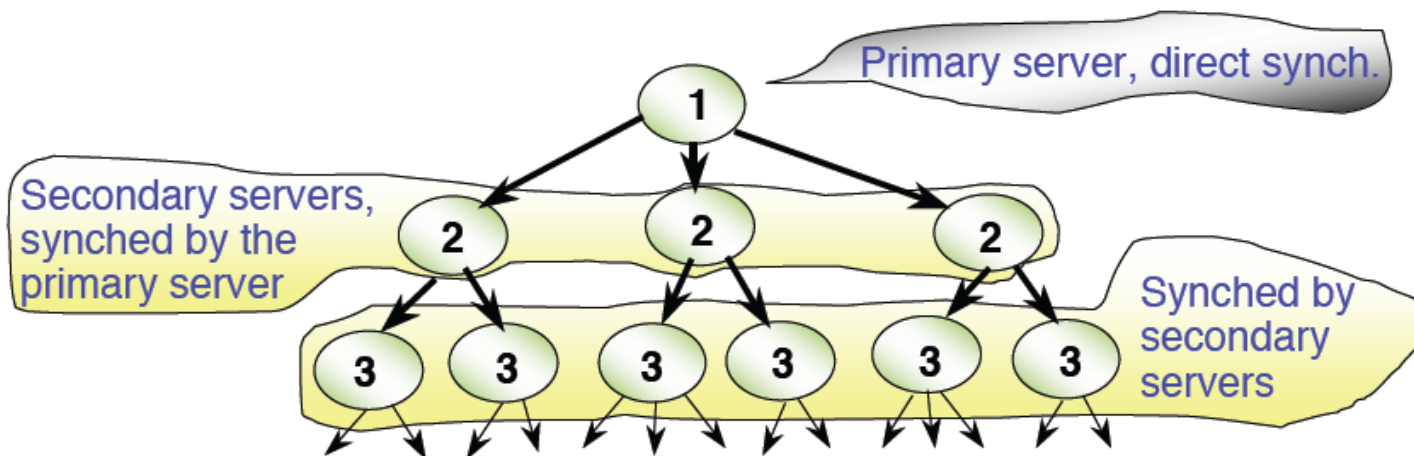
Berkeley Algorithm

- ❑ The time daemon asks all the other machines for their clock values (+ Election when daemon failure)
- ❑ The machines answer
- ❑ The time daemon tells everyone how to adjust their clock



NTP (Network Time Protocol)

- ❑ NTP provides a service enabling clients across the Internet to be synchronized accurately using one of UTC sources.
- ❑ Uses a network of time servers to synchronize all processors on a network.
- ❑ Time servers are connected by a synchronization subnet tree. The root is adjusted directly . Each node synchronizes its children nodes.



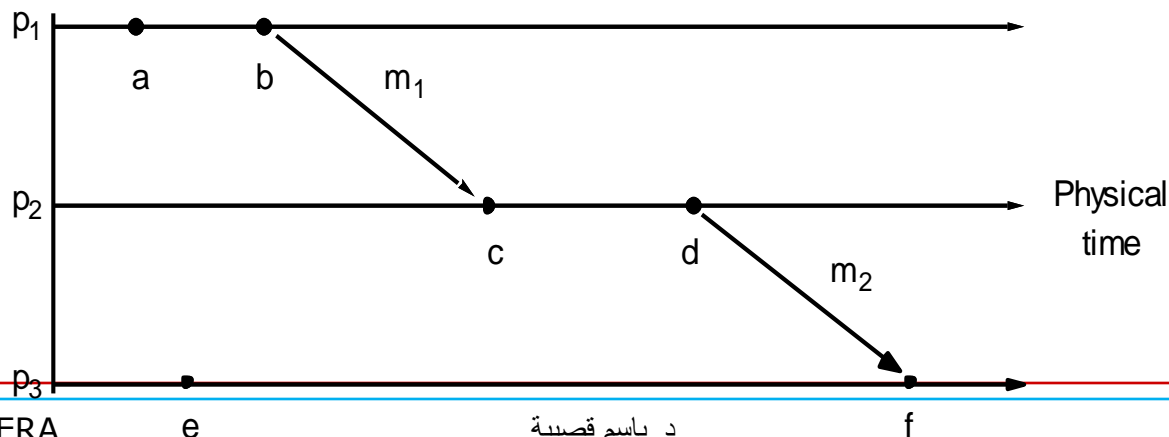


Logical Time



Logical Time

- We can not synchronize perfectly physical clocks. So, we can not give absolute time to events.
- So, we use physical causality (happened-before relationship) to order events (between different processes) :
 - On the same process : $\text{time}(a) < \text{time}(b)$. Then, $a \longrightarrow b$
 - **P1** sends **m** to **p2**. Then, $\text{send}(m) \longrightarrow \text{receive}(m)$
 - (Transitivity) if $a \longrightarrow b \ \&\& \ b \longrightarrow c$. Then, $a \longrightarrow c$



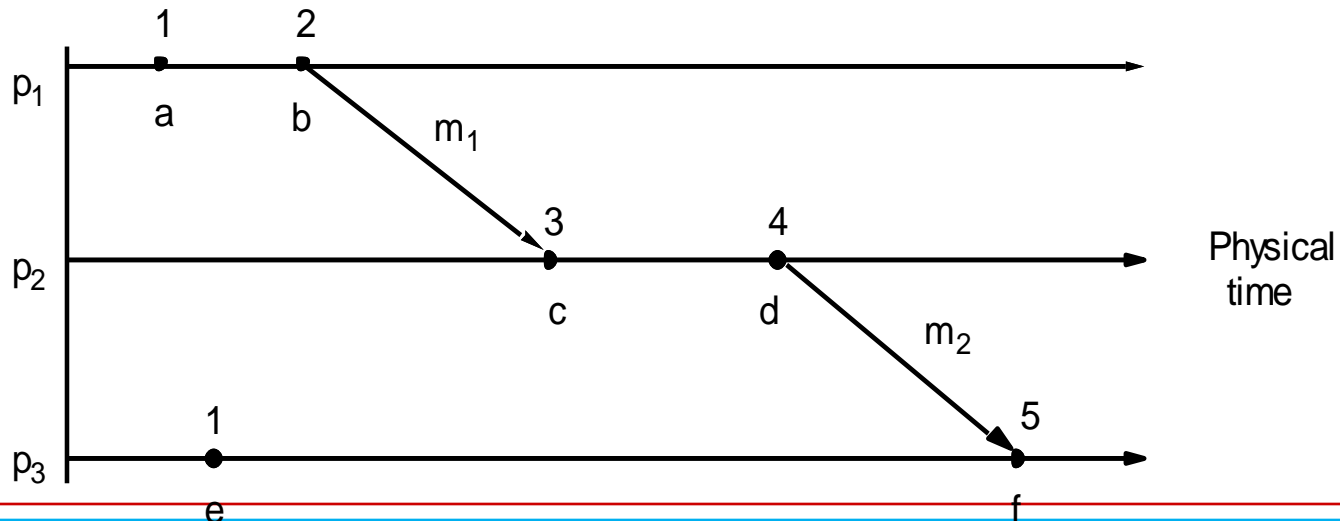


Lamport Method



Logical Clocks (Lamport Method)

- Each process P_i has a logical clock L_i . Initialized $L_i = 0$.
- **R1** : L_i is incremented before each event issued at P_i : $L_i = L_i + 1$
- **R2** : A) when P_i sends message m . It add to m the value $t = L_i$
B) the receiving process P_j calculates $L_j = \max(L_j, t)$ and apply **R1** before timestamp the **receive(m)** event.





Lamport Method Consequences

- If $\mathbf{a} \longrightarrow \mathbf{b}$. Then, $\mathbf{L}(\mathbf{a}) < \mathbf{L}(\mathbf{b})$.

