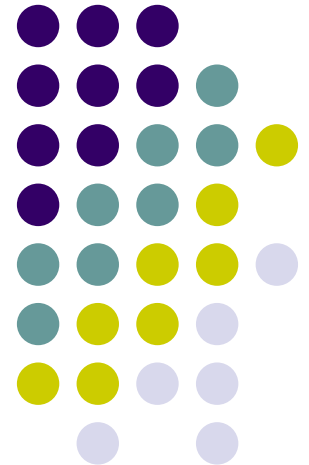


Distributed Algorithms for Building Reliable Systems

Seif Haridi

Causal Broadcast



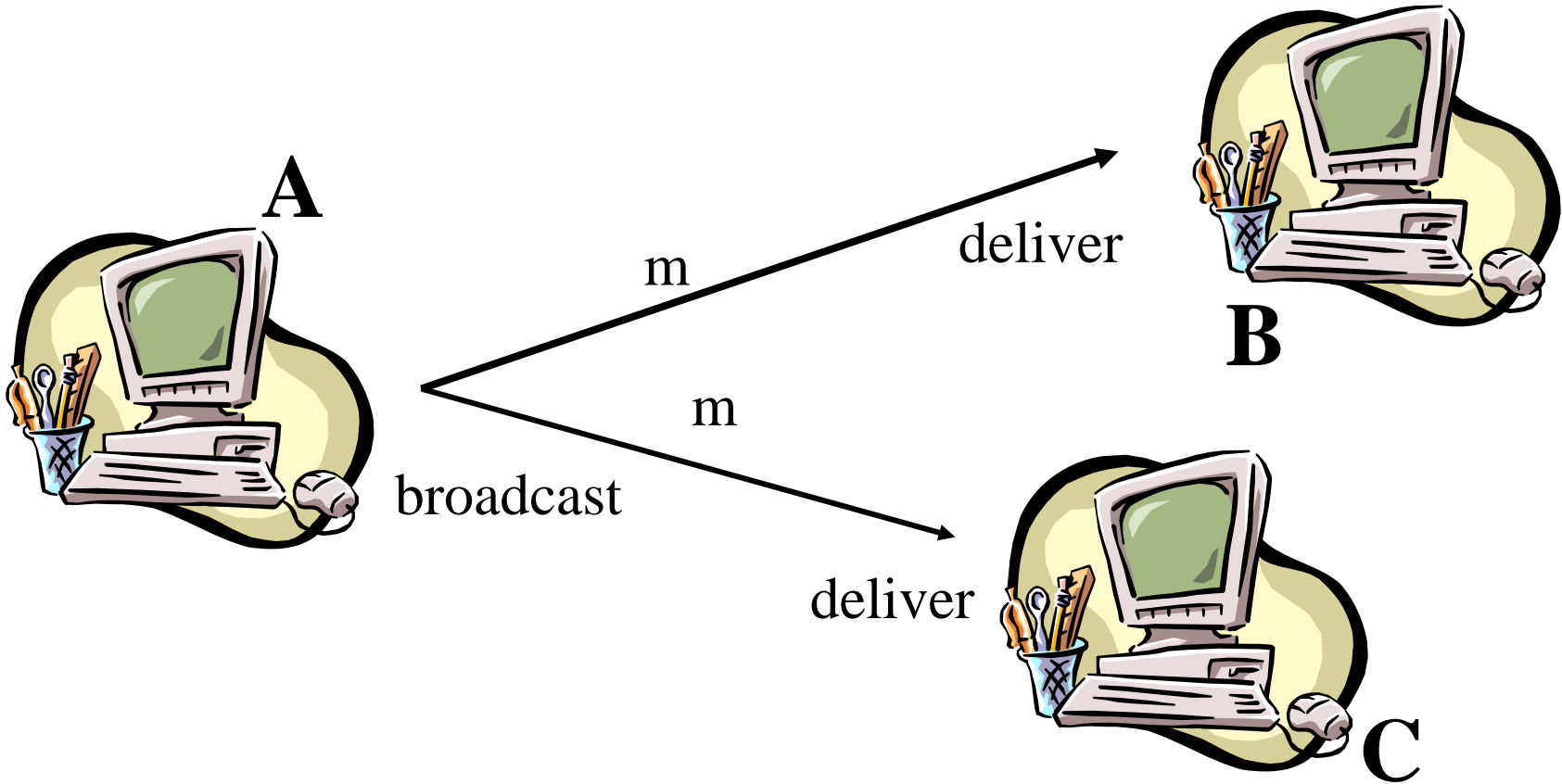


Overview

- **Intuitions: why causal broadcast?**
- **Specifications of *causal broadcast***
- **Algorithms:**
 - A *non-blocking* algorithm using the *past* and
 - A *blocking* algorithm using *vector clocks*



Broadcast



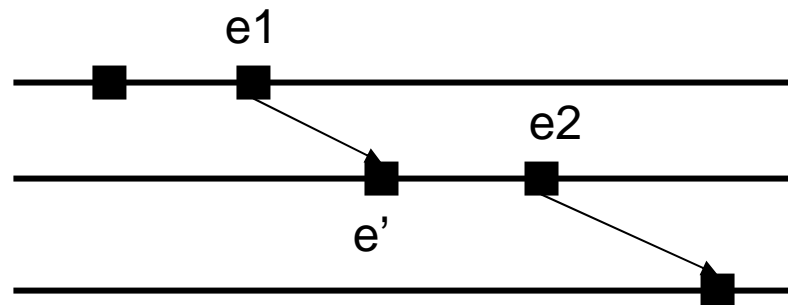
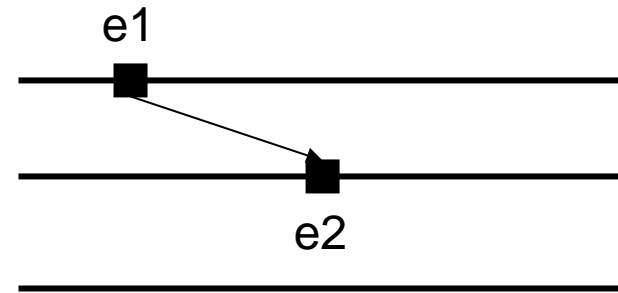
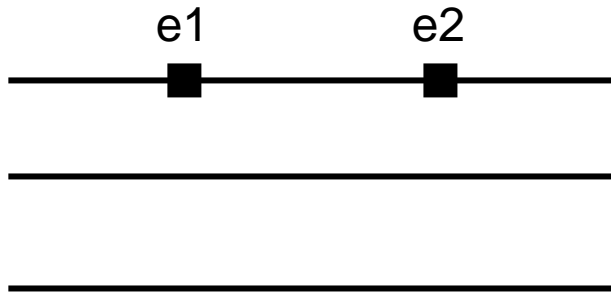
Cause-effect relations in message passing systems



- An event e_1 may potentially have caused another event e_2 if the following relation, called, *happens-before* and denote by $e_1 \rightarrow e_2$ holds
 - e_1 and e_2 occurs at the same process p , and e_1 occurs before e_2
 - e_1 is the transmission of a message m at process p and e_2 is the reception of the same message at process q
 - There exist some event e' such that $e_1 \rightarrow e'$ and $e' \rightarrow e_2$



Happens-before relation





Intuitions (1)

- So far, we did not consider ordering among messages; In particular, we considered messages to be independent
- Two messages from the same process might not be delivered in the order they were broadcast
- A message m_1 that causes a message m_2 might be delivered by some process after m_2

