

Distributed Systems

CS425/ECE428

Feb 1 2023

Instructor: Radhika Mittal

Logistics Related

- HW1 released! Due on Feb 15th
 - You should be able to solve the first three questions right-away.
 - You should be able to solve the fourth question by the end of this class (hopefully).
 - You should be able to solve the fifth question by the end of next class.
- Newly registered students:
 - Please make sure you have access to Campuswire and Gradescope
 - If you are in 4 credits, make sure you have been allocated a VM cluster for the MPs.
 - Email Manoj (netid: gmk6) to get the required access.
- *Please say your name before speaking up in class 😊*

Recap: Logical timestamps

- How to reason about ordering of events across processes without synchronized clocks?
- Happened-before Relationship
- Lamport Logical Clock
- Vector Clock

Today's agenda

- **Global State**
 - Chapter 14.5
 - Goal: reason about how to capture the state across all processes of a distributed system without requiring time synchronization.

Process, state, events

- Consider a system with n processes: $\langle p_1, p_2, p_3, \dots, p_n \rangle$.
- Each process p_i is associated with state s_i .
 - State includes values of all local variables, affected files, etc.
- Each channel can also be associated with a state.
 - Which messages are currently *pending* on the channel.
 - Can be computed from process' state:
 - Record when a process sends and receives messages.
 - if p_i sends a message that p_j has not yet received, it is pending on the channel.
- State of a process (or a channel) gets transformed when an event occurs. 3 types of events:
 - local computation, sending a message, receiving a message.

Capturing a global snapshot

- Useful to capture a global snapshot of the system:
 - *Checkpointing* the system state.
 - Reasoning about unreferenced objects (for garbage collection).
 - Deadlock detection.
 - Distributed debugging.

Capturing a global snapshot

- Global state or global snapshot is state of each process (and each channel) in the system at a given *instant of time*.
- Difficult to capture a global snapshot of the system.
- Strawman:
 - Each process records its state at 2:05pm.
 - We get the global state of the system at 2:05pm.
 - *But precise clock synchronization is difficult to achieve.*
- **How do we capture global snapshots without precise time synchronization across processes?**

Some more notations and definitions

- State of a process (or a channel) gets transformed when an *event* occurs.
- 3 types of events:
 - local computation, sending a message, receiving a message.
- e_i^n is the n^{th} event at p_i .

Some more notations and definitions

- For a process p_i , where events e_i^0, e_i^1, \dots occur:

$$\text{history}(p_i) = h_i = \langle e_i^0, e_i^1, \dots \rangle$$

$$\text{prefix history}(p_i^k) = h_i^k = \langle e_i^0, e_i^1, \dots, e_i^k \rangle$$

s_i^k : p_i 's state immediately after k^{th} event.

- For a set of processes $\langle p_1, p_2, p_3, \dots, p_n \rangle$:

$$\text{global history: } H = \cup_i (h_i)$$

Some more notations and definitions

- For a process p_i , where events e_i^0, e_i^1, \dots occur:

$$\text{history}(p_i) = h_i = \langle e_i^0, e_i^1, \dots \rangle$$

$$\text{prefix history}(p_i^k) = h_i^k = \langle e_i^0, e_i^1, \dots, e_i^k \rangle$$

s_i^k : p_i 's state immediately after k^{th} event.

- For a set of processes $\langle p_1, p_2, p_3, \dots, p_n \rangle$:

$$\text{global history: } H = \cup_i (h_i)$$

Some more notations and definitions

- For a process p_i , where events e_i^0, e_i^1, \dots occur:

$$\text{history}(p_i) = h_i = \langle e_i^0, e_i^1, \dots \rangle$$

$$\text{prefix history}(p_i^k) = h_i^k = \langle e_i^0, e_i^1, \dots, e_i^k \rangle$$

s_i^k : p_i 's state immediately after k^{th} event.

- For a set of processes $\langle p_1, p_2, p_3, \dots, p_n \rangle$:

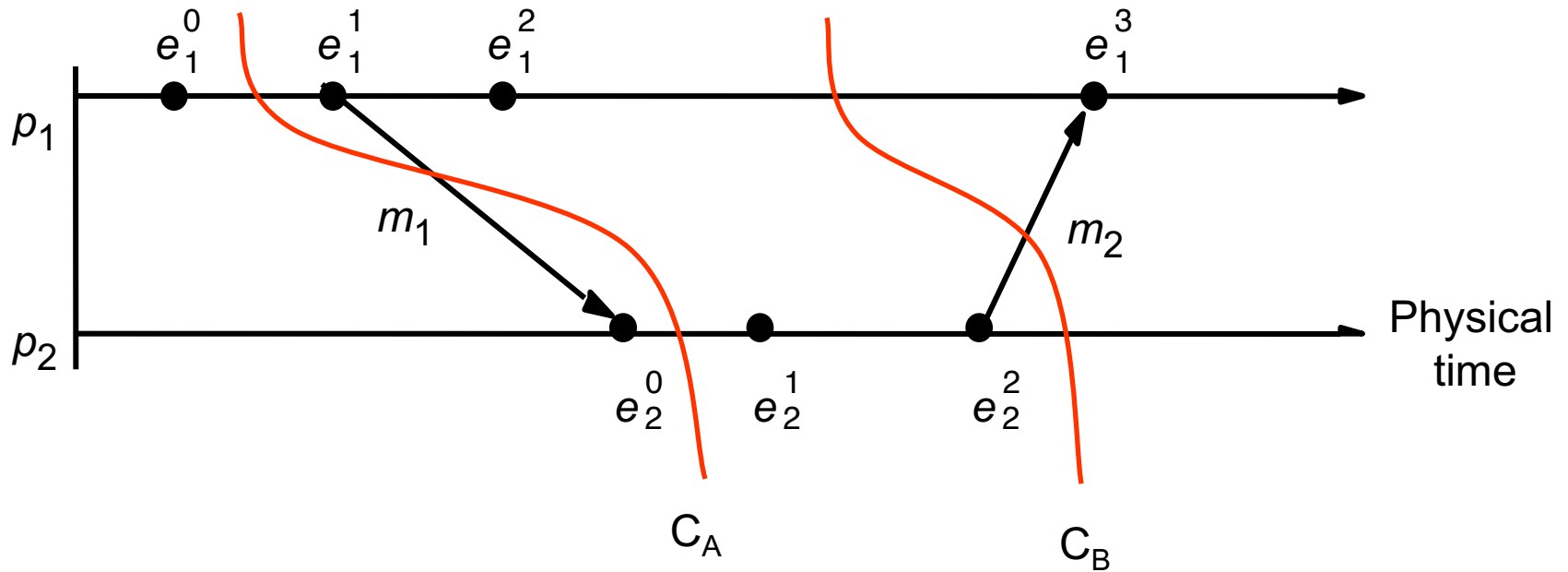
$$\text{global history: } H = \cup_i (h_i)$$

$$\text{a cut } C \subseteq H = h_1^{c_1} \cup h_2^{c_2} \cup \dots \cup h_n^{c_n}$$

$$\text{the frontier of } C = \{e_i^{c_i}, i = 1, 2, \dots, n\}$$

$$\text{global state } S \text{ that corresponds to cut } C = \cup_i (s_i^{c_i})$$

Example: Cut



$$C_A: \langle e_1^0, e_2^0 \rangle$$

Frontier of C_A :

$$C_B: \langle e_1^0, e_1^1, e_1^2, e_2^0, e_2^1, e_2^2 \rangle$$

Frontier of C_B :

Some more notations and definitions

- For a process p_i , where events e_i^0, e_i^1, \dots occur:

$$\text{history}(p_i) = h_i = \langle e_i^0, e_i^1, \dots \rangle$$

$$\text{prefix history}(p_i^k) = h_i^k = \langle e_i^0, e_i^1, \dots, e_i^k \rangle$$

s_i^k : p_i 's state immediately after k^{th} event.

- For a set of processes $\langle p_1, p_2, p_3, \dots, p_n \rangle$:

$$\text{global history: } H = \cup_i (h_i)$$

$$\text{a cut } C \subseteq H = h_1^{c_1} \cup h_2^{c_2} \cup \dots \cup h_n^{c_n}$$

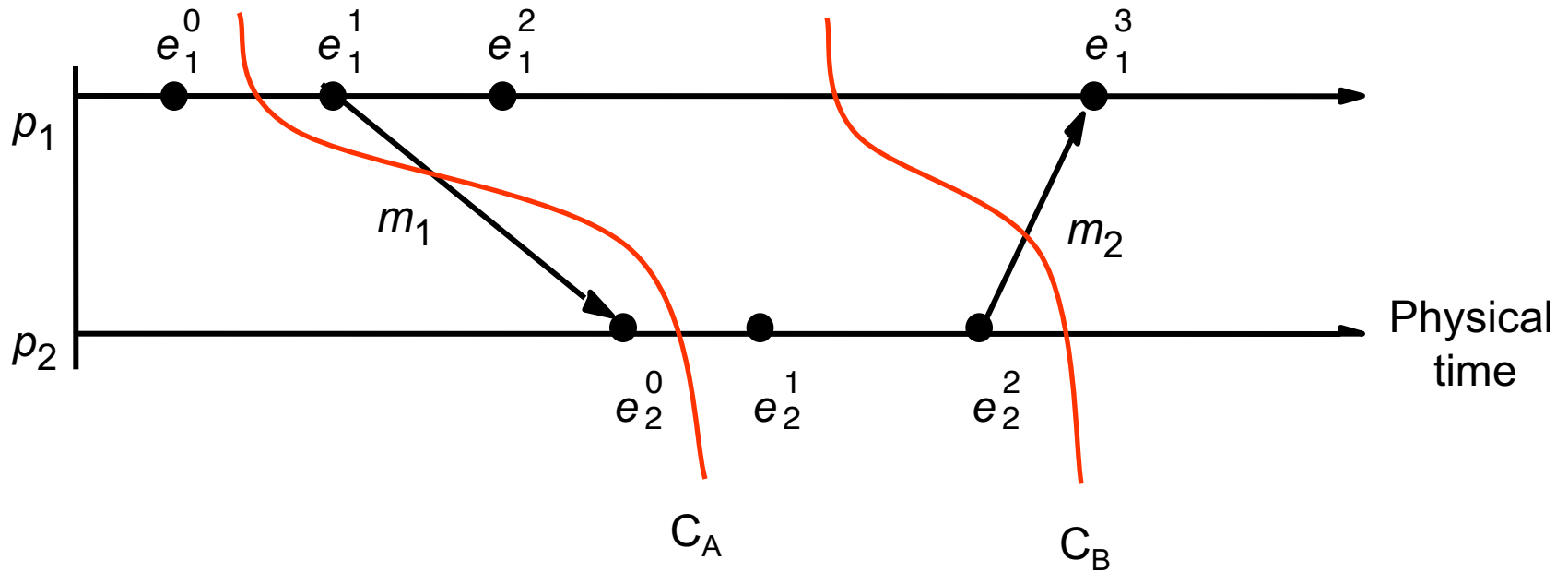
$$\text{the frontier of } C = \{e_i^{c_i}, i = 1, 2, \dots, n\}$$

$$\text{global state } S \text{ that corresponds to cut } C = \cup_i (s_i^{c_i})$$

