

Algorithmes Distribués

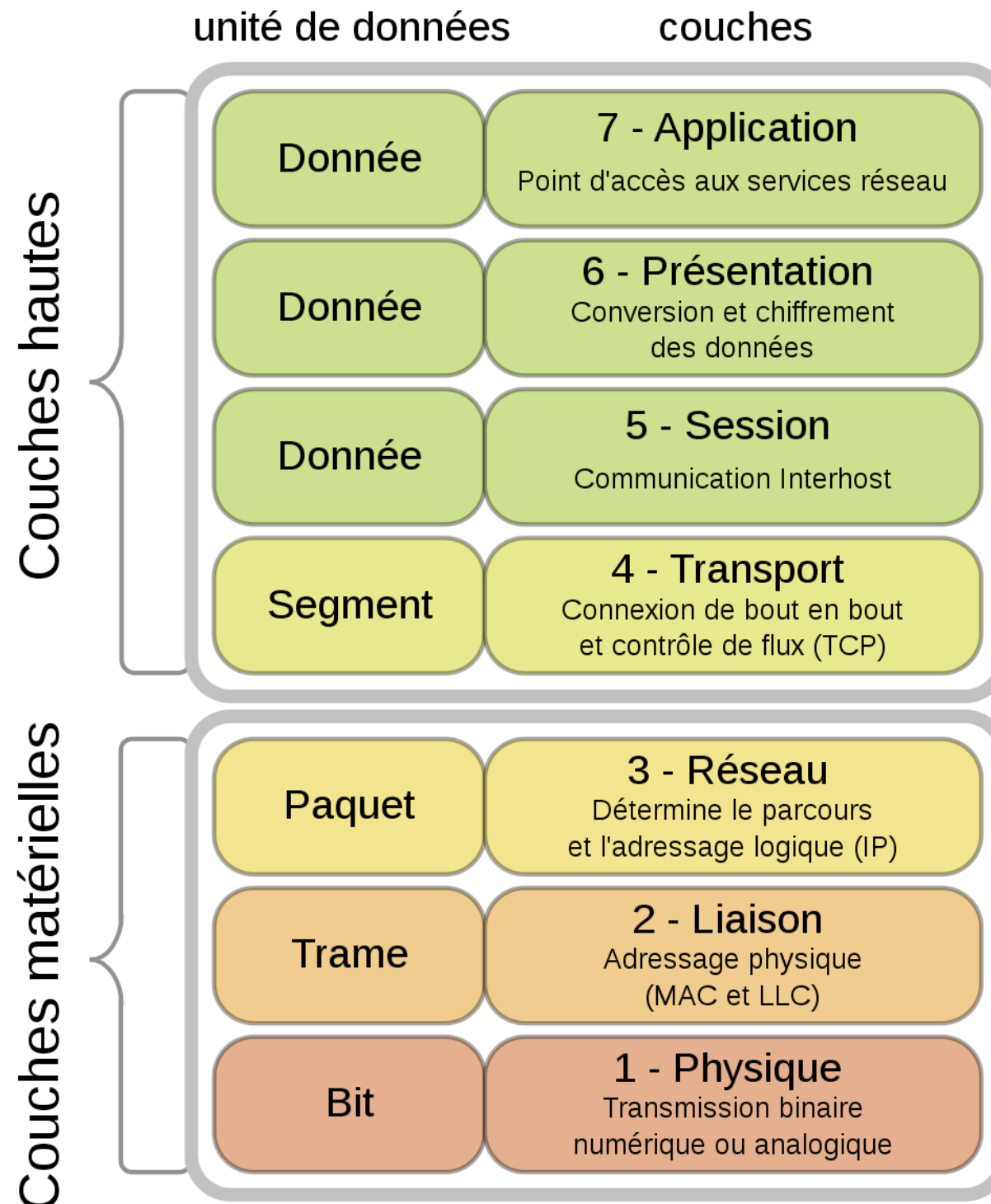
Application réseau

L3 - AMU

2021-2022

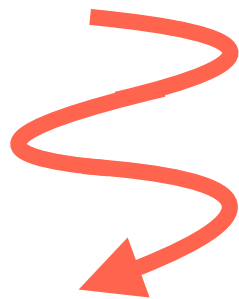
corentin travers

Modèle en couche



Abstraction

processus



processus



canal point-à-point

Couche

Couches matérielles

Segment

4 - Transport

Connexion de bout en bout
et contrôle de flux (TCP)

Paquet

3 - Réseau

Détermine le parcours
et l'adressage logique (IP)

Trame

2 - Liaison

Adressage physique
(MAC et LLC)

Bit