- But, if $\mathbf{L}(\mathbf{x}) < \mathbf{L}(\mathbf{y})$. Then, we can not infer that $\mathbf{x} \longrightarrow \mathbf{y}$
- Not all events can be ordered by the \longrightarrow relation
- $\mathbf{a} \parallel \mathbf{e}$ are concurrent if $\mathbf{a} \not\rightarrow \mathbf{e}$ and $\mathbf{e} \not\rightarrow \mathbf{a}$.
 - EX: \mathbf{a} and $\mathbf{e} = 2$ events at 2 different processes and there is no messages between these processes



Total order of Lamport clocks

- Different events at different processes can have the same timestamps
- Event total ordering by taking into account the process ID
 - $(T_i, i) < (T_j, j)$ if $T_i < T_j$ or $T_i == T_j \ \&\& \ I < J$
 - This order (without physical significance)
 - Sometimes, (useful)
 - EX: to order the process entry to a critical section



Vector Clock



Vector Clock

- ❑ Vector clock = solution of Lamport problems
- ❑ Vector clock is an array of N integer for N processes
- ❑ Processes use vector clock to timestamp its events.
- ❑ processes add vector clock on the sent messages.



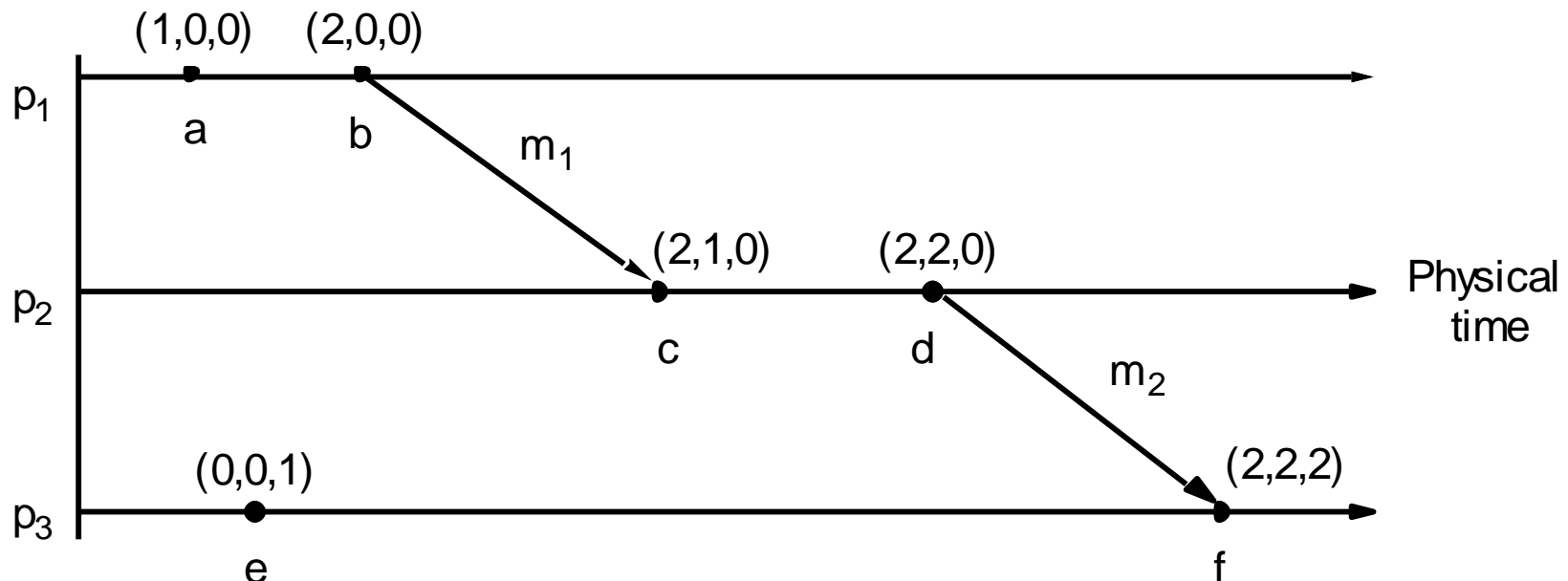
Vector Clock Algorithm

- Each process P_i has its vector clock V_i .
 - Initially, $V_i[j] = 0$ for $j = 1, 2, \dots, N$
- **R1** : $V_i[i]$ is incremented before each event issued at P_i :
$$V_i[i] = V_i[i] + 1$$
- **R2** :
 - **A)** when P_i sends message m . It adds the value $t = V_i$
 - **B)** the receiving process P_j calculates its V_j as :
$$V_j[j] = \max(V_j[j], t[j]) \text{ for } j = 1, 2, \dots, N$$
 - **C)** apply **R1** before timestamp the **receive(m)** event.

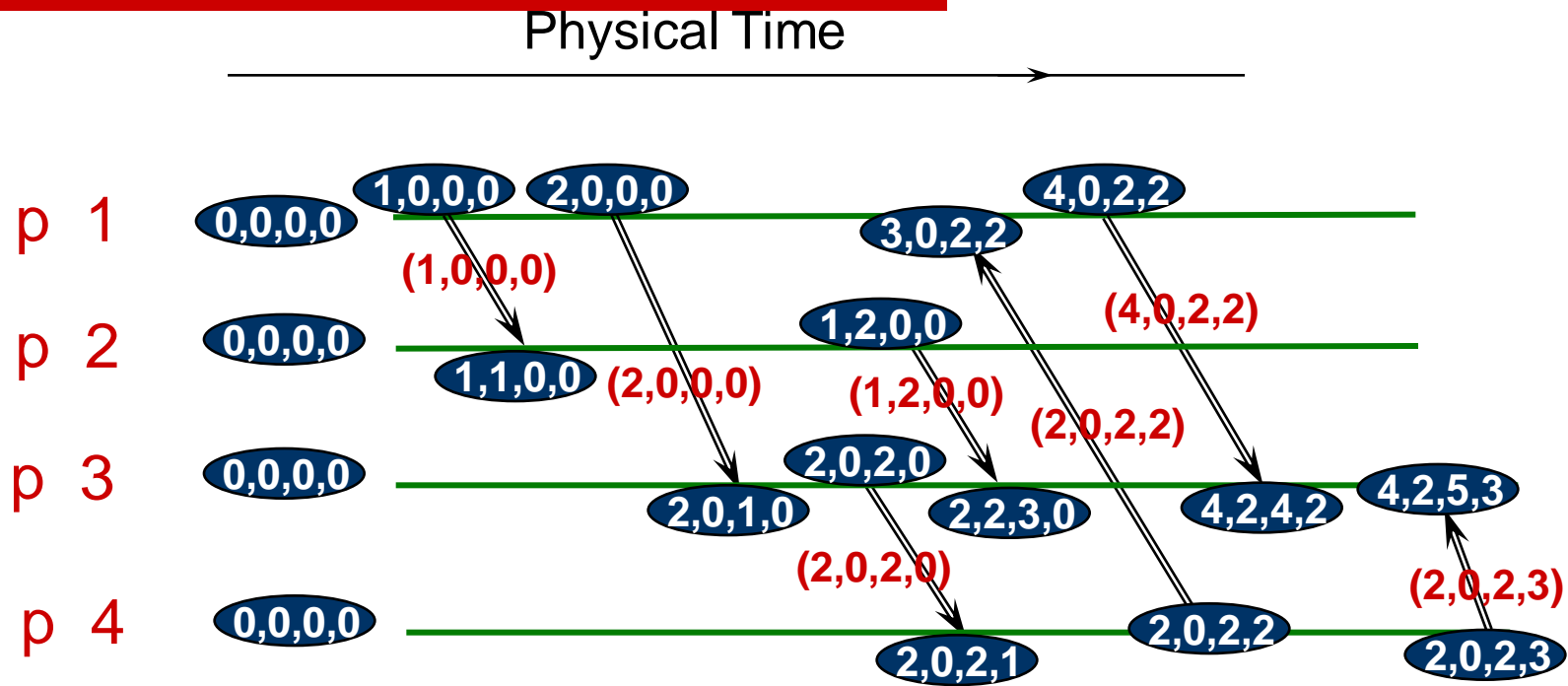


Vector Clock Algorithm

- $V_i[i]$: Nb of events that P_i has timestamped
- $V_i[j]$ ($j \neq i$) : Nb of events occurred at P_j and P_i was affected by them.