Consistent cuts and snapshots

- A cut C is **consistent** if and only if
$$\forall e \in C \text{ (if } f \rightarrow e \text{ then } f \in C)$$

Example: Cut



$C_A: \langle e_1^0, e_2^0 \rangle$
 Frontier of $C_A: \{e_1^0, e_2^0\}$

Inconsistent cut.

$C_B: \langle e_1^0, e_1^1, e_1^2, e_2^0, e_2^1, e_2^2 \rangle$
 Frontier of $C_B: \{e_1^2, e_2^2\}$

Consistent cut.

Consistent cuts and snapshots

- A cut \mathbf{C} is **consistent** if and only if
$$\forall e \in \mathbf{C} \text{ (if } f \rightarrow e \text{ then } f \in \mathbf{C}\text{)}$$
- A global state \mathbf{S} is consistent if and only if it corresponds to a consistent cut.

How to capture global state?

- Ideally: state of each process (and each channel) in the system *at a given instant of time*.
 - Difficult to capture -- requires precisely synchronized time.
- Relax the problem: find a consistent global state.
 - For a system with n processes $\langle p_1, p_2, p_3, \dots, p_n \rangle$, capture the state of the system after the c_i^{th} event at process p_i .
 - State corresponding to the *cut* defined by frontier events $\{e_i^{c_i}, \text{ for } i = 1, 2, \dots, n\}$.
 - We want the state to be consistent.
 - Must correspond to a consistent cut.

How to find a consistent global state that corresponds to a consistent cut?

Chandy-Lamport Algorithm

- Goal:
 - Record a global snapshot
 - Process state (and channel state) for a set of processes.
 - The recorded global state is consistent.
- Identifies a consistent cut.
- Records corresponding state locally at each process.

Chandy-Lamport Algorithm

- *System model and assumptions:*
 - System of n processes: $\langle p_1, p_2, p_3, \dots, p_n \rangle$.
 - There are two uni-directional communication channels between each ordered process pair : p_j to p_i and p_i to p_j .
 - Communication channels are FIFO-ordered (first in first out).
 - if p_i sends m before m' to p_j , then p_j receives m before m' .
 - All messages arrive intact, and are not duplicated.
 - No failures: neither channel nor processes fail.

Chandy-Lamport Algorithm

- *Requirements:*
 - Snapshot should not interfere with normal application actions, and it should not require application to stop sending messages.
 - Any process may initiate algorithm.

Chandy-Lamport Algorithm Intuition

- First, initiator p_i :
 - records its own state.
 - creates a special **marker** message.
 - sends the **marker** to all other process.

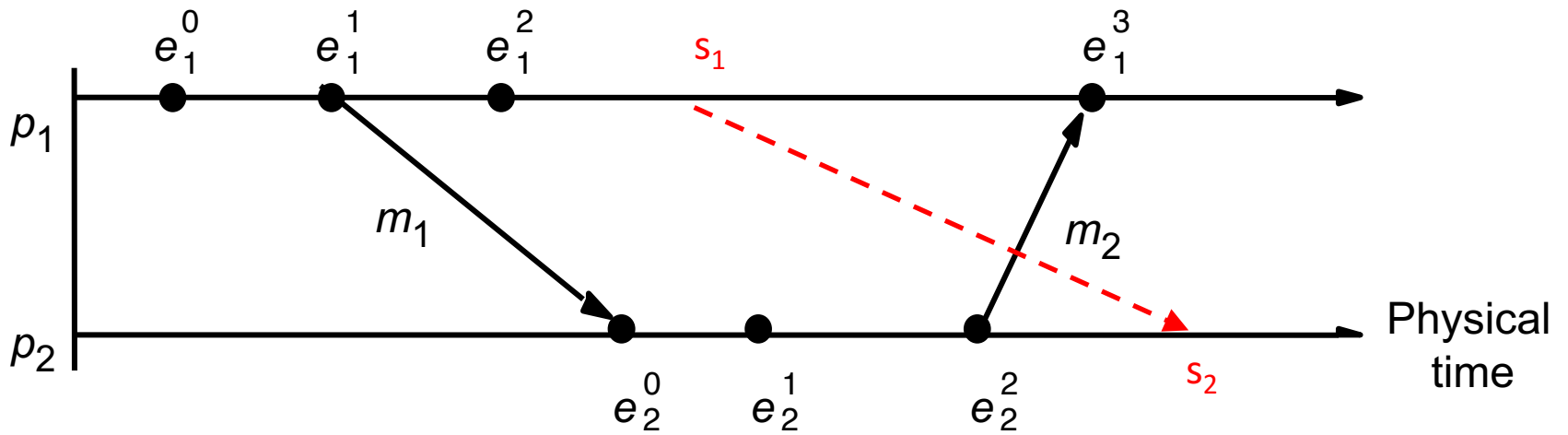
- When a process receives a **marker**.
 - records its own state.

Chandy-Lamport Algorithm Intuition

- First, initiator p_i :
 - records its own state.
 - creates a special **marker** message.
 - sends the **marker** to all other process.

- When a process receives a **marker**.
 - records its own state.

Chandy-Lamport Algorithm Intuition



Cut frontier: $\{e_1^2, e_2^2\}$

Chandy-Lamport Algorithm Intuition

- First, initiator p_i :
 - records its own state.
 - creates a special **marker** message.
 - sends the **marker** to all other process.

- When a process receives a **marker**.
 - records its own state.

This captures the local state at each process.

How do we ensure the state is consistent?

What about the channel state?

Chandy-Lamport Algorithm Intuition

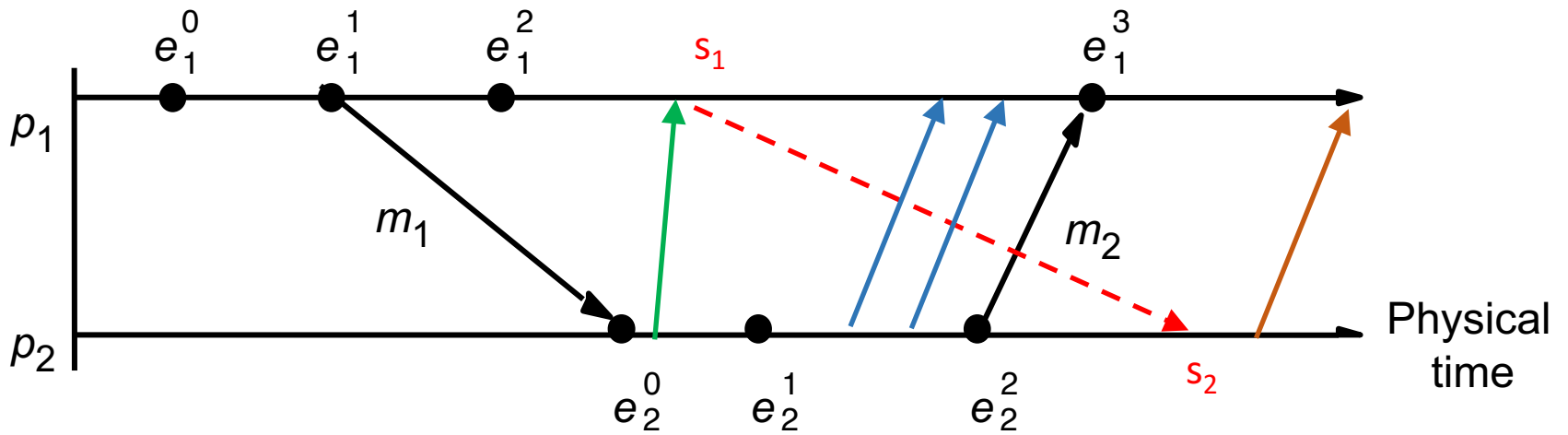
- First, initiator p_i :
 - **records** its own state.
 - creates a special **marker** message.
 - sends the **marker** to all other process.

- When a process receives a **marker**.
 - If marker is received for the first time.
 - **records** its own state.
 - sends **marker** on all other channels.