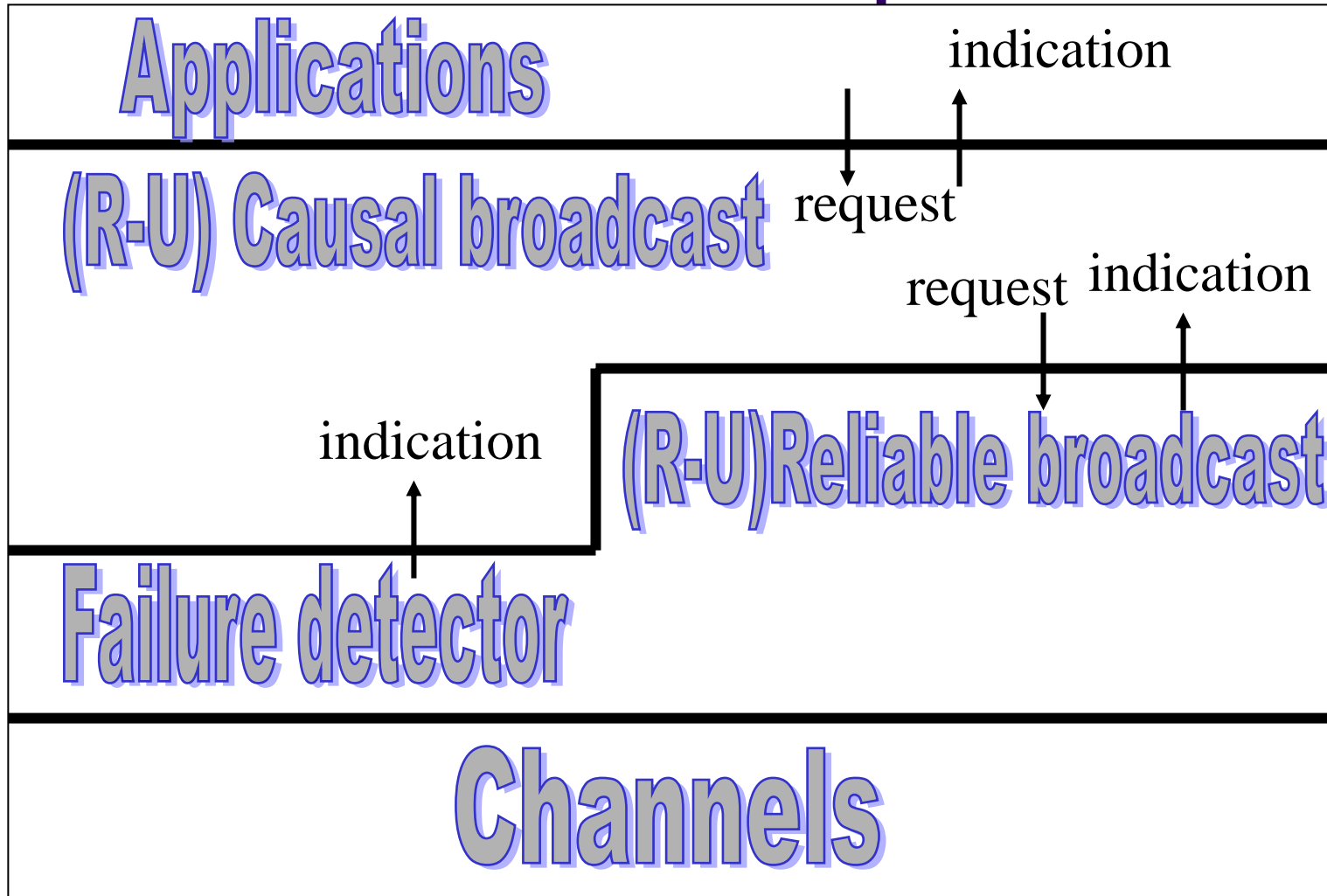
Intuitions (2)



- ☞ Consider a replicated distributed bulletin board system that manages two types of news
 - ☞ Proposals and comments
 - ☞ Every new event that is displayed contains a reference to the event that *caused* it, e.g., a comment includes a reference to the actual proposal
- ☞ Even uniform reliable broadcast does not guarantee such a dependency of delivery
- ☞ Causal broadcast alleviates the need for the application to deal with such dependencies



Modules of a process



Overview



- **Intuitions: why causal broadcast?**
- **Specifications of *causal broadcast***
- **Algorithms:**
 - A *non-blocking* algorithm using the *past* and
 - A *blocking* algorithm using *vector clocks*

Causal broadcast



• *Events*

- Request: $\langle \text{coBroadcast } m \rangle$
- Indication: $\langle \text{coDeliver src, } m \rangle$