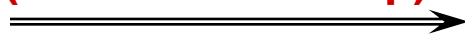


Example: Vector Logical Time



(n,m,p,q) Vector logical clock

(vector timestamp)



Message



Rules for Vector Clock

- $V1 = V2$ iff $V1[j] = V2[j]$ for $j = 1, 2, 3, \dots, N$
- $V1 \leq V2$ iff $V1[j] \leq V2[j]$ for $j = 1, 2, 3, \dots, N$
- $V1 < V2$ iff $V1 \leq V2$ && $V1 \neq V2$



Vector Clock Consequences

- $a \longrightarrow b$ \rightarrow $V(a) < V(b)$
- $V(a) < V(b)$ \rightarrow $a \longrightarrow b$
- $a \parallel b$ **if** neither $V(a) \not\leq V(b)$ nor $V(b) \not\leq V(a)$

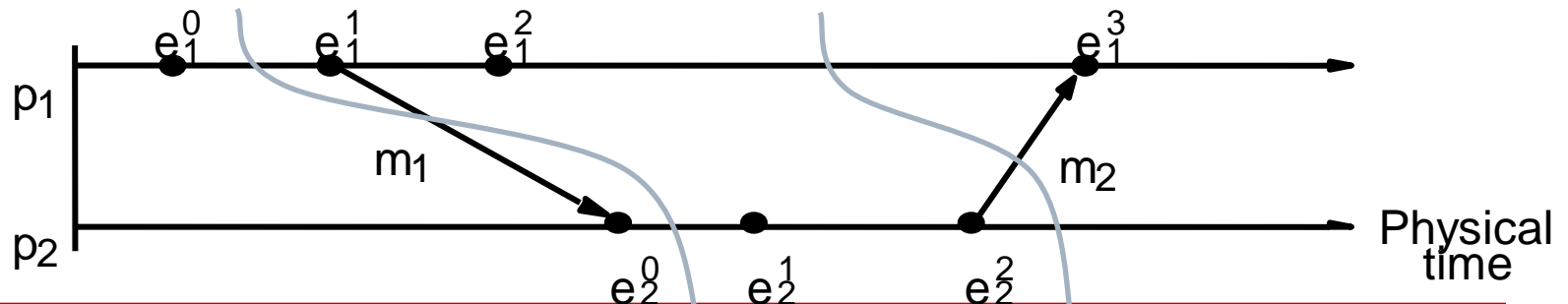


Consistent Cuts



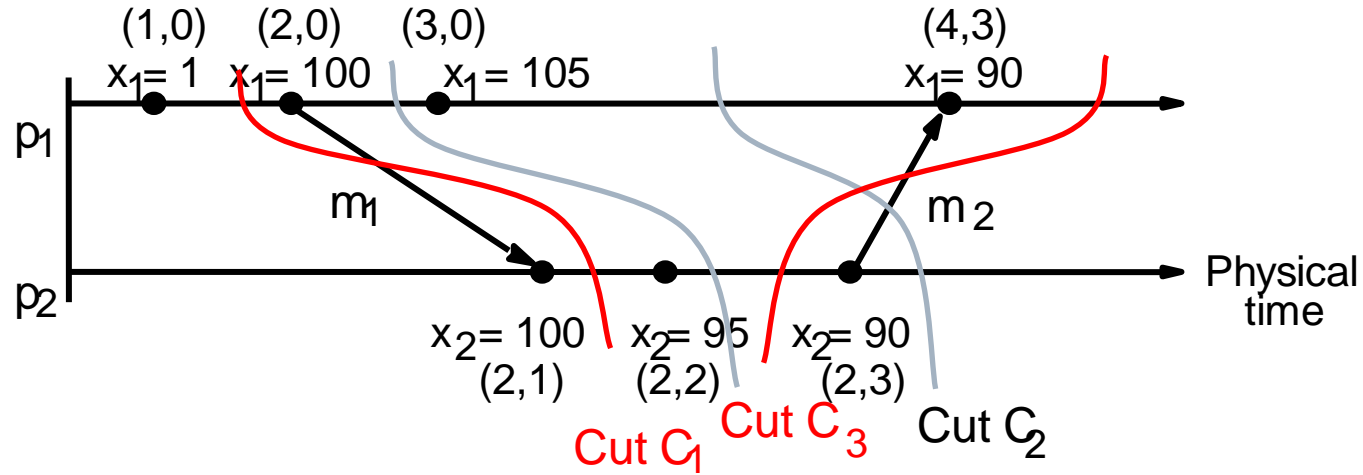
Global states & Consistent cuts

- ❑ If all processes had perfectly synchronized clocks → we could agree a time at which processes record their local states.
- ❑ **Cut** = method to assemble a meaningful global state from local states recorded **at different real times.**
- ❑ Consistent cut = for each event it contains, it also contains all the events that happened-before that event.
- ❑ Consistent Global State = recorded states before a consistent cut.
- ❑ Distributed system execution = transitions between system global states



Distributed Debugging (Logical Time & Global States)

P1	P2
$x_1 = 1$ (1,0)	$(2,1)$ $x_2 = 100$
$(2,0)$ $x_1 = 100$	$x_2 = 95$ (2,2)
$(3,0)$ $x_1 = 105$	$x_2 = 95$ (2,2)
$x_1 = 90$ (4,3)	$x_2 = 95$ (2, 2)
$(3,0)$ $x_1 = 105$	$x_2 = 90$ (2, 3)



- إن كل فعالية تبعث إلى المراقب حالة متحولاتها المحلية + شعاعها الزمني
- المراقب يحترم بأن تكون عدد أحداث P_j و المعروفة لـ P_i عندما تم إرسال حالة P_j ليست بأكبر من تلك الحاصلة لدى P_i عندما أرسلت حالتها

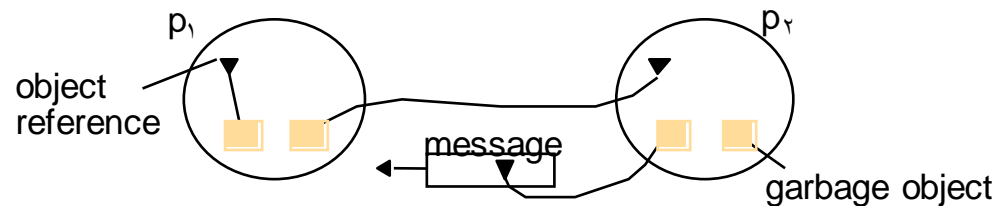


Global States

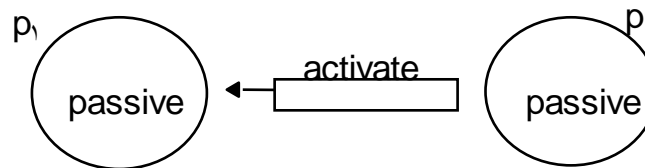
Global States

- Global state = assembly of processes local states + communication channel states (in absence of global time)
- Global states allow us to verify a property in Distributed System Execution like :