Leads to a consistent cut (we'll get back to it)

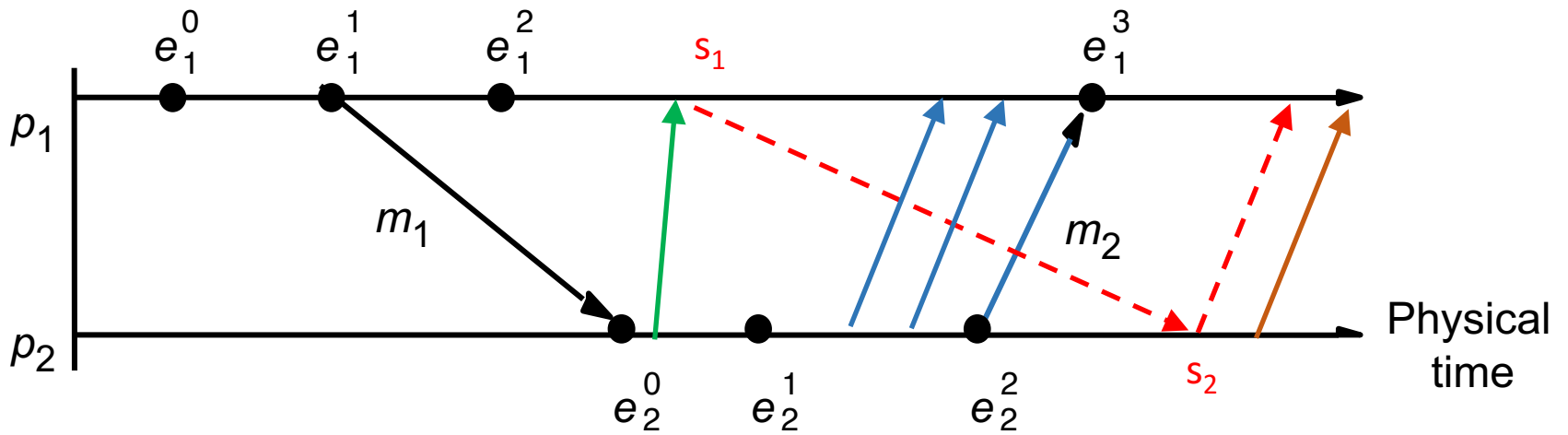
What about the channel state?

Chandy-Lamport Algorithm Intuition



Cut frontier: $\{e_1^2, e_2^2\}$

Chandy-Lamport Algorithm Intuition



Cut frontier: $\{e_1^2, e_2^2\}$

Chandy-Lamport Algorithm Intuition

- First, initiator p_i :
 - records its own state.
 - creates a special **marker** message.
 - sends the **marker** to all other process.
 - start recording messages received on other channels.
 - until a marker is received on a channel.
- When a process receives a **marker**.
 - If marker is received for the first time.
 - records its own state.
 - sends **marker** on all other channels.
 - start recording messages received on other channels.
 - until a marker is received on a channel.

Chandy-Lamport Algorithm

- First, initiator p_i :
 - **records** its own state.
 - creates a special **marker** message.
 - for $j=1$ to n except i
 - p_i **sends** a **marker** message on outgoing channel c_{ij}
 - **starts recording** the incoming messages on each of the incoming channels at $p_i : c_{ji}$ (for $j=1$ to n except i).

Chandy-Lamport Algorithm

Whenever a process p_i receives a **marker** message on an incoming channel c_{ki}

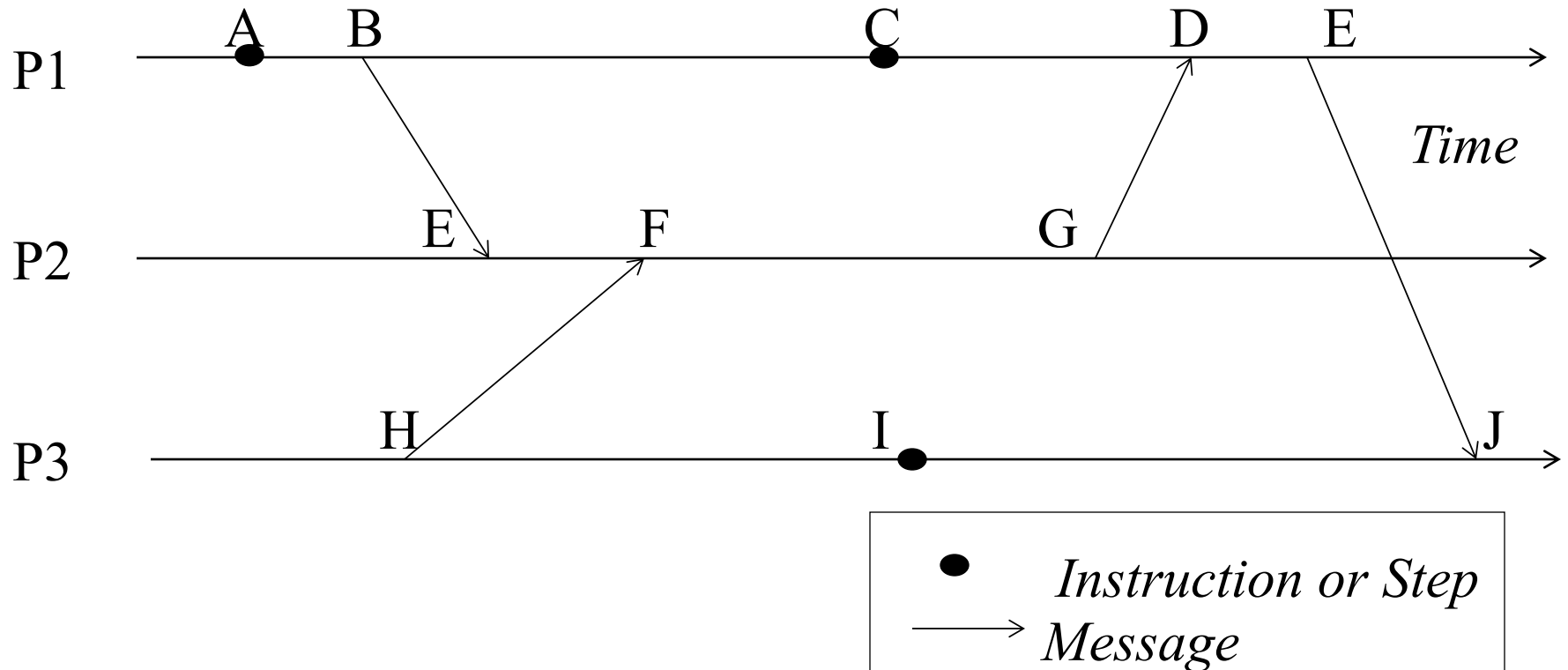
- if (this is the first **marker** p_i is seeing)
 - p_i **records** its own state first
 - **marks the state of channel c_{ki} as “empty”**
 - for $j=1$ to n except i
 - p_i **sends** out a **marker** message on outgoing channel c_{ij}
 - **starts recording** the incoming messages on each of the incoming channels at $p_i : c_{ji}$ (for $j=1$ to n except i and k).
- else // already seen a **marker** message
 - **mark** the state of channel c_{ki} as all the messages that have arrived on it **since recording was turned on for c_{ki}**

Chandy-Lamport Algorithm

The algorithm terminates when

- All processes have received a **marker**
 - To record their own state
- All processes have received a **marker** on all the $(n-1)$ incoming channels
 - To record the state of all channels