• *Property:*

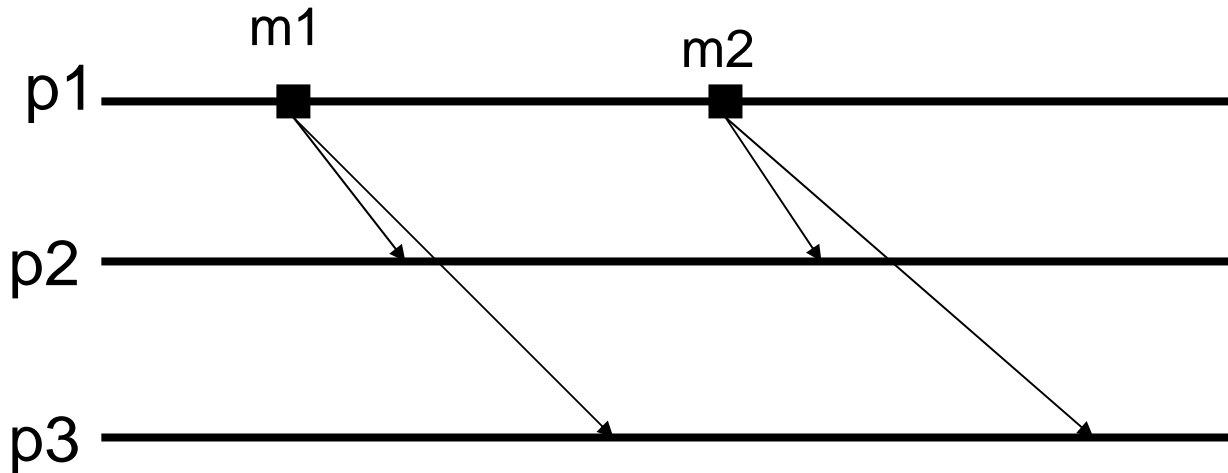
- *Causal Order (CO)*

Causality



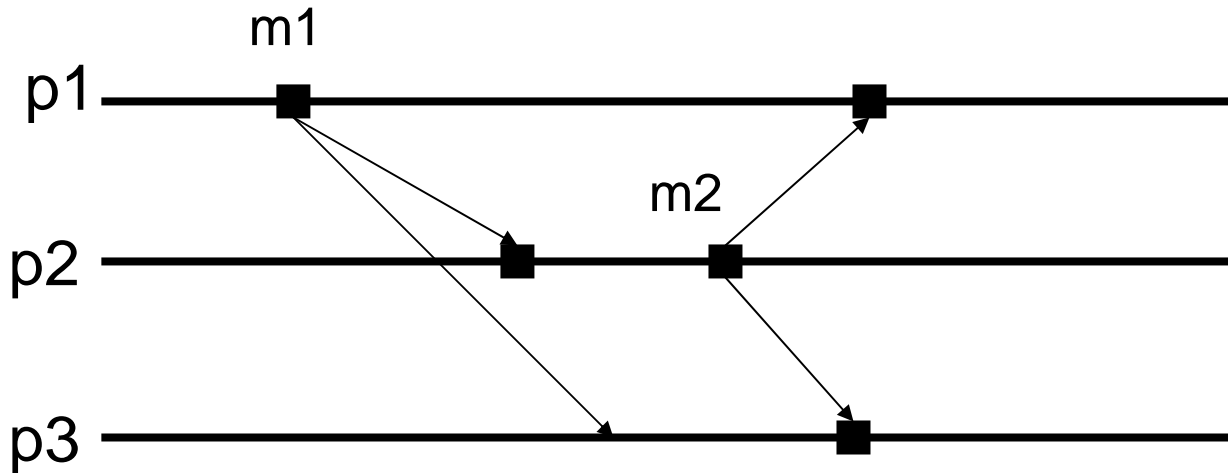
- Let m_1 and m_2 be any two messages:
 $m_1 \rightarrow m_2$ (m_1 causally precedes m_2) iff
 - C1 (FIFO order)**. Some process p_i broadcasts m_1 before broadcasting m_2
 - C2 (Local order)**. Some process p_i delivers m_1 and then broadcasts m_2
 - C3 (Transitivity)**. There is a message m' such that $m_1 \rightarrow m'$ and $m' \rightarrow m_2$

Causality



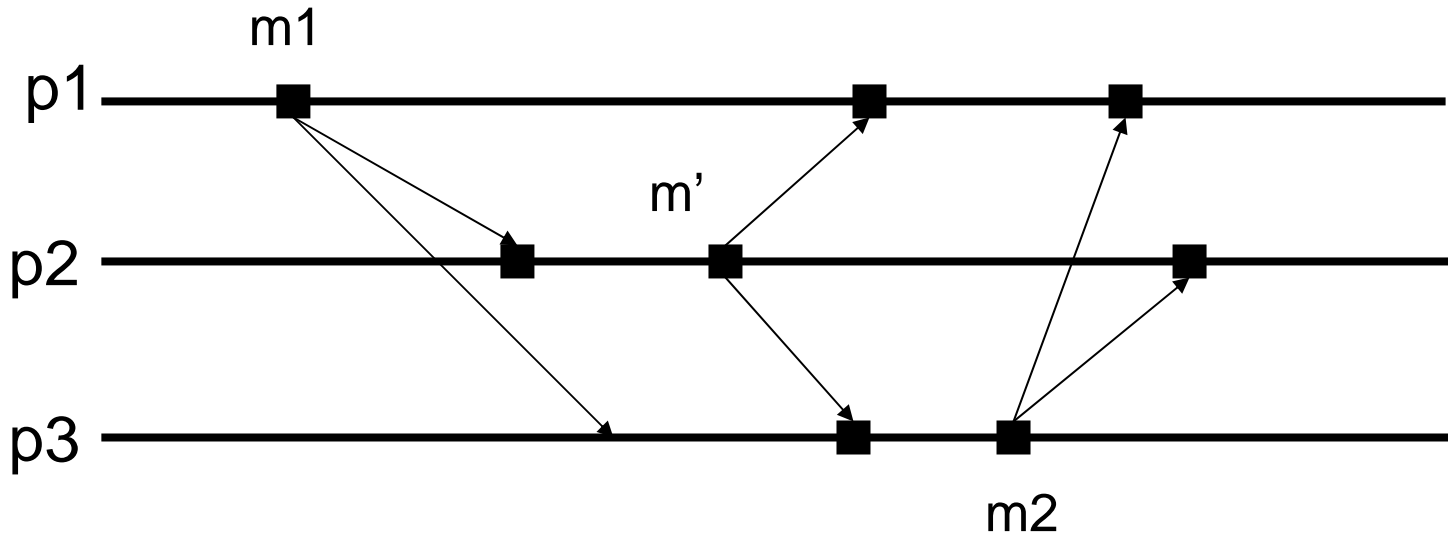
 **C1 (FIFO order).**

Causality



 **C2 (Local order).**

Causality



❁ **C3 (Transitivity).**

Causal broadcast



• *Events*

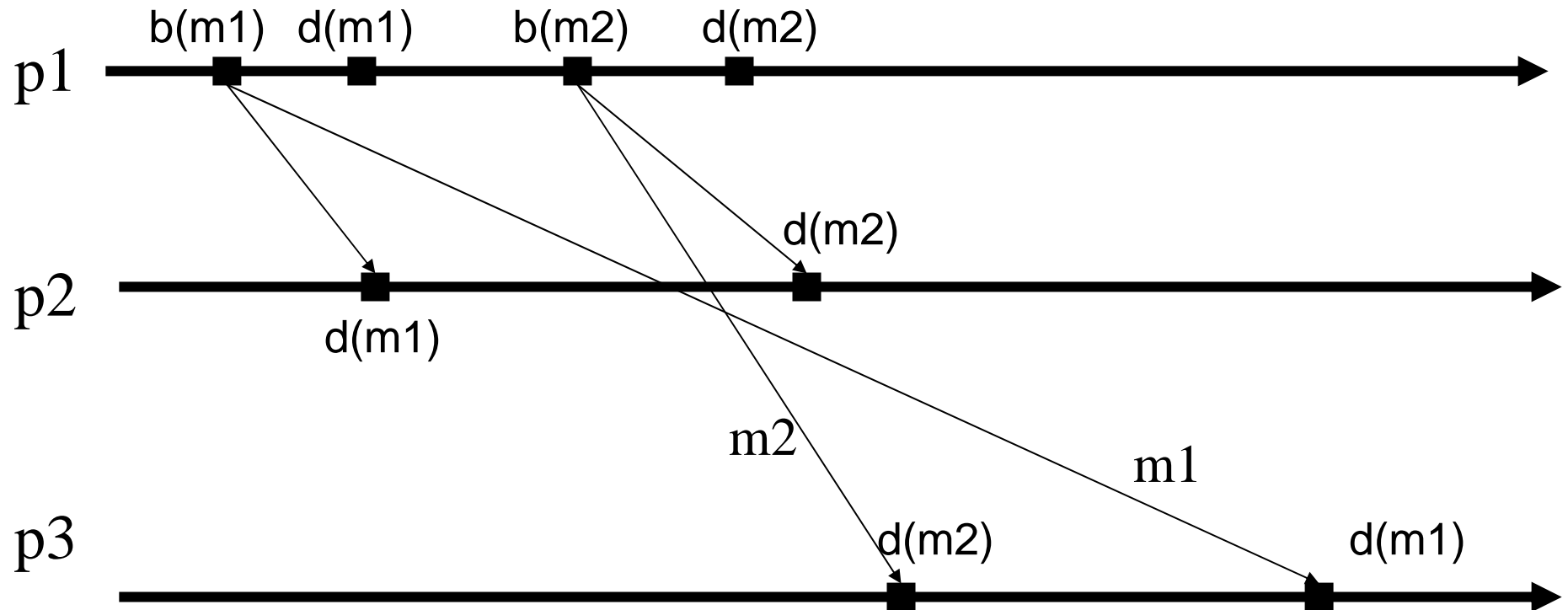
- Request: $\langle \text{coBroadcast } m \rangle$
- Indication: $\langle \text{coDeliver src, } m \rangle$