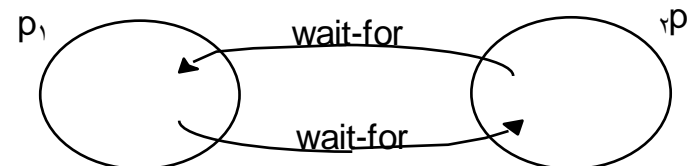
- Garbage Collection



- Termination



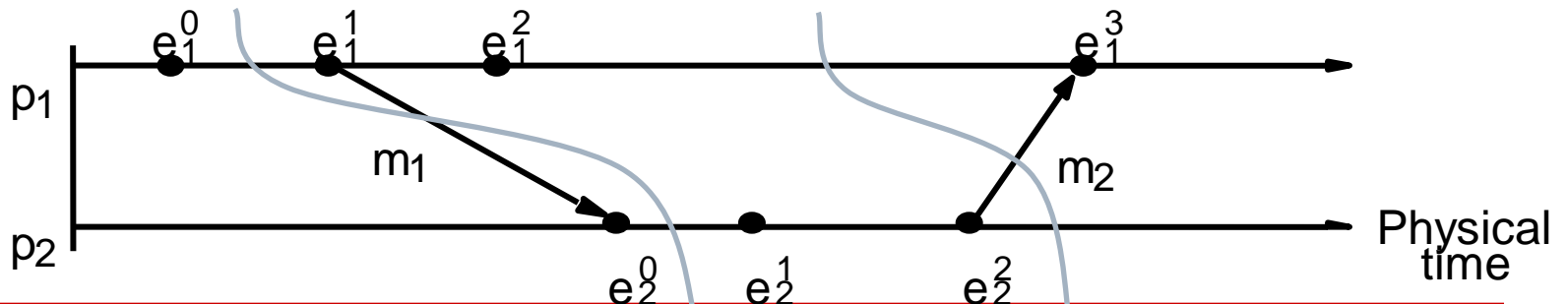
- Deadlock





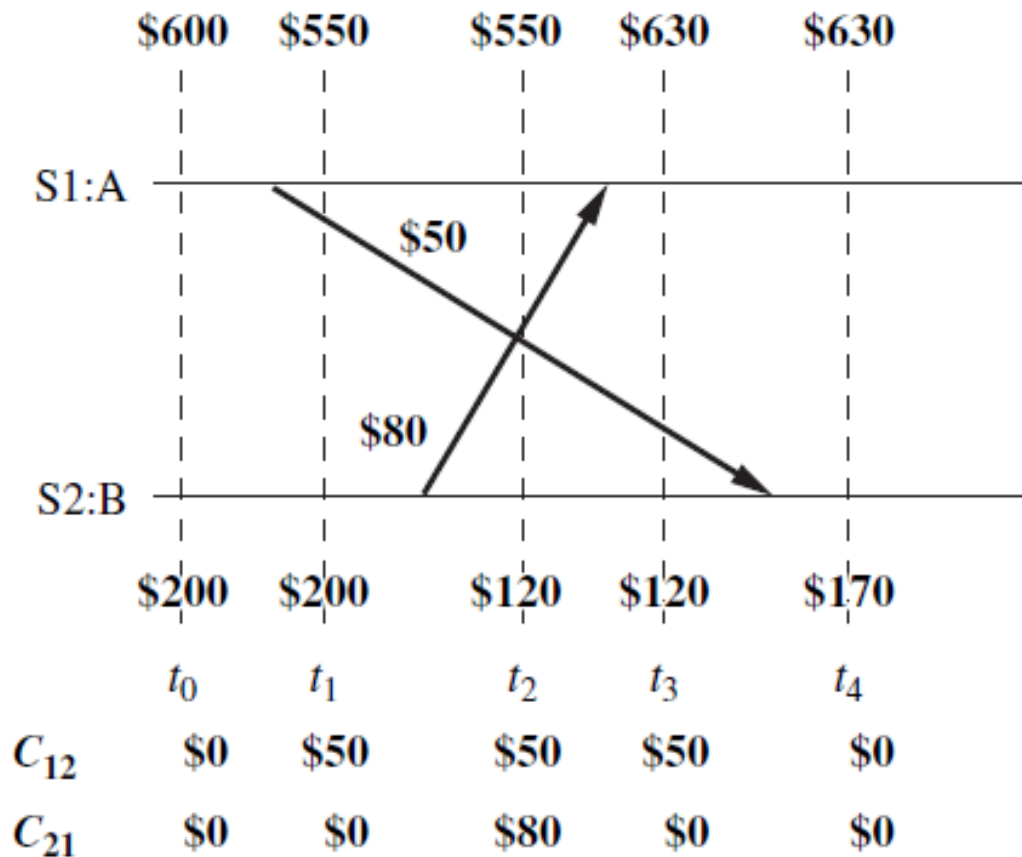
Consistent cuts

- ❑ If all processes had perfectly synchronized clocks → we could agree a time at which processes record their local states.
- ❑ **Cut** = method to assemble a meaningful global state from local states recorded **at different real times.**
- ❑ Consistent cut = for each event it contains, it also contains all the events that happened-before that event.
- ❑ Consistent Global State = recorded states before a consistent cut.
- ❑ Distributed system execution = transitions between system global states





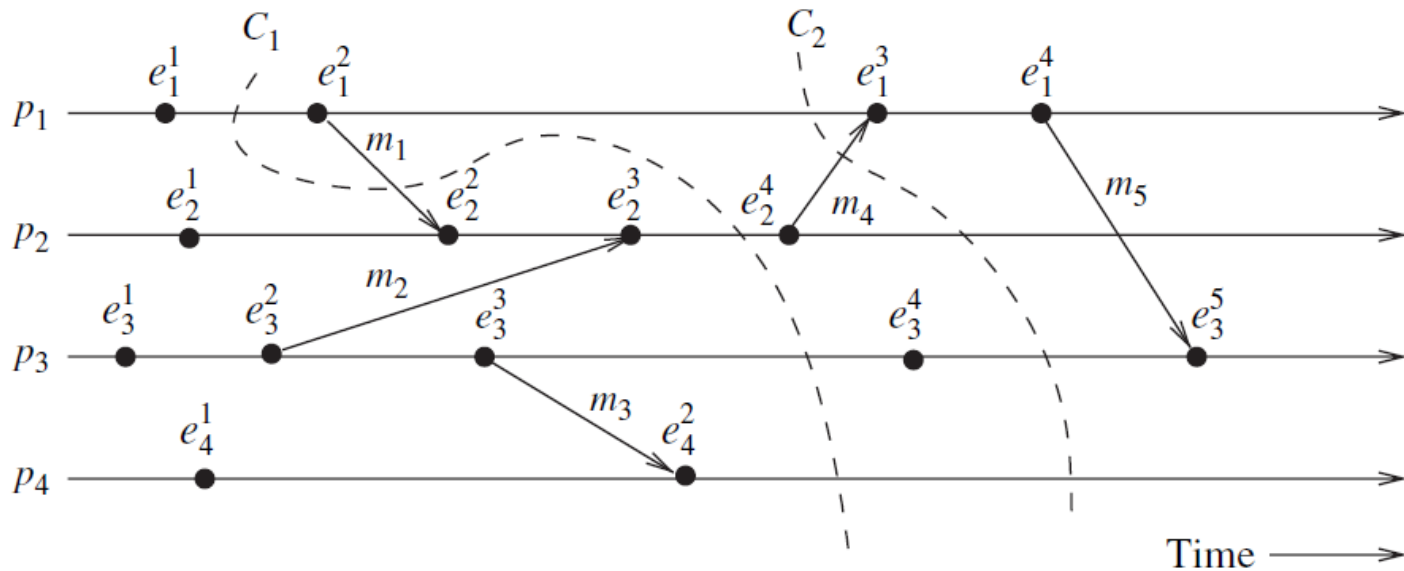
Ex : Banking Sys





Recording Consistent State (Consistent Cut)

- $C1 = send(m_{ij}) \in LSi$
 - $\Rightarrow mij \in Scij \oplus rec(m_{ij}) \in LSj$
- $C2 = send(m_{ij}) \notin LSi \Rightarrow mij \notin Scij \wedge rec(m_{ij}) \notin LSj$