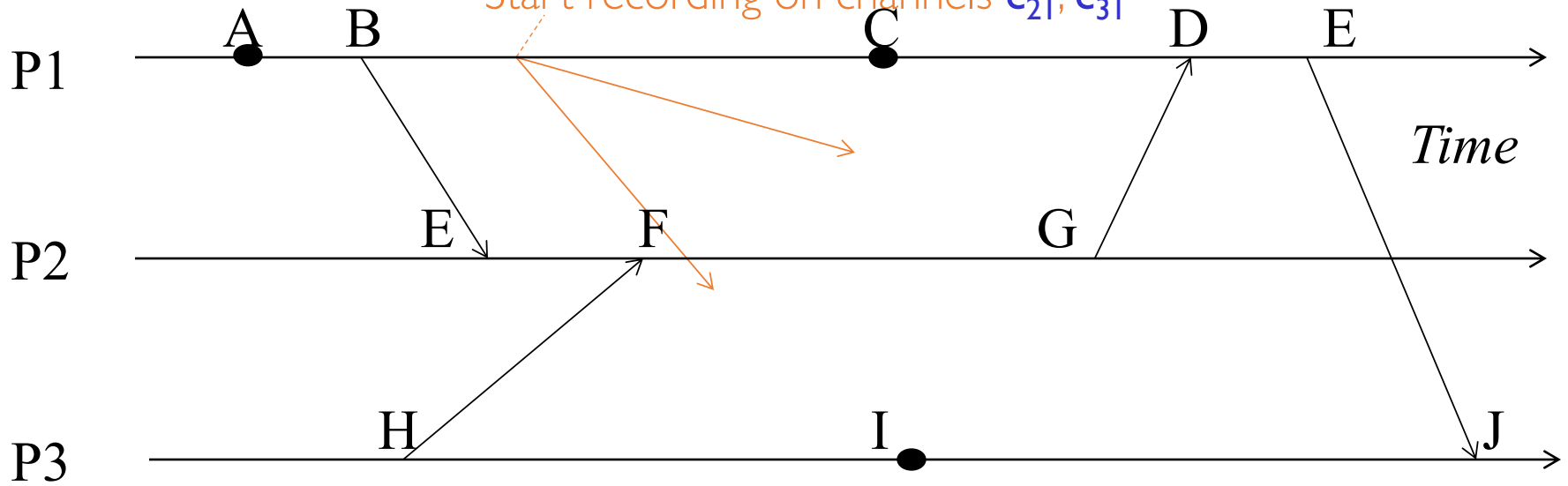
Example



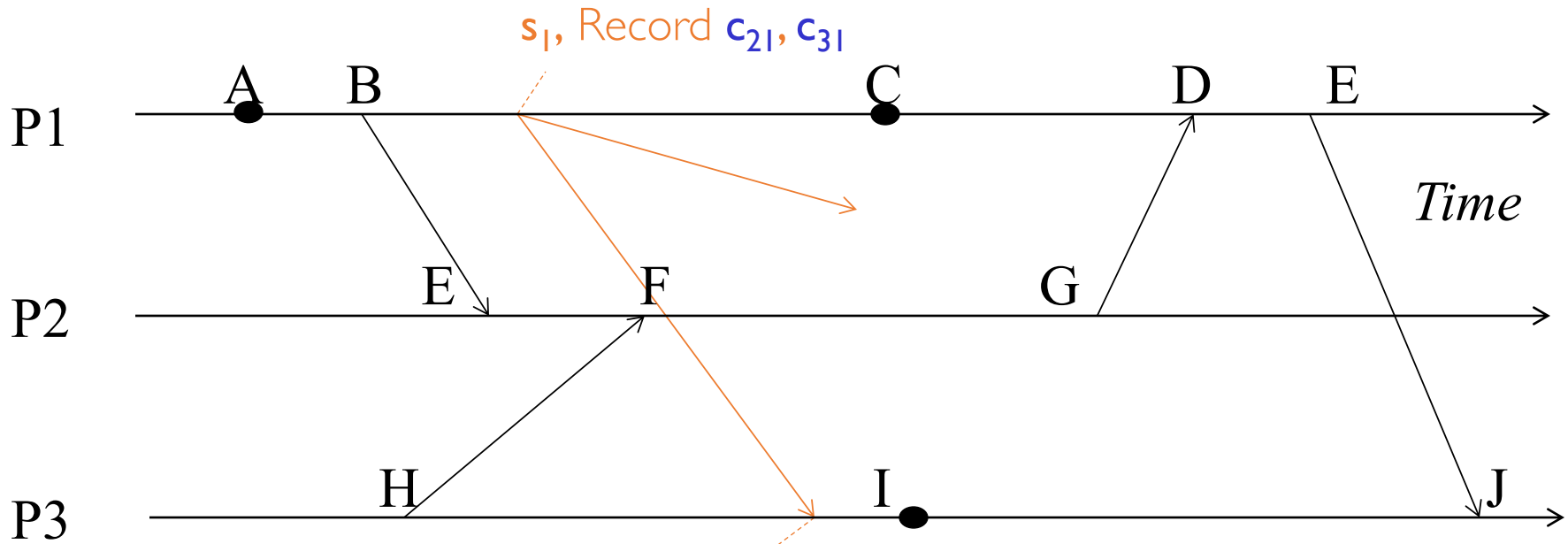
Example

p_1 is initiator:

- Record local state s_1 ,
- Send out markers
- Start recording on channels c_{21}, c_{31}

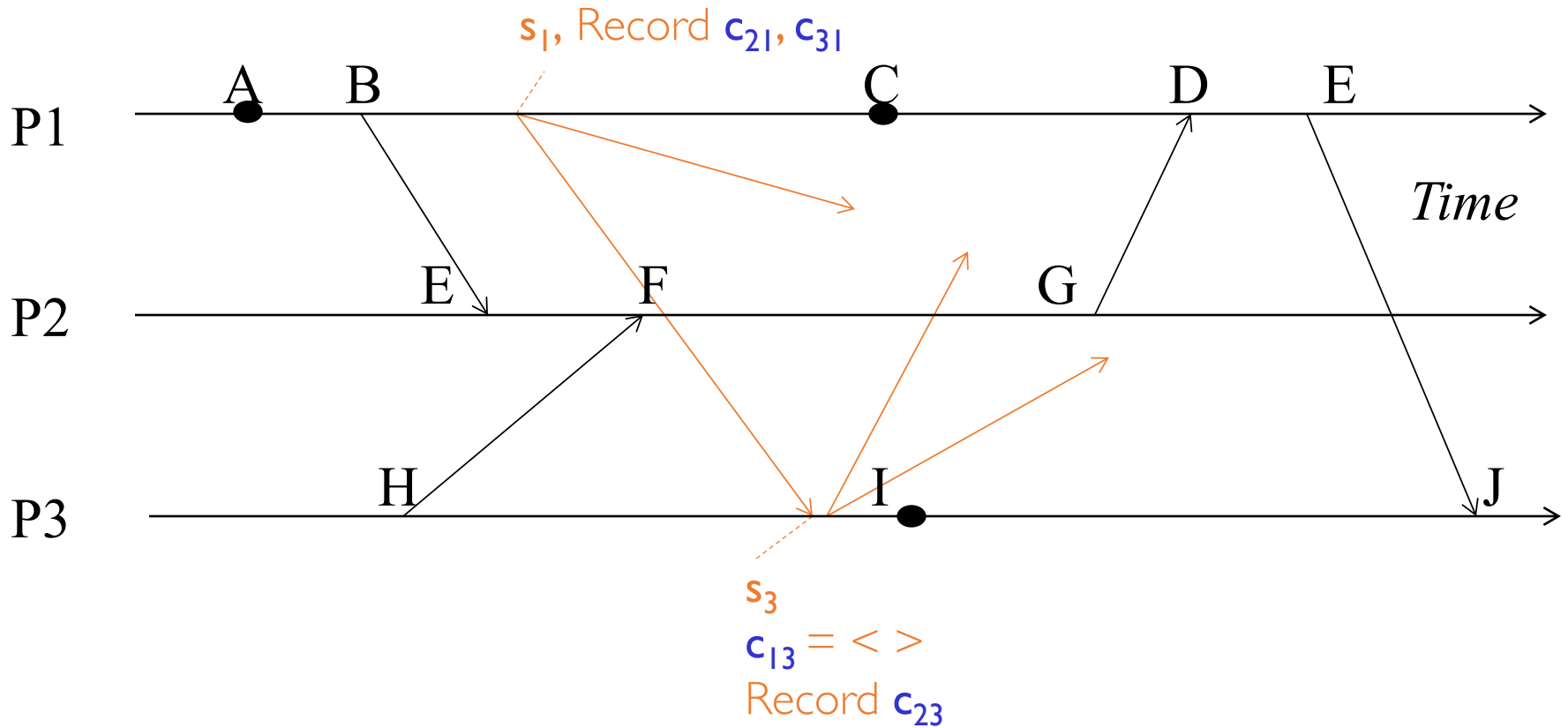


Example

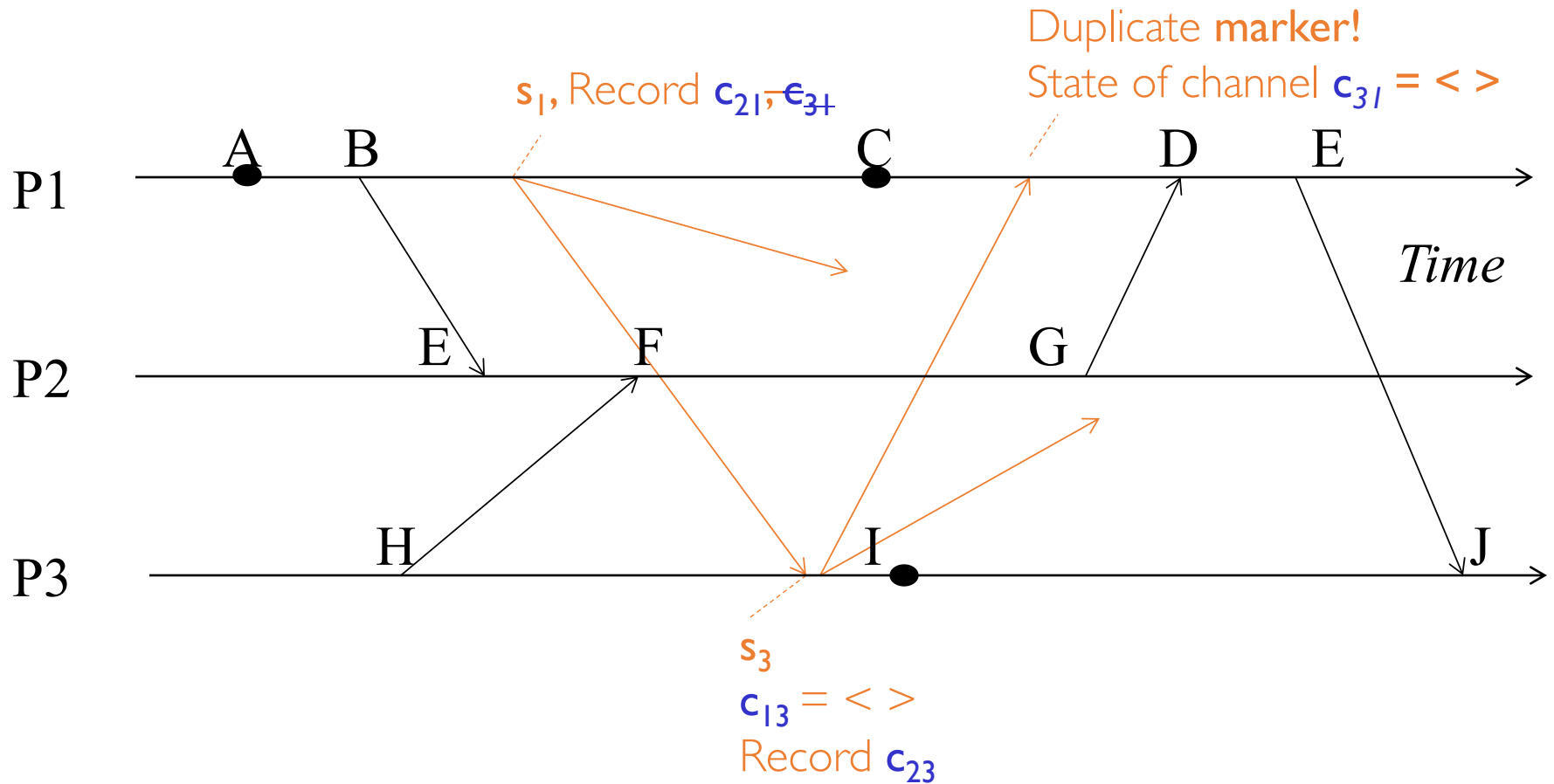


- First **marker!**
- Record own state as s_3
- Mark c_{13} state as empty
- Start recording on other incoming c_{23}
- Send out markers