• *Property:*

- **CO:** If any process p_i delivers a message m_2 , then p_i must have delivered every message m_1 such that $m_1 \rightarrow m_2$

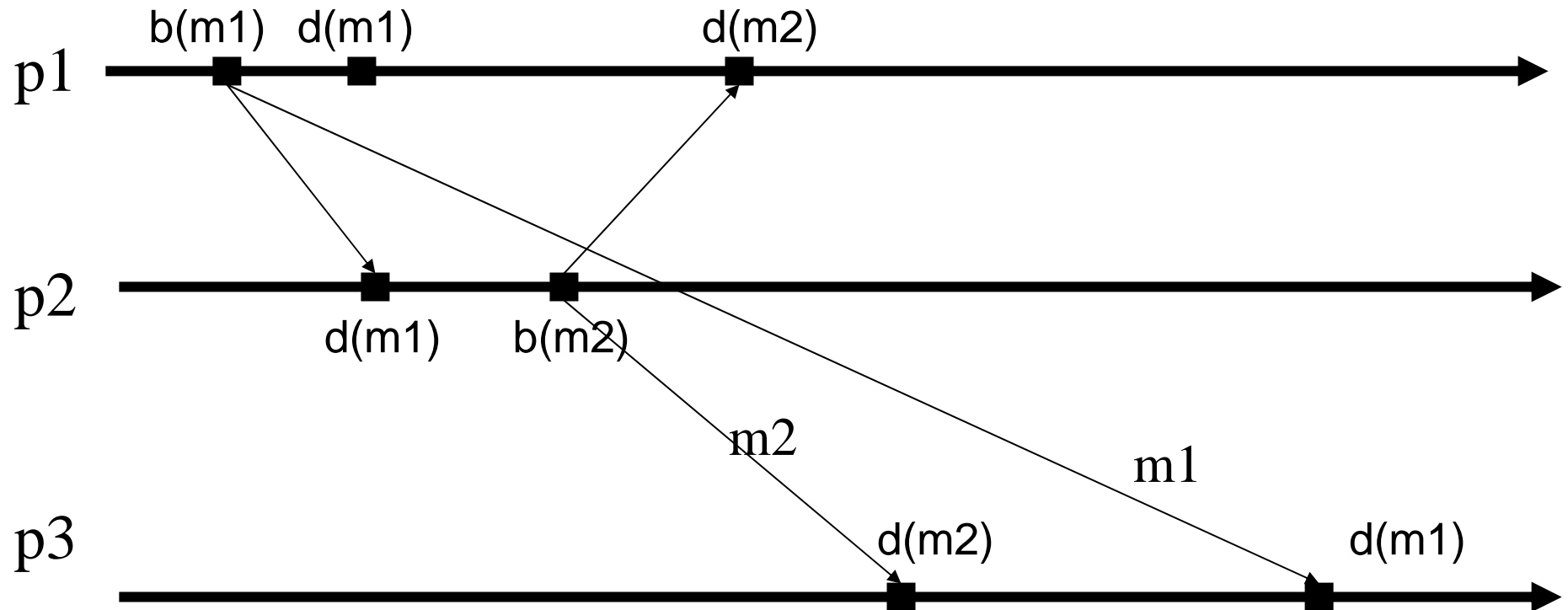
Causality ?

d(elivery) b(broadcast)



Causality ?

d(elivery) b(broadcast)



Reliable causal broadcast (rcb)



• *Events*

• Request: $\langle \text{rcoBroadcast } m \rangle$

• Indication: $\langle \text{rcoDeliver src, } m \rangle$

• *Properties:*

• *RB1, RB2, RB3, RB4 +*

• *CO*

Uniform causal broadcast (ucb)



• *Events*

- Request: $\langle \text{ucoBroadcast}, m \rangle$
- Indication: $\langle \text{ucoDeliver}, \text{src}, m \rangle$

• *Properties:*

- *URB1, URB2, URB3, URB4 +*
- *CO*

Reliable broadcast (rb)



☛ *Properties*

- ☛ ***RB1. Validity.*** If a correct process p_i broadcasts a message m , then p_i eventually delivers m .
- ☛ ***RB2 = BEB2.***
- ☛ ***RB3 = BEB3.***
- ☛ ***RB4. Agreement:*** For any message m , if a **correct process delivers** m , then every correct process delivers m

Uniform broadcast (urb)



☛ *Properties*

☛ ***URB1 = BEB1.***

☛ ***URB2 = BEB2.***

☛ ***URB3 = BEB3.***

☛ ***URB4. Uniform Agreement:*** For any message m , if a **process delivers** m , then every correct process delivers m



Overview

- ☞ Intuitions: why causal broadcast?
- ☞ Specifications of *causal broadcast*
- ☞ Algorithms:
 - ☞ A *non-blocking* algorithm using the *past* and
 - ☞ A *blocking* algorithm using *vector clocks*

Algorithms



- ☛ We present **reliable causal broadcast** algorithms using **reliable broadcast** underlying abstractions
- ☛ We obtain **uniform causal broadcast** algorithms by using instead an underlying **uniform reliable broadcast**

Algorithm1 Fail-Silent No-Waiting Causal Broadcast



- Whenever a process `rbDelivers` m , it `coDelivers` m (if not already `coDelivered`)
- Each message m carries an ordered list of messages past_m that causally precede m
- Before m is `coDelivered`, past_m is inspected
- Messages in past_m that are not `coDelivered` are `coDelivered` before m