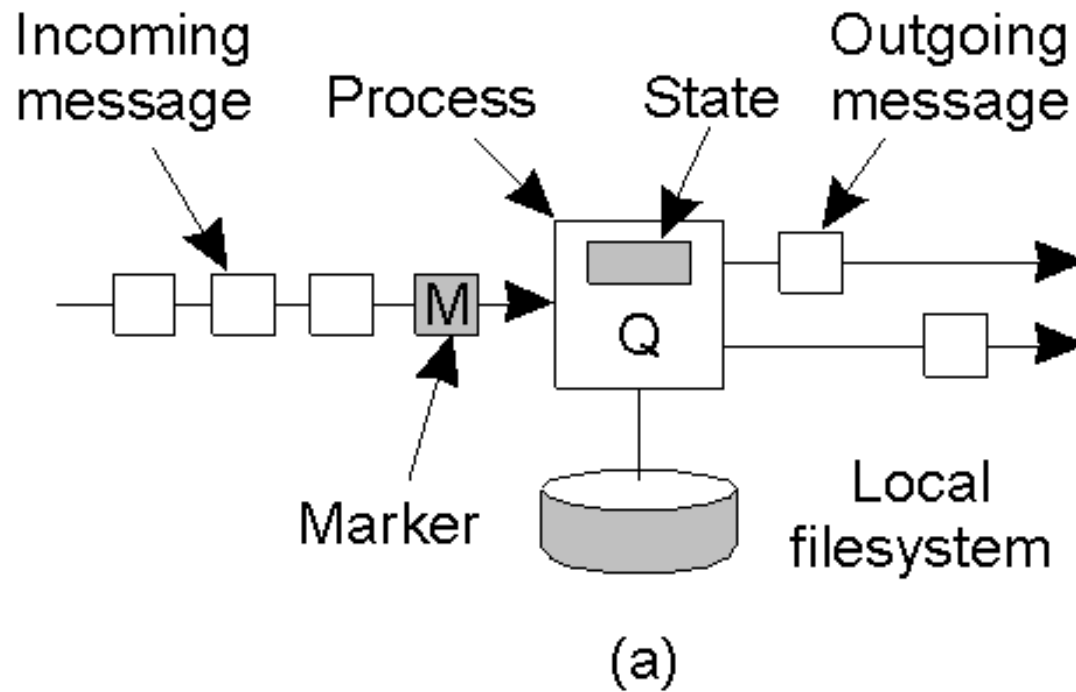


Snapshot Algorithm Bases to Recording Consistent State

- **B1** = Which messages to/or not to be recorded in the snapshot (in a channel state or in a process state)
 - Any message that is sent by a process before recording its snapshot must be recorded
 - Any message that is sent by a process after recording its snapshot, must not be recorded in the global snapshot
- **B2** = How to determine when a process takes its snapshot
 - A process **p_j** must record its snapshot before processing a message **m_{ij}** that was sent by process **p_i** after recording its snapshot

“snapshot” algorithm

- Any process may initiate the algorithm
- There is a path between any two processes





“snapshot” algorithm

□ ***Marker receiving rule for process p_i***

On p_i 's receipt of a *marker* message over channel c :

if (p_i has not yet recorded its state) *it*

records its process state now;

records the state of c as the empty set;

turns on recording of messages arriving over other incoming channels;

else

p_i records the state of c as the set of messages it has received over c

since it saved its state.

end if

□ ***Marker sending rule for process p_i***

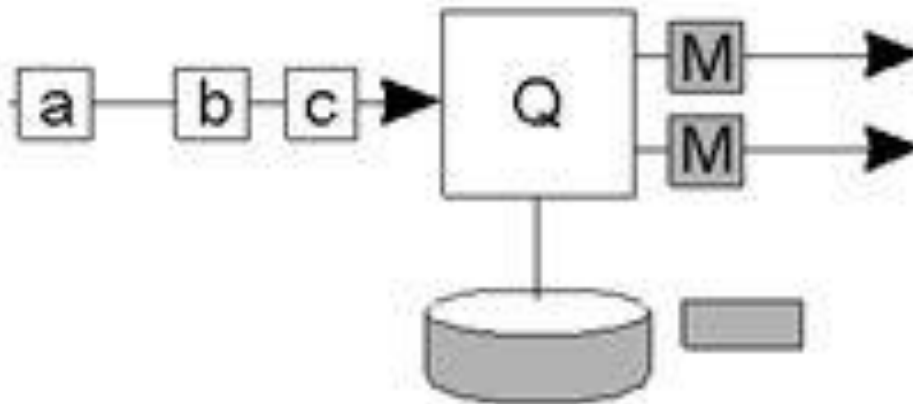
After p_i has recorded its state, for each outgoing channel c :

p_i sends one marker message over c

(before it sends any other message over c).

“snapshot” algorithm (step 1)

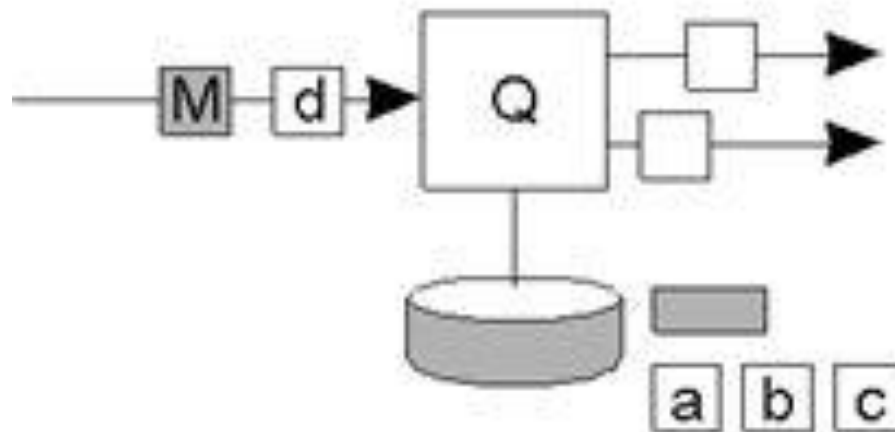
- Processes can continue their execution and send / receive messages while the algorithm runs.



(b)

“snapshot” algorithm (step 2)

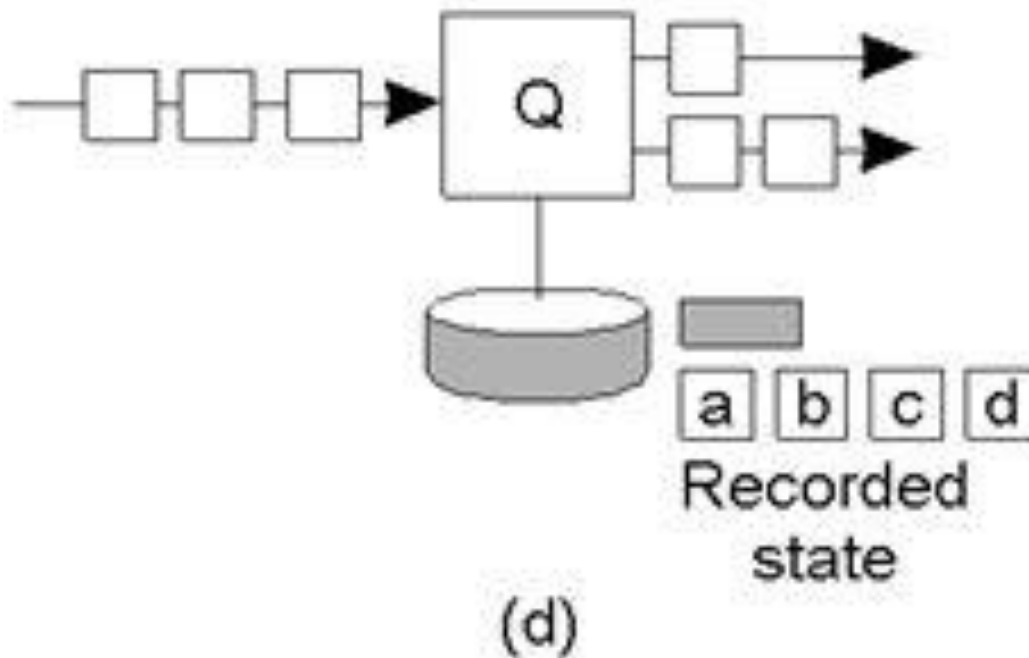
- while any process can initiate the algorithm (multiple executions). So, we need ID for M



(c)