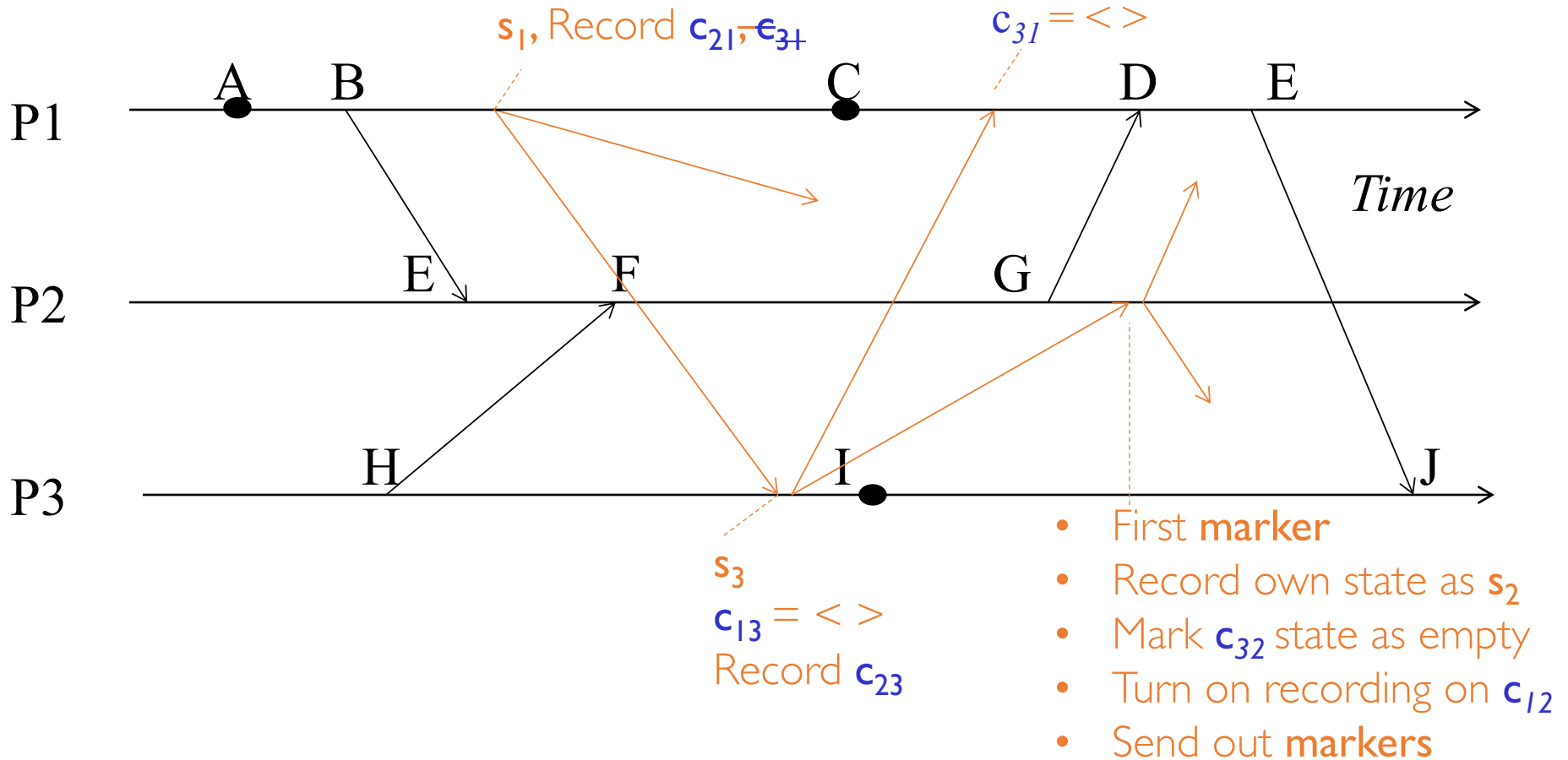
Example



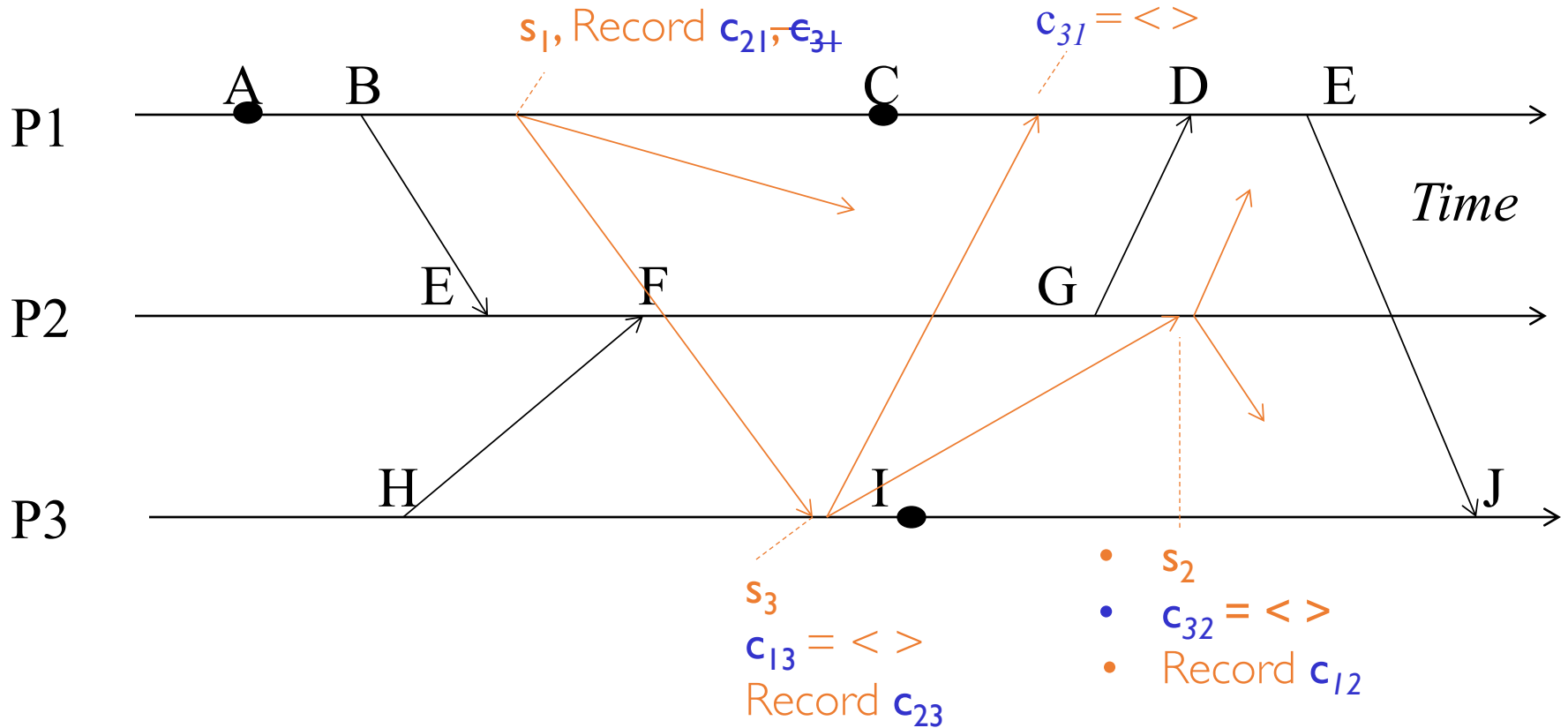
Example



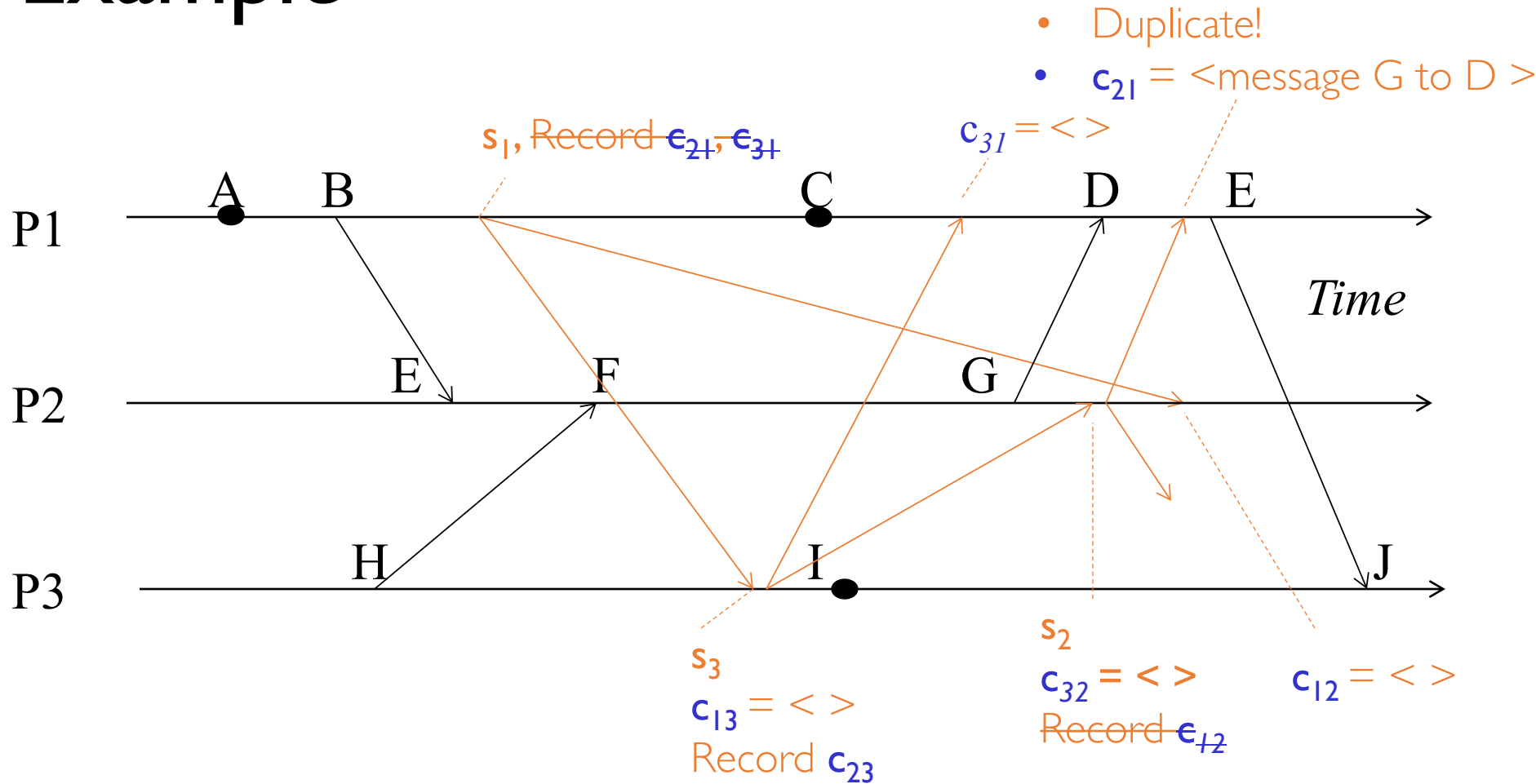
Example



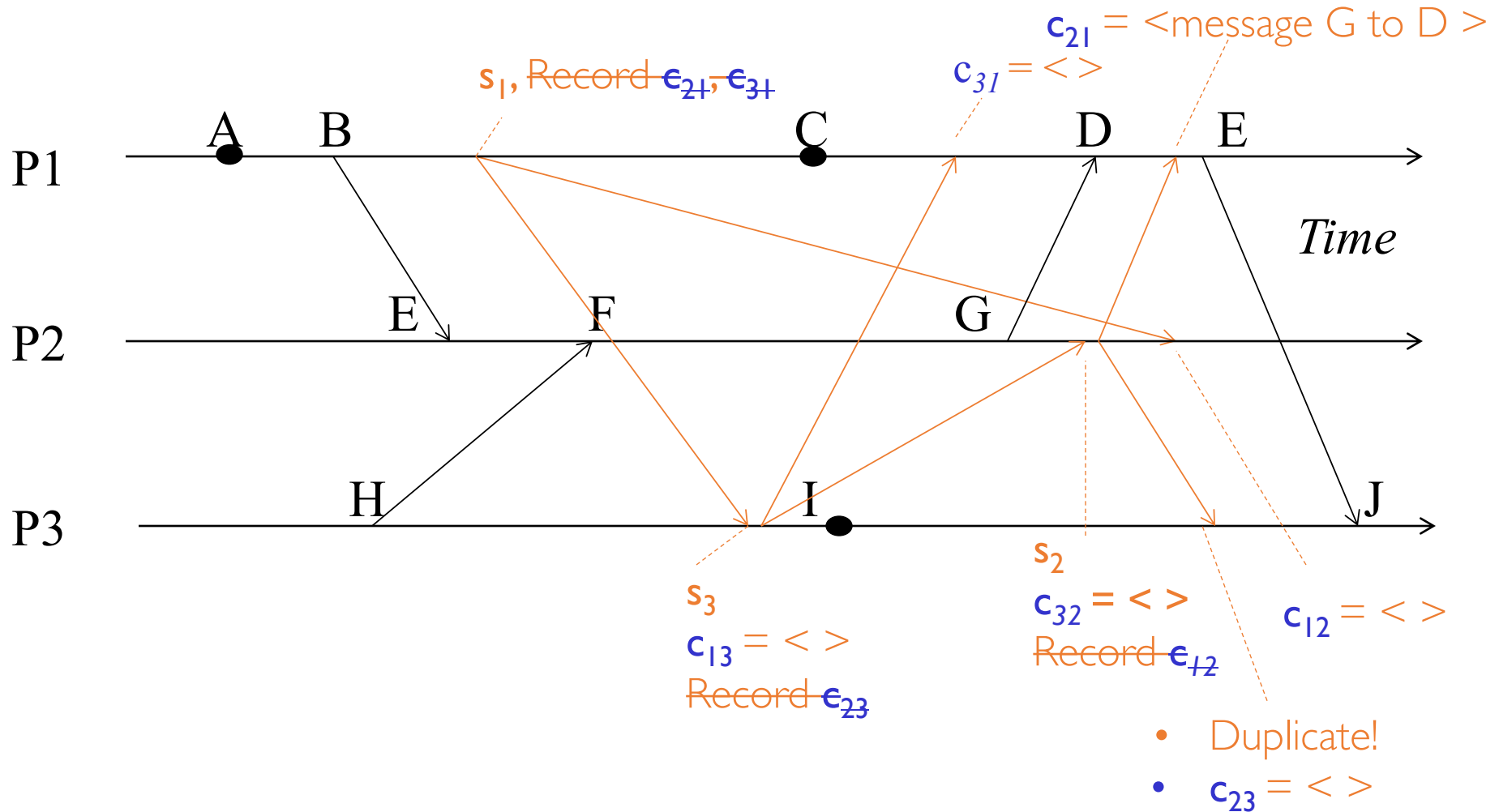
Example



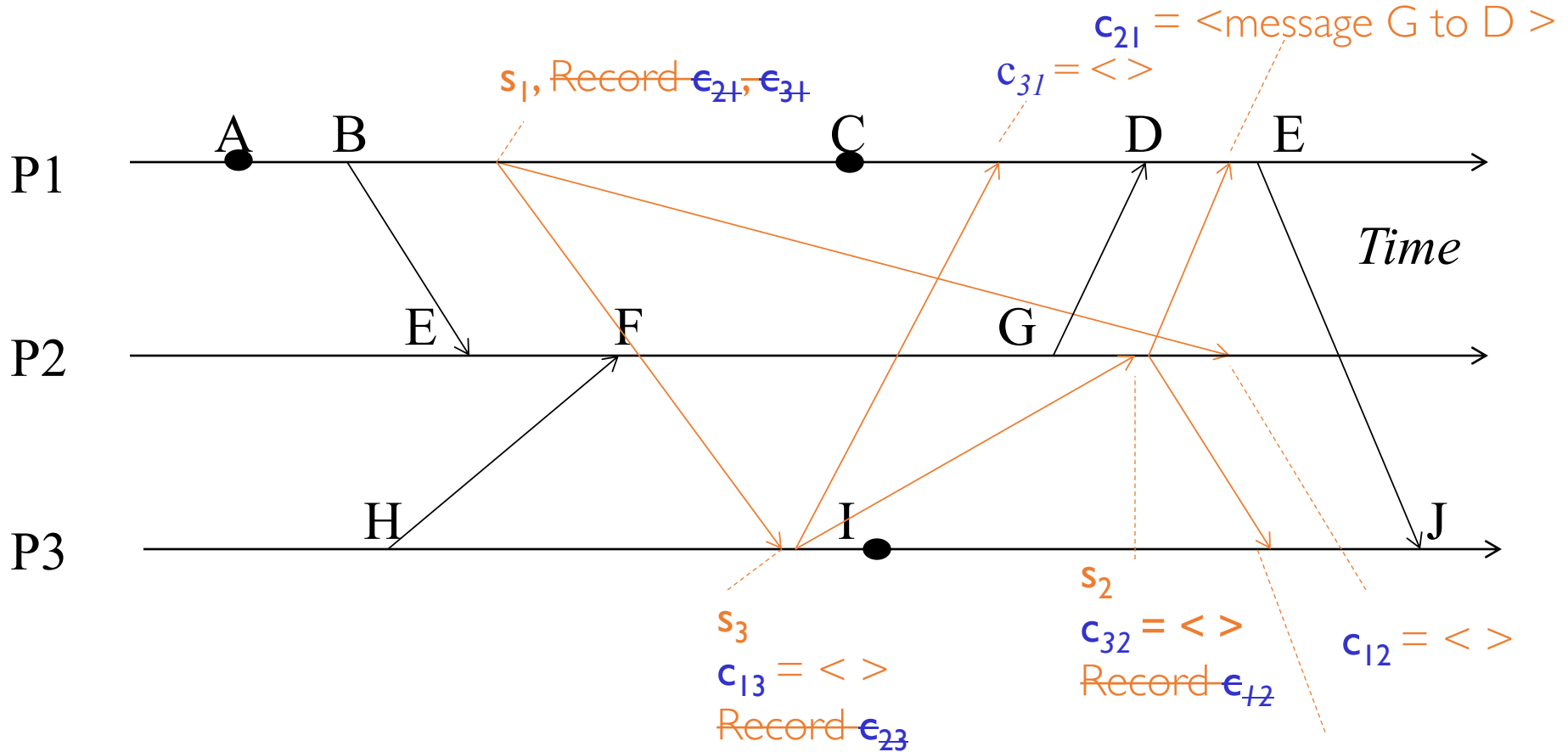
Example



Example



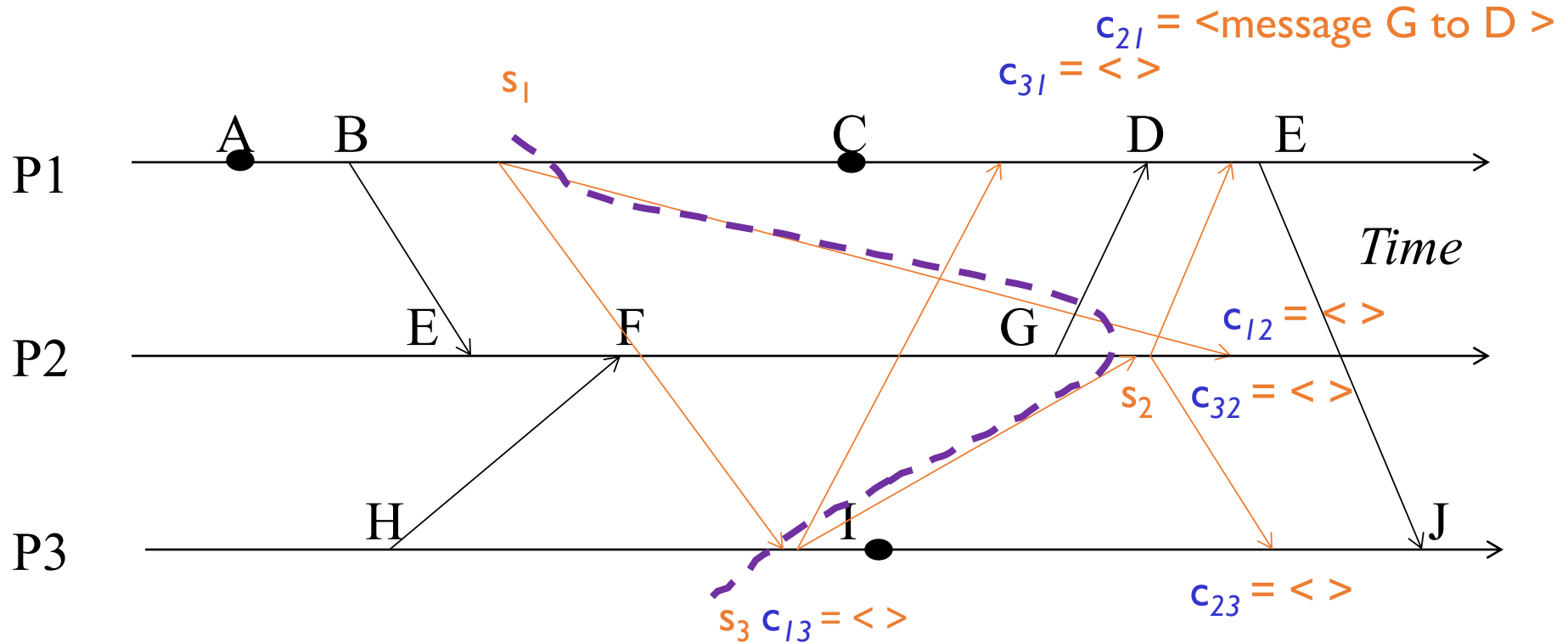
Example



Algorithm has terminated!

- Duplicate!
- $c_{23} = \langle \rangle$