Algorithm 1

Fail-silent No Waiting Causal Broadcast



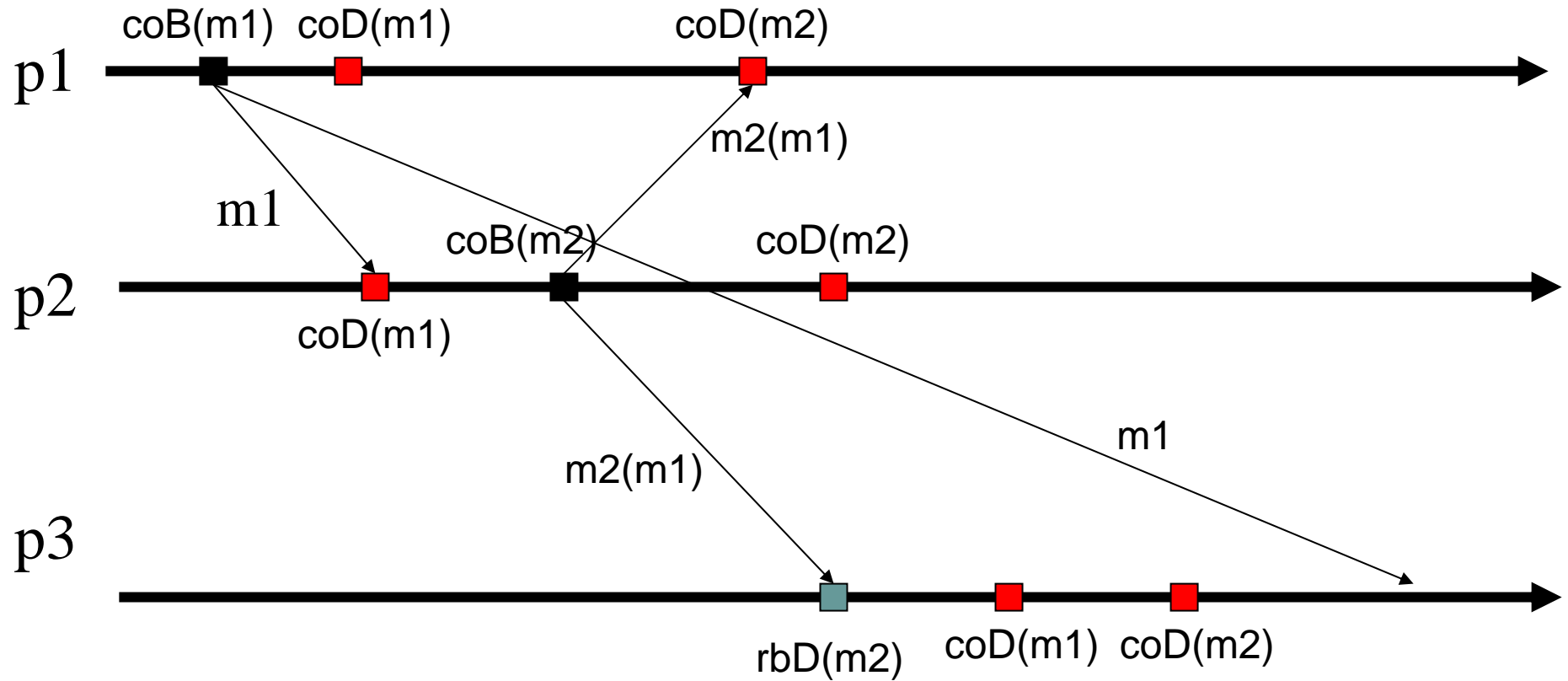
- ☛ **Implements:** ReliableCausalOrderBroadcast (rco)
- ☛ **Uses:** ReliableBroadcast (rb).
- ☛ **upon event** $\langle \text{Init} \rangle$ **do**
 - ☛ $\text{delivered} := \emptyset; \text{past} := \text{nil}$
- ☛ **upon event** $\langle \text{rcoBroadcast } m \rangle$ **do**
 - ☛ **trigger** $\langle \text{rbBroadcast} (\text{DATA}, \text{past}, m) \rangle$
 - ☛ **past := past ++ (pi, m)**

Algorithm 1 (cont'd)

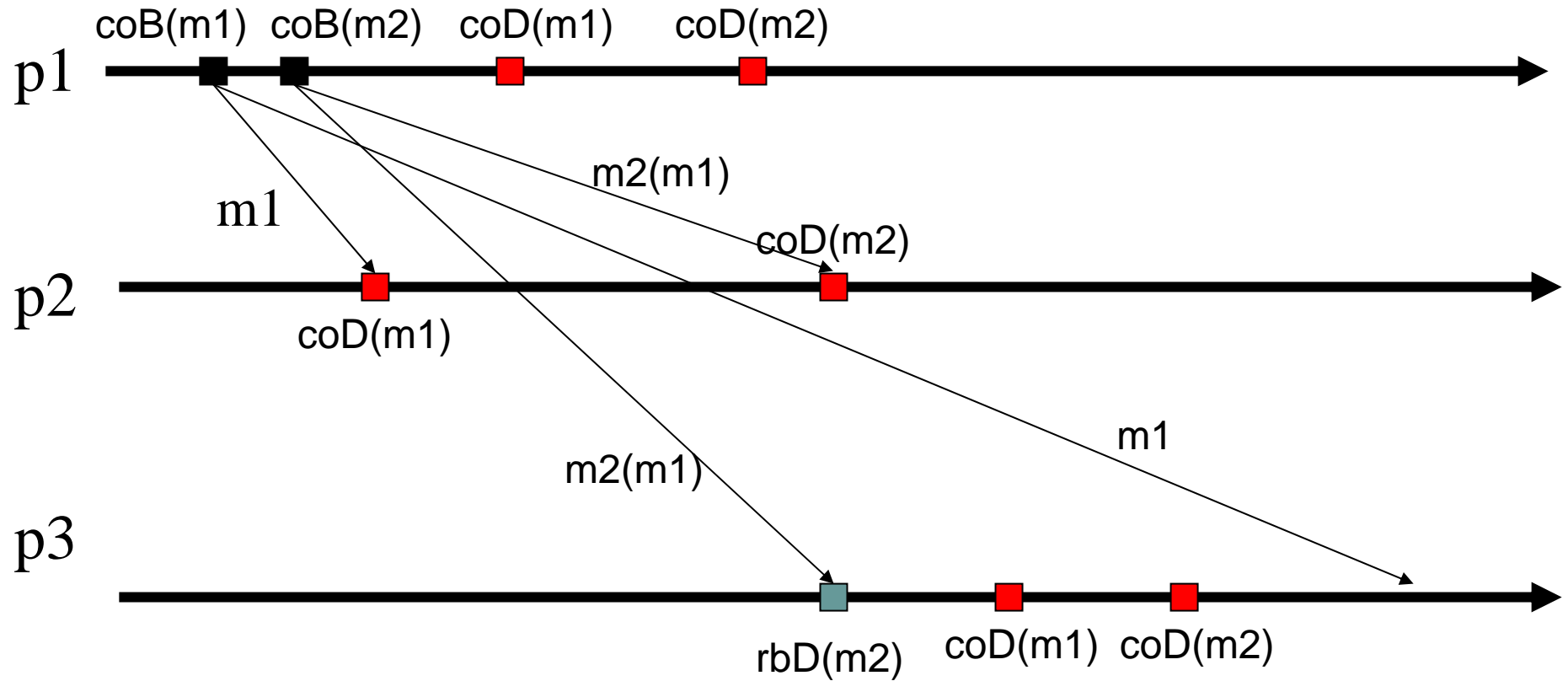


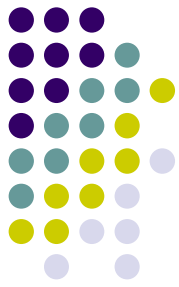
- ☞ **upon event** $\langle \text{rbDeliver } p_i, (\text{DATA}, \text{past}_m, m) \rangle$ **do**
 - ☞ **if** $m \notin \text{delivered}$ **then**
 - ☞ **forall** $(s_n, n) \in \text{past}_m$ **do** // in ascending order
 - ☞ **if** $n \notin \text{delivered}$ **then**
 - ☞ **trigger** $\langle \text{rcoDeliver } s_n, n \rangle$
 - ☞ $\text{delivered} := \text{delivered} \cup \{n\}$
 - ☞ $\text{past} := \text{past} ++ (s_n, n)$ // sequence (list) append
 - ☞ **trigger** $\langle \text{rcoDeliver } p_i, m \rangle$
 - ☞ $\text{delivered} := \text{delivered} \cup \{m\}$
 - ☞ $\text{past} := \text{past} ++ (p_i, m)$ // sequence (list) append

Algorithm 1



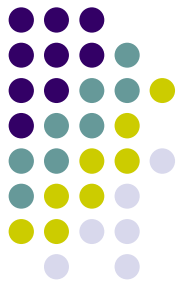
Algorithm 1





Algorithm 1 comments

- Consider the sequence of events at p_1
 - $\langle \text{rcoB } m_1 \rangle_{p_1}, \langle \text{rcoB } m_2 \rangle_{p_1},$
 $\langle \text{rbD } p_1, (\text{DATA}, \text{past}_{m_2}, m_2) \rangle_{p_1},$
 $\langle \text{rbD } p_1, (\text{DATA}, \text{past}_{m_1}, m_1) \rangle_{p_1}$
 - What is the order of co-delivery of m_1 and m_2 at p_1
- We know that reliable broadcast (rb) has the no-duplicate property why in this case we need the test
if $m \notin$ delivered then ... ?