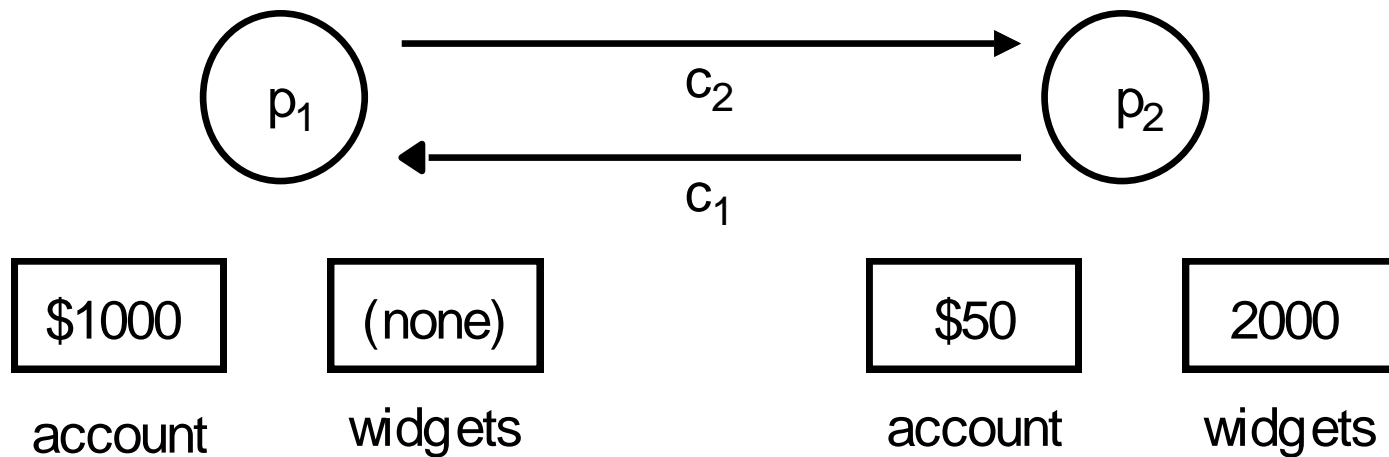
“snapshot” algorithm (step 3)