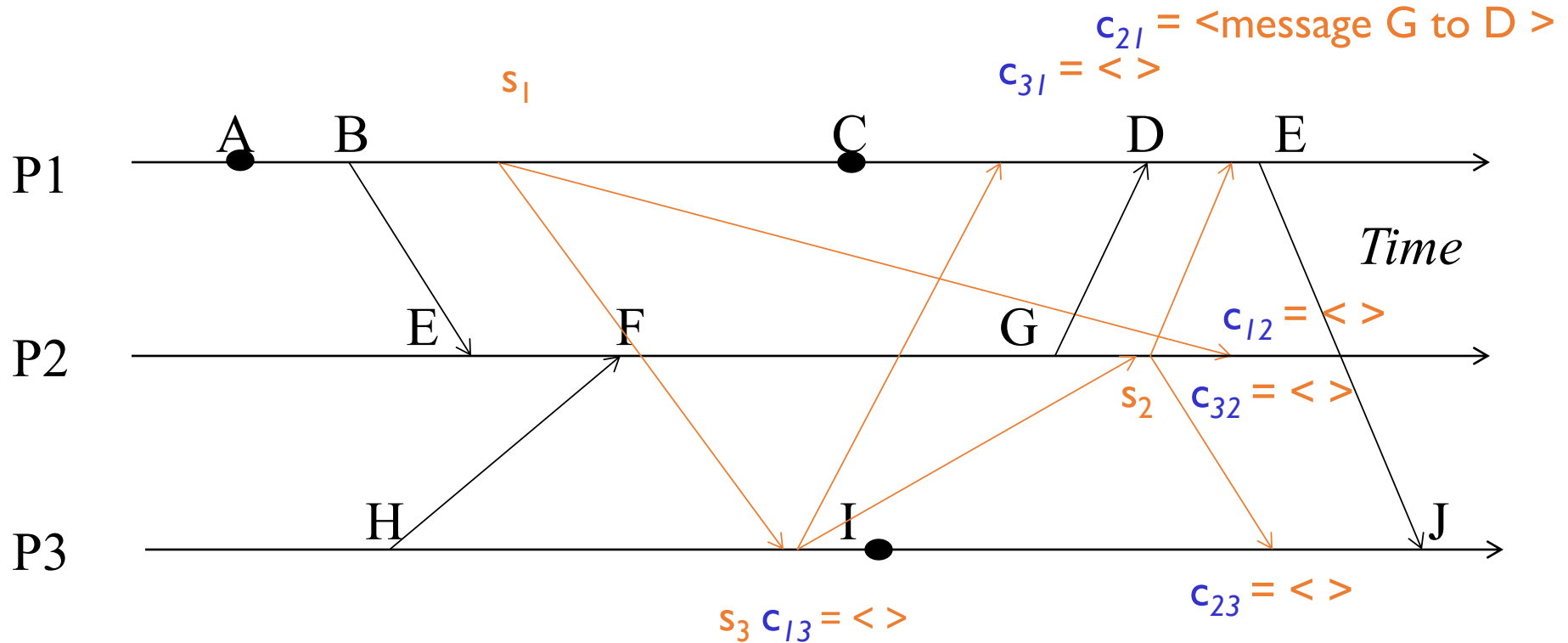
Example



Frontier for the resulting cut:
 $\{B, G, H\}$

Channel state for the cut:
Only c_{21} has a pending message.

Example



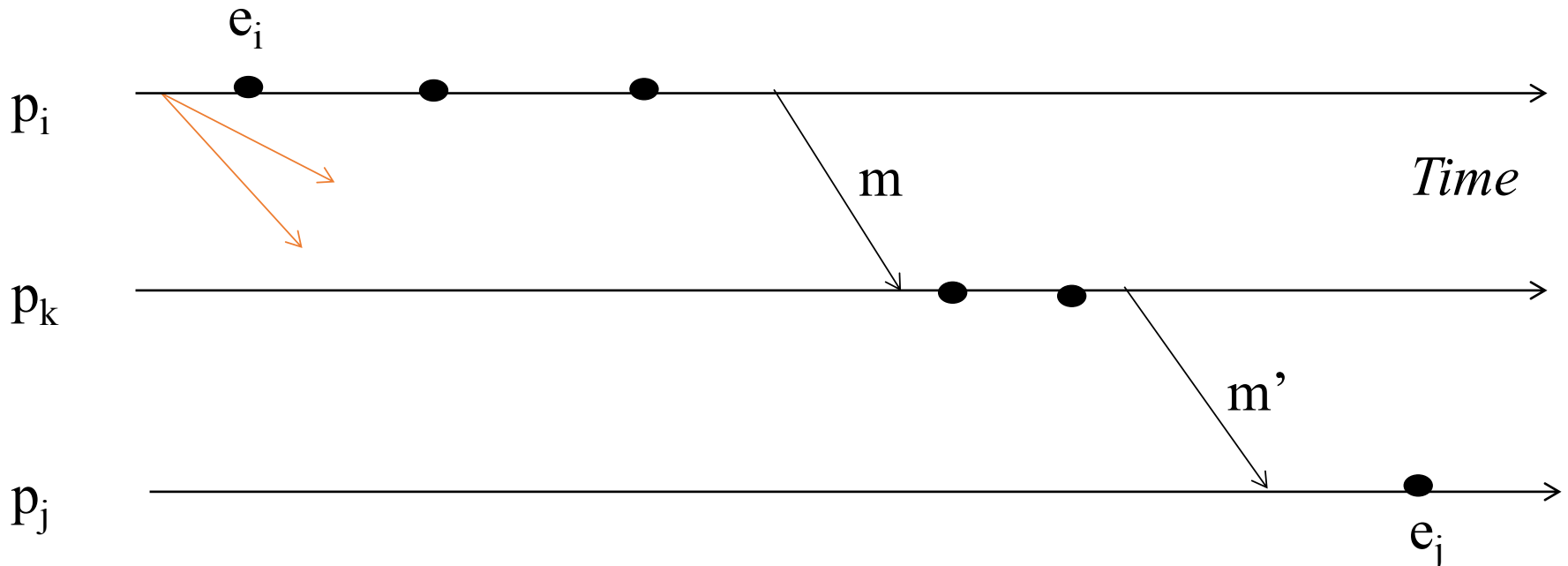
Global snapshots pieces can be collected at a central location.

Chandy-Lamport Algorithm: Properties

- Any run of the Chandy-Lamport Global Snapshot algorithm creates a consistent cut.
- Let e_i and e_j be events occurring at p_i and p_j , respectively such that
 - $e_i \rightarrow e_j$ (e_i happens before e_j)
- The snapshot algorithm ensures that
 - if e_j is in the cut then e_i is also in the cut.
- That is: if $e_j \rightarrow \langle p_j \text{ records its state} \rangle$, then it must be true that $e_i \rightarrow \langle p_i \text{ records its state} \rangle$.