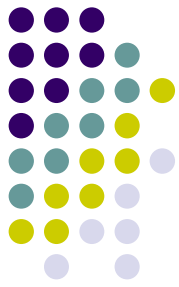


Uniformity

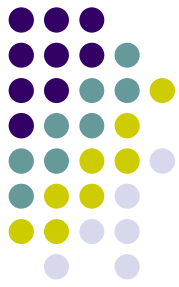
- Algorithm 1 ensures causal reliable broadcast
- If we replace reliable broadcast with uniform reliable broadcast, Algorithm 1 would ensure uniform causal broadcast

Algorithm 1'

Garbage collection of past

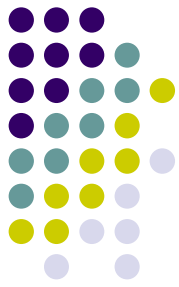


- Assume a fail-stop model, and use a perfect failure detector
- When a process $rbDelivers$ m , it $rbBroadcasts$ an ACK_m to all processes
- When ACK_m is $rbDelivered$ from all correct processes m is purged from *past*



Algorithm 1' (gc)

- **Implements:** GarbageCollection (+ Algo 1)
- **Extends:** Algorithm 1
- **Uses:**
 - ReliableBroadcast (rb)
 - PerfectFailureDetector(P)
- **upon event** $\langle \text{Init} \rangle$ **do**
 - $\text{delivered} := \emptyset; \text{past} := \text{nil}$
 - $\text{correct} := \Pi$
 - $\text{ack}[m] := \emptyset$ (for all m)



Algorithm 1' (gc – cont'd)

- ☛ **upon event** $\langle \text{crash } p_i \rangle$ **do**
 - ☛ $\text{correct} := \text{correct} \setminus \{p_i\}$
 - ☛ **if exists** m **such that** $\text{correct} \subseteq \text{ack}[m]$ **then**
 - ☛ $\text{past} := \text{past} \setminus \{(sm, m)\}$
- **upon exists** $m \in \text{delivered}$ **such that** $\text{self} \notin \text{ack}[m]$ **do**
 - $\text{ack}[m] := \text{ack}[m] \cup \{\text{self}\}$
 - **trigger** $\langle \text{rbBroadcast}(\text{ACK}, m) \rangle$



Algorithm 1' (gc – cont'd)

- **upon event** $\langle \text{rbDeliver } p_i, (\text{ACK}, m) \rangle$ **do**
 - $\text{ack}[m] := \text{ack}[m] \cup \{p_i\}$
 - **if** $\text{correct} \subseteq \text{ack}[m]$ **then**
 - $\text{past} := \text{past} \setminus \{(s_m, m)\}$

Algorithm 1' (gc – cont'd)

Alternative formulation



- ☞ **upon event** $\langle \text{crash } p_i \rangle$ **do**
 - ☞ $\text{correct} := \text{correct} \setminus \{p_i\}$
 - ☞ **if** $\text{correct} \subseteq \text{ack}[m]$ **then**
 - ☞ $\text{past} := \text{past} \setminus \{(sm, m)\}$
- **upon** $\langle \text{rcoDeliver } p_i, m \rangle$ **do**
 - **trigger** $\langle \text{rbBroadcast}(\text{ACK}, m) \rangle$



Algorithm 1' (gc – cont'd)

Alternative formulation

- upon event $\langle \text{rbDeliver } p_i, (\text{ACK}, m) \rangle$ do
 - if $p_i \in \text{correct}$ then
 - $\text{ack}[m] := \text{ack}[m] \cup \{p_i\}$
 - if $\text{correct} \subseteq \text{ack}[m]$ then
 - $\text{past} := \text{past} \setminus \{(s_{m'}, m)\}$
 - delete $\text{ack}[m]$

Fail-Silent Algorithm

Waiting Causal Broadcast



- Represent the past by a vector of sequence numbers (array of integers) called vector clock (VC)
- $\text{rank}(p_i) = i$, convert the process identifier to an integer that is used to index the vector
- At process p_i
 - $\text{VC}[i]$: the number of messages p_i coBroadcasted
 - $\text{VC}[j]$, $j \neq i$: the number of messages p_i coDelivered from p_j

Fail-Silent Algorithm

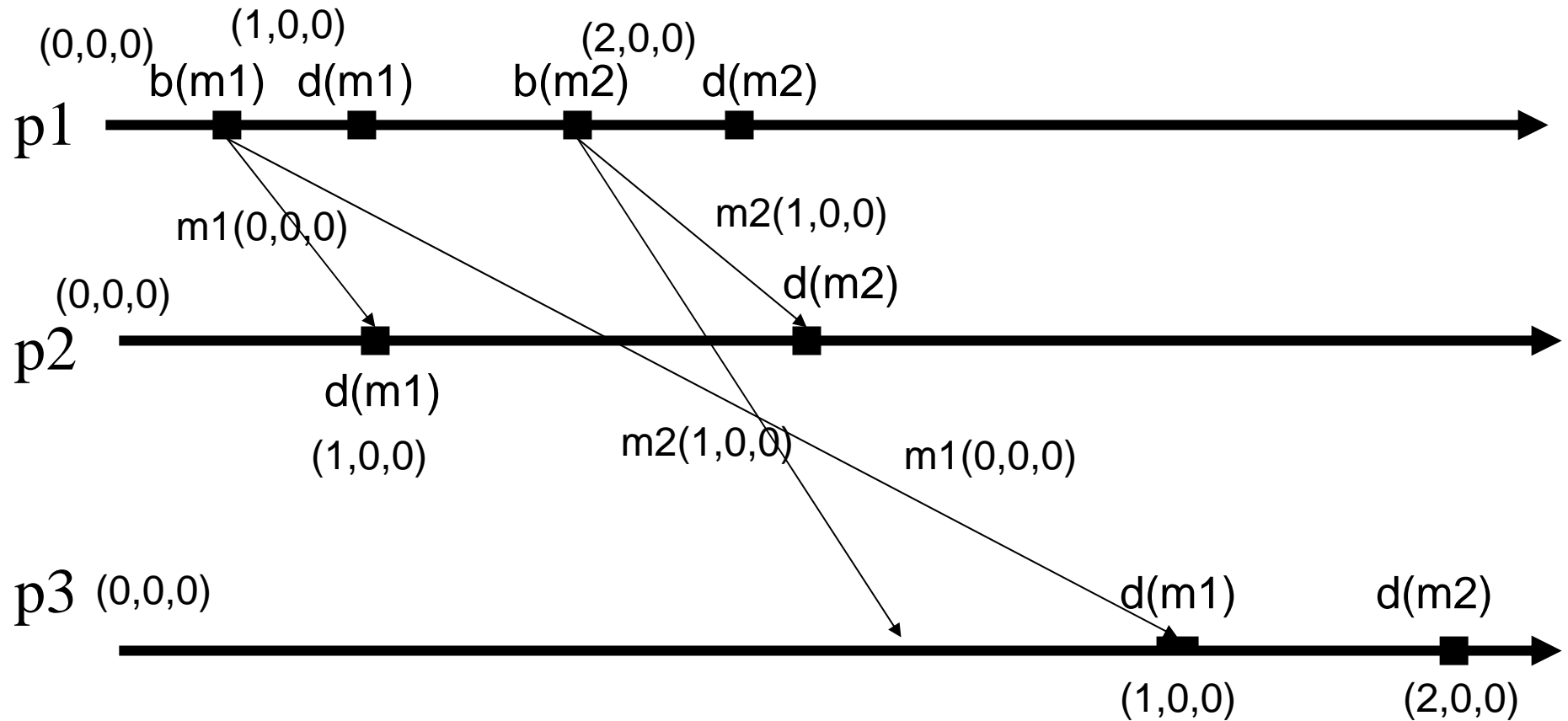
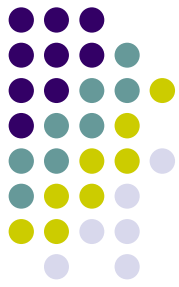
Waiting Causal Broadcast