- Process finishes when it receives M on all its incoming channels



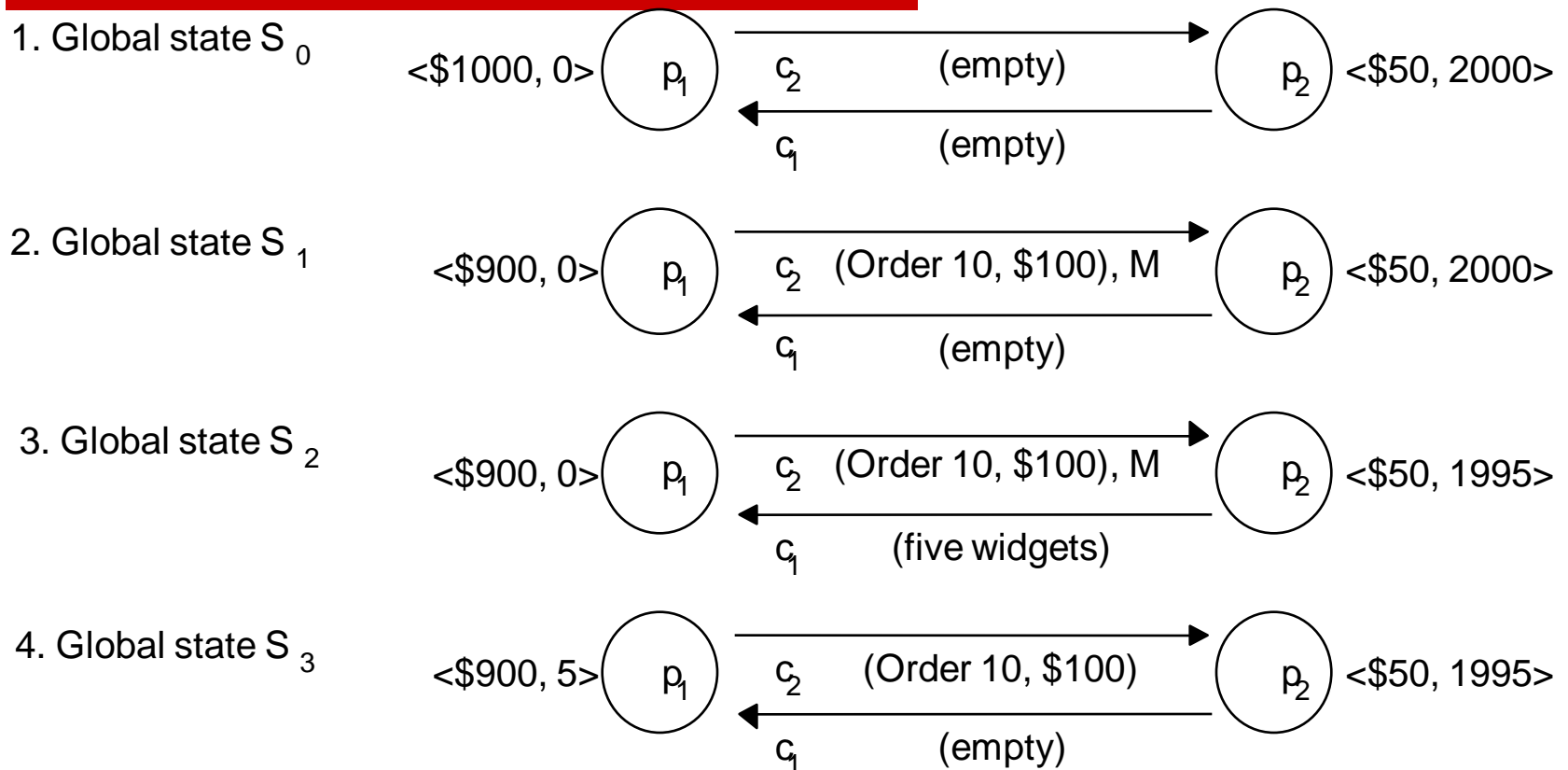


“snapshot” algorithm - Example

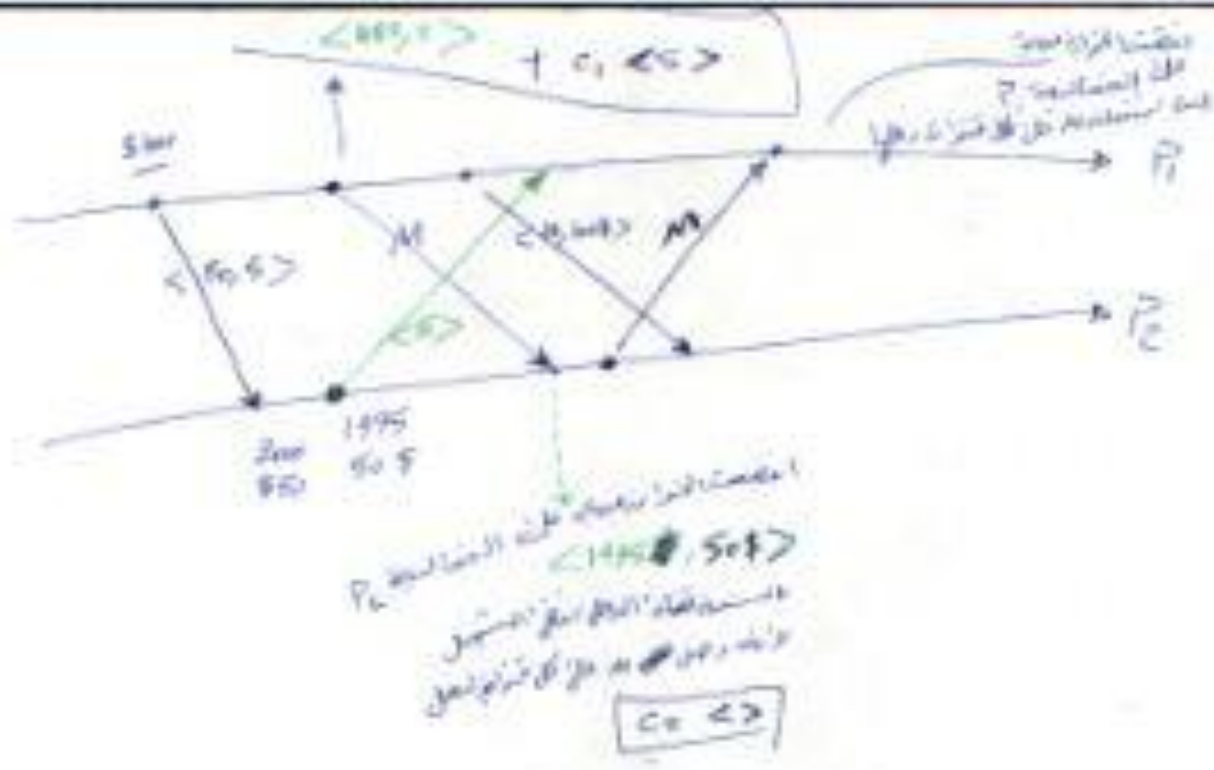


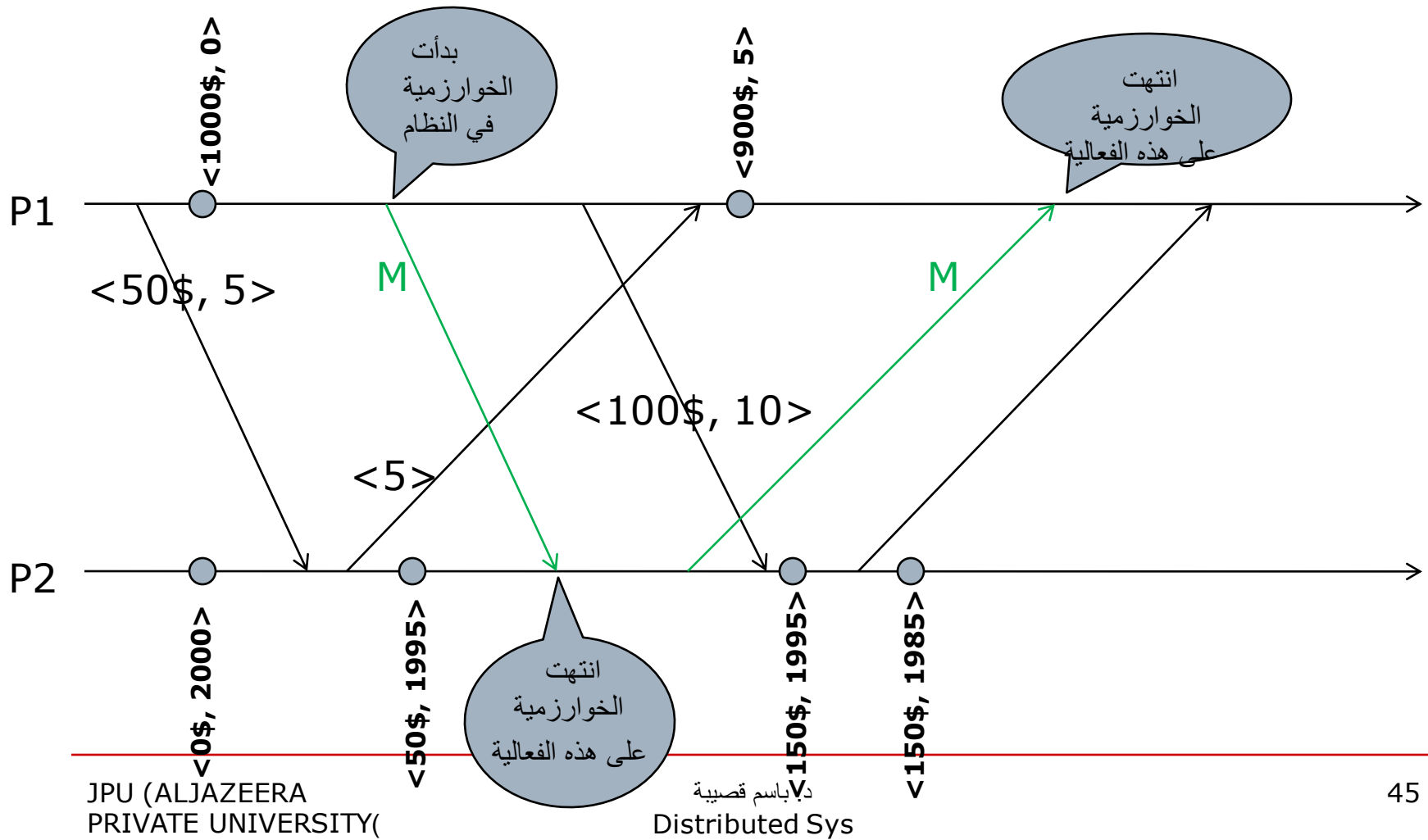


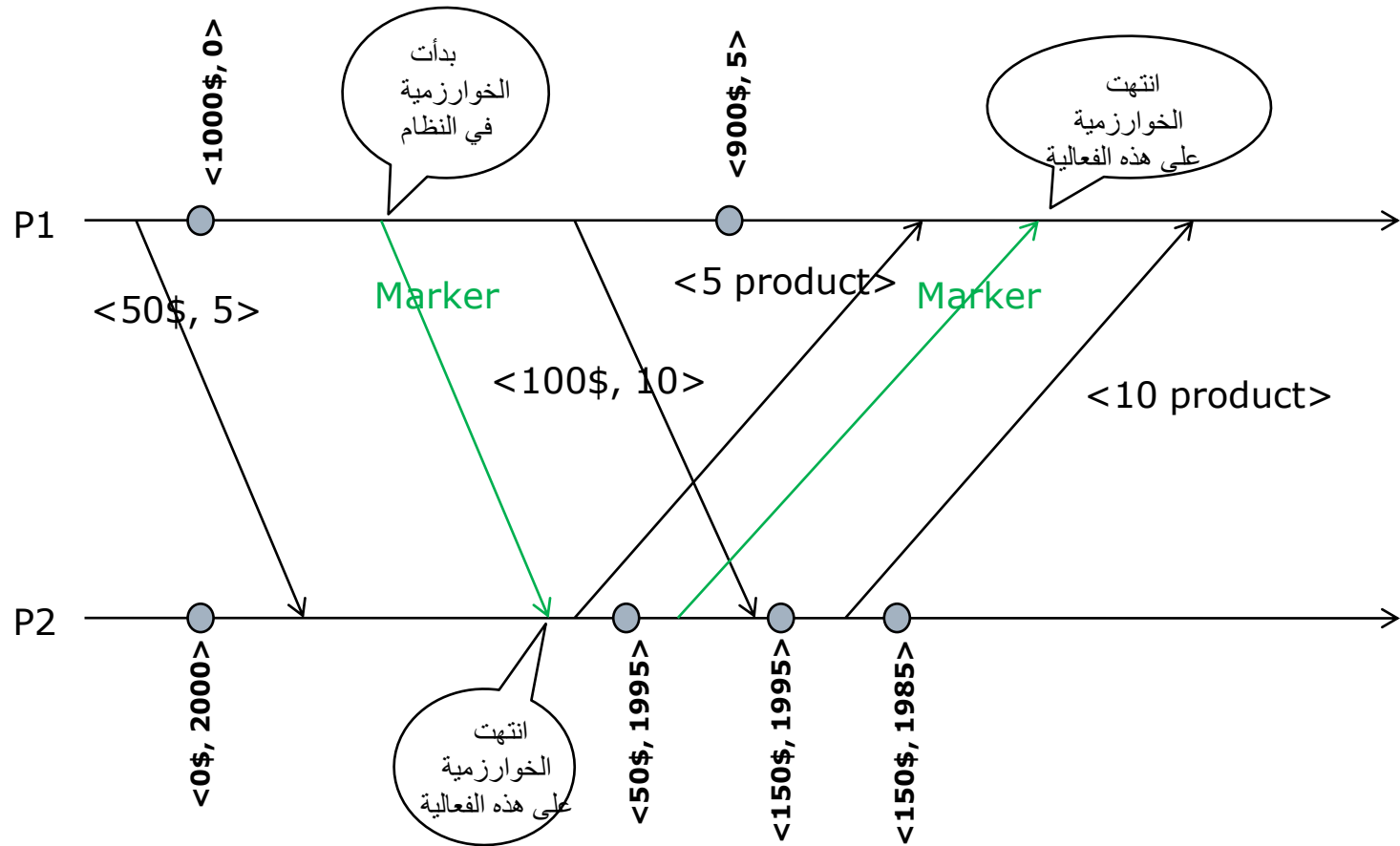