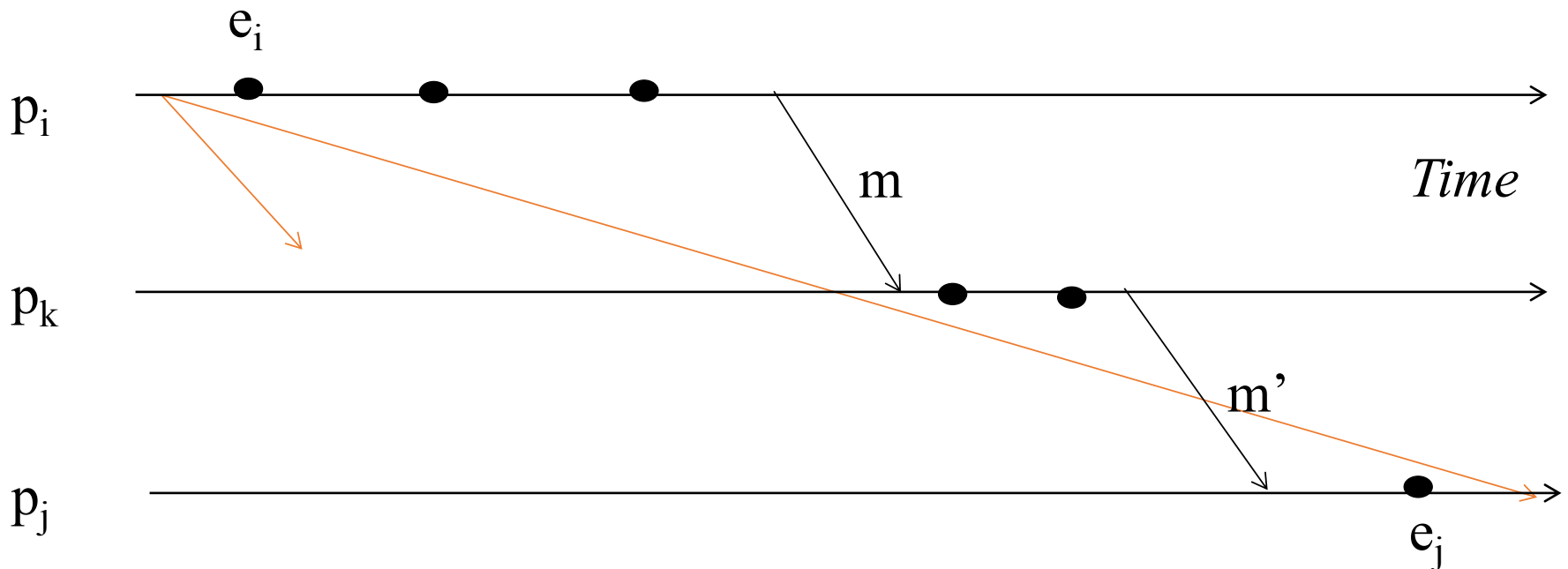
Chandy-Lamport Algorithm: Properties

- If $e_j \rightarrow \langle p_j \text{ records its state} \rangle$, then it must be true that $e_i \rightarrow \langle p_i \text{ records its state} \rangle$.
- By contradiction, suppose $e_j \rightarrow \langle p_j \text{ records its state} \rangle$, and $\langle p_i \text{ records its state} \rangle \rightarrow e_i$.



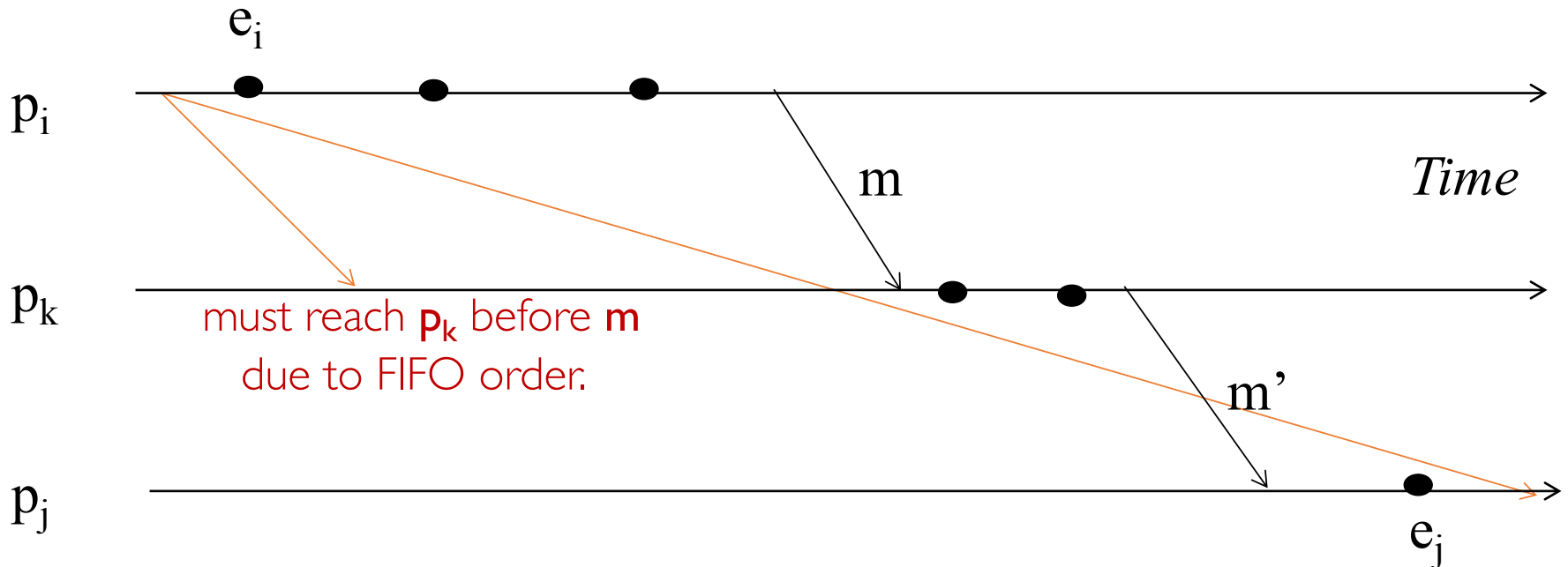
Chandy-Lamport Algorithm: Properties

- If $e_j \rightarrow \langle p_j \text{ records its state} \rangle$, then it must be true that $e_i \rightarrow \langle p_i \text{ records its state} \rangle$.
- By contradiction, suppose $e_j \rightarrow \langle p_j \text{ records its state} \rangle$, and $\langle p_i \text{ records its state} \rangle \rightarrow e_i$.



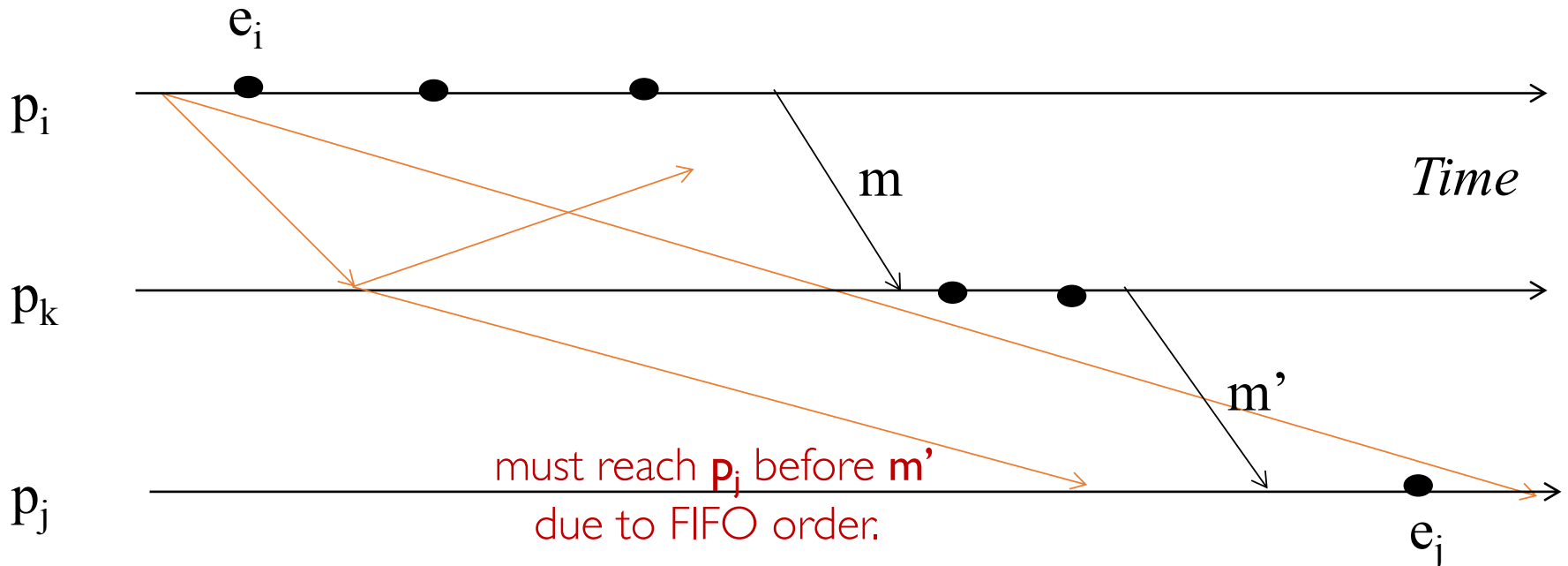
Chandy-Lamport Algorithm: Properties

- If $e_j \rightarrow \langle p_j \text{ records its state} \rangle$, then it must be true that $e_i \rightarrow \langle p_i \text{ records its state} \rangle$.
- By contradiction, suppose $e_j \rightarrow \langle p_j \text{ records its state} \rangle$, and $\langle p_i \text{ records its state} \rangle \rightarrow e_i$.



Chandy-Lamport Algorithm: Properties

- If $e_j \rightarrow \langle p_j \text{ records its state} \rangle$, then it must be true that $e_i \rightarrow \langle p_i \text{ records its state} \rangle$.
- By contradiction, suppose $e_j \rightarrow \langle p_j \text{ records its state} \rangle$, and $\langle p_i \text{ records its state} \rangle \rightarrow e_i$.



Chandy-Lamport Algorithm: Properties

- If $e_j \rightarrow \langle p_j \text{ records its state} \rangle$, then it must be true that $e_i \rightarrow \langle p_i \text{ records its state} \rangle$.
- By contradiction, suppose $e_j \rightarrow \langle p_j \text{ records its state} \rangle$, and $\langle p_i \text{ records its state} \rangle \rightarrow e_i$.
- Consider the path of app messages (through other processes) that go from e_i to e_j .
- Due to FIFO ordering, markers on each link in above path will precede regular app messages.
- Thus, since $\langle p_i \text{ records its state} \rangle \rightarrow e_i$, it must be true that p_j received a marker before e_j .
- Thus e_j is not in the cut \Rightarrow contradiction.

Global Snapshot Summary

- The ability to calculate global snapshots in a distributed system is very important.
- But don't want to interrupt running distributed application.
- Chandy-Lamport algorithm calculates global snapshot.
- Obeys causality (creates a consistent cut).
- Can be used to detect global properties.
 - Safety vs. Liveness (next class)