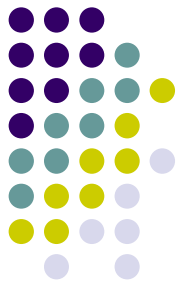
- This vector is attached to each message m that p coBroadcasts
- A process q that $rbDelivers$ m compares this vector with its own vector to determine which messages are missing, and from which processes
- Process q does not $coDeliver$ m until the missing messages are $coDelivered$ (Waiting)

Algorithm 2

rcod(elivery) rcob(broadcast)



Algorithm 2



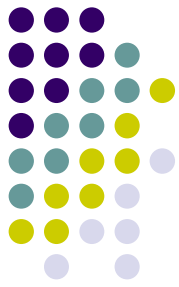
- ☛ **Implements:** ReliableCausalOrderBroadcast (rco)
- ☛ **Uses:** ReliableBroadcast (rb)
- ☛ **upon event** $\langle \text{Init} \rangle$ **do**
 - ☛ **forall** $p_i \in \Pi$ **do** $\text{VC}[\text{rank}(p_i)] := 0$
- ☛ **upon event** $\langle \text{rcoBroadcast } m \rangle$ **do**
 - ☛ **trigger** $\langle \text{rbBroadcast } (\text{DATA}, \text{VC}, m) \rangle$
 - ☛ $\text{VC}[\text{self}] := \text{VC}[\text{self}] + 1$
 - ☛ **trigger** $\langle \text{rcoDeliver self}, m \rangle$

Algorithm 2



- ☛ **upon event** $\langle \text{rbDeliver } p_j, (\text{DATA}, \text{VC}_m, m) \rangle$ **do**
 - ☛ **if** $p_j \neq \text{self}$ **then**
 - ☛ $\text{pending} := \text{pending} \cup (p_j, (\text{DATA}, \text{VC}_m, m))$
 - ☛ **deliver-pending**

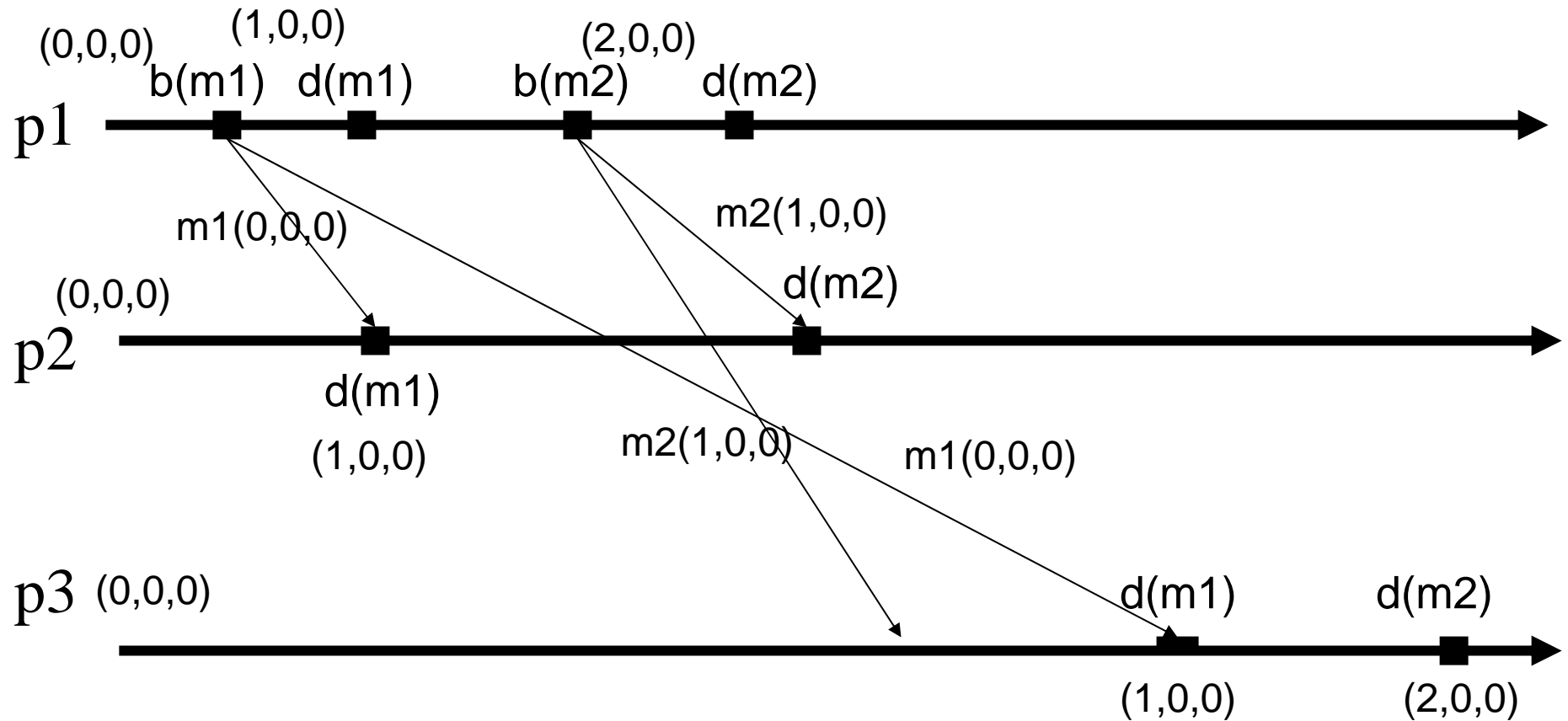
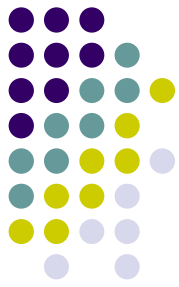
Algorithm 2



- ☛ procedure deliver-pending
 - ☛ while exists $x=(s_m, (DATA, VC_m, m)) \in \text{pending}$ such that $VC \geq VC_m$ do
 - ☛ pending := pending \ (s_m, (DATA, VC_m, m))
 - ☛ VC[rank(s_m)] := VC[rank(s_m)] + 1
 - ☛ trigger $\langle \text{rcoDeliver } s_m, m \rangle$

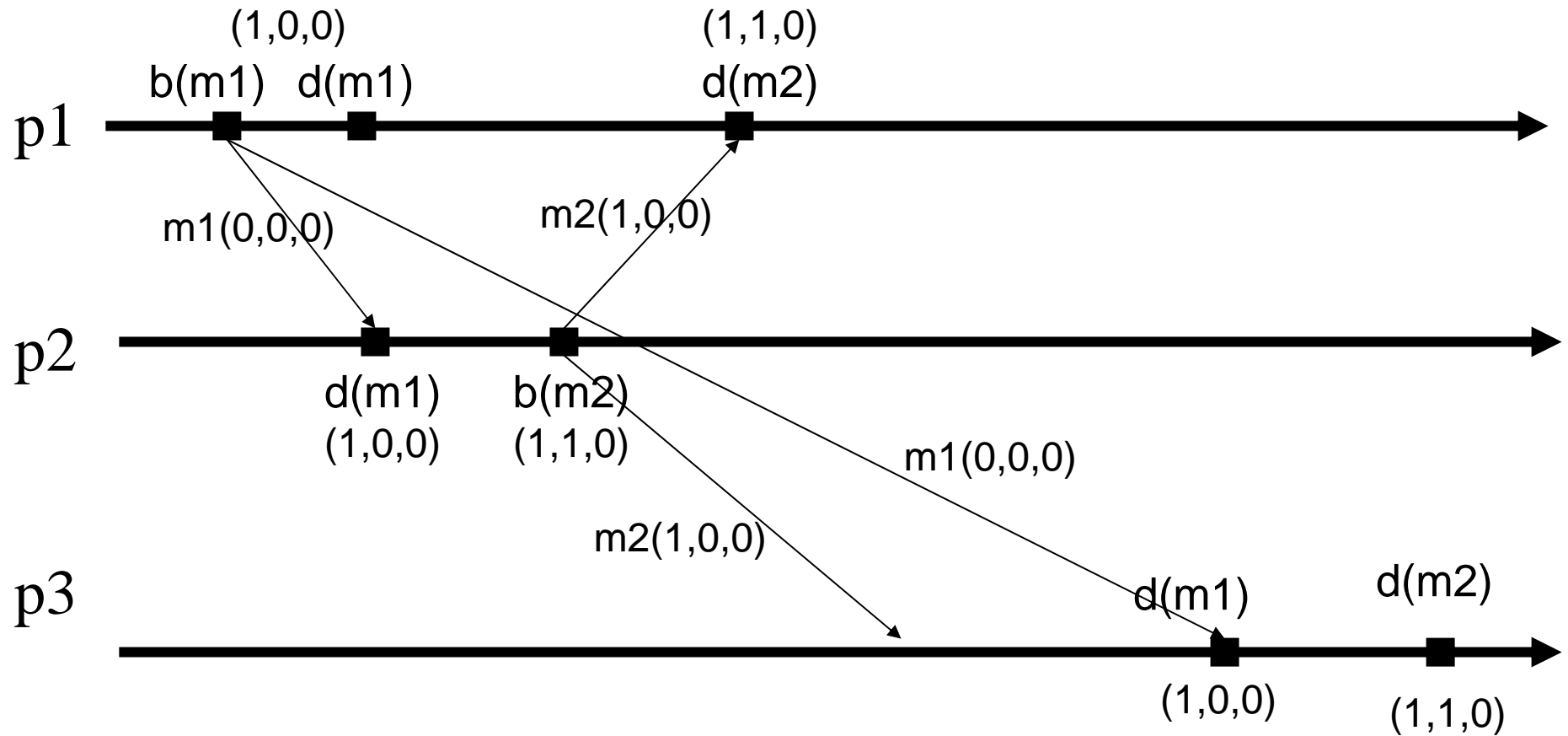
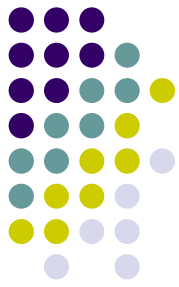
Algorithm 2

rcod(elivery) rcob(broadcast)



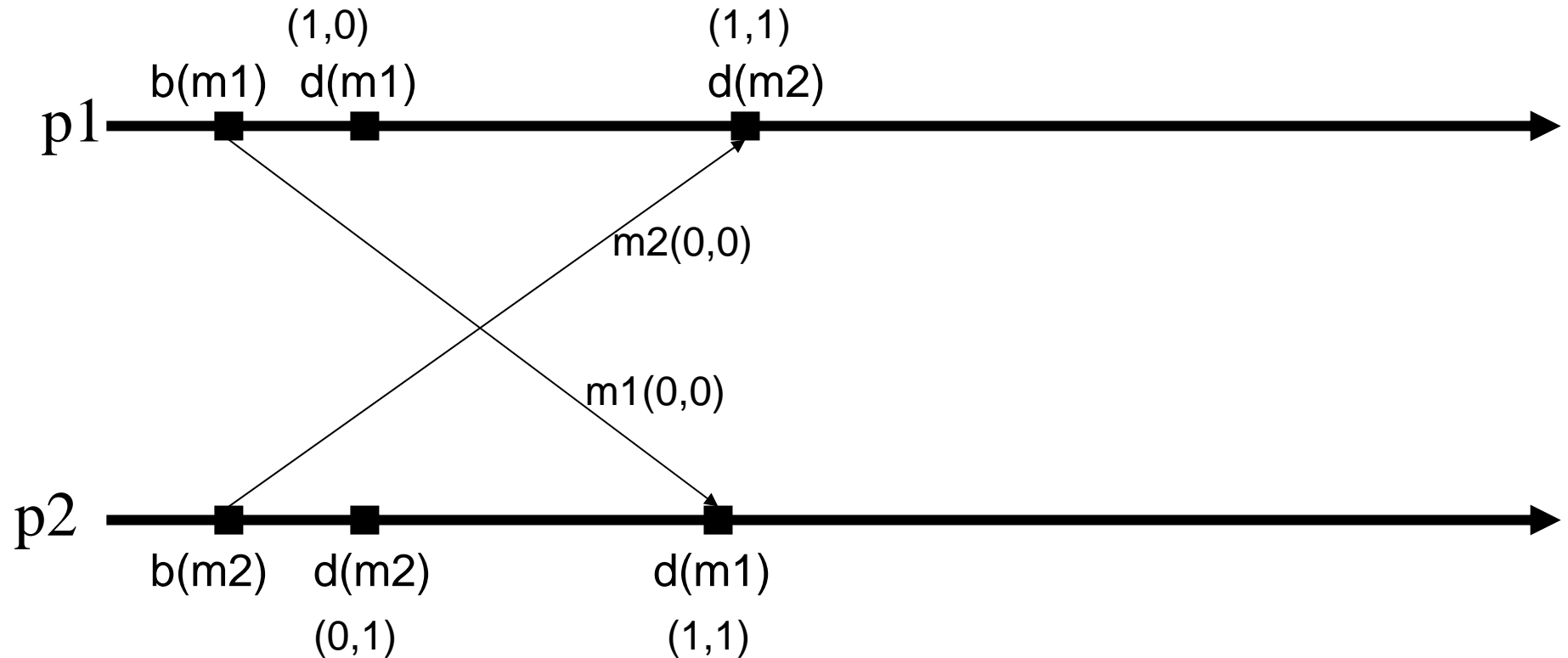
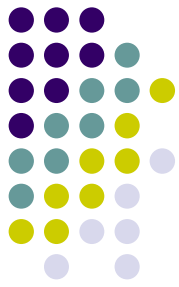
Algorithm 2

rcod(elivery) rcob(broadcast)



Algorithm 2

rcod(elivery) rcob(broadcast)



Notice different delivery order