“snapshot” algorithm - Example

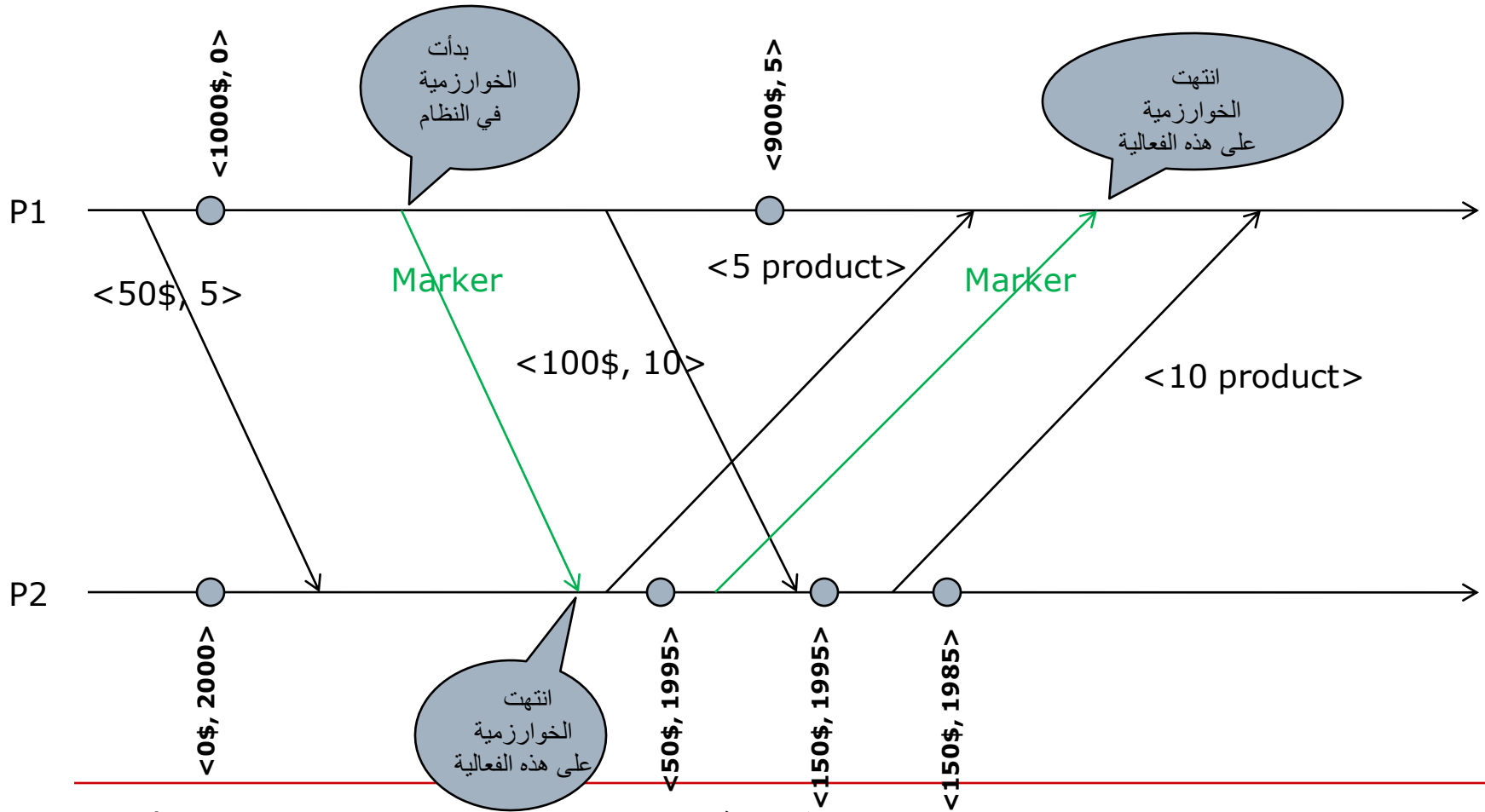


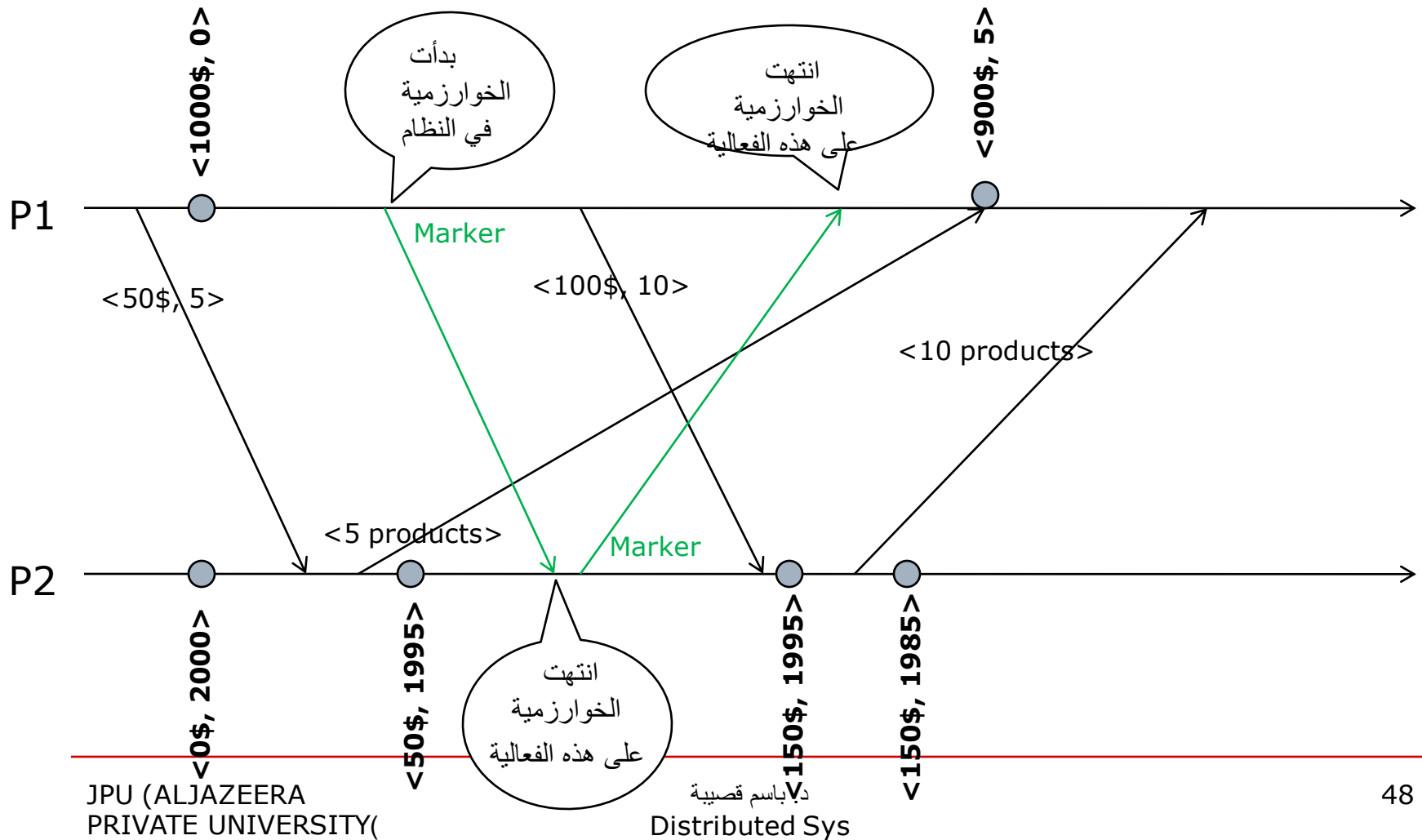
‘snap’ Global State = P1 $\langle \$1000, 0 \rangle$, P2 $\langle \$50, 1995 \rangle$, C1 $\langle 5 \rangle$, C2 $\langle \rangle$













Questions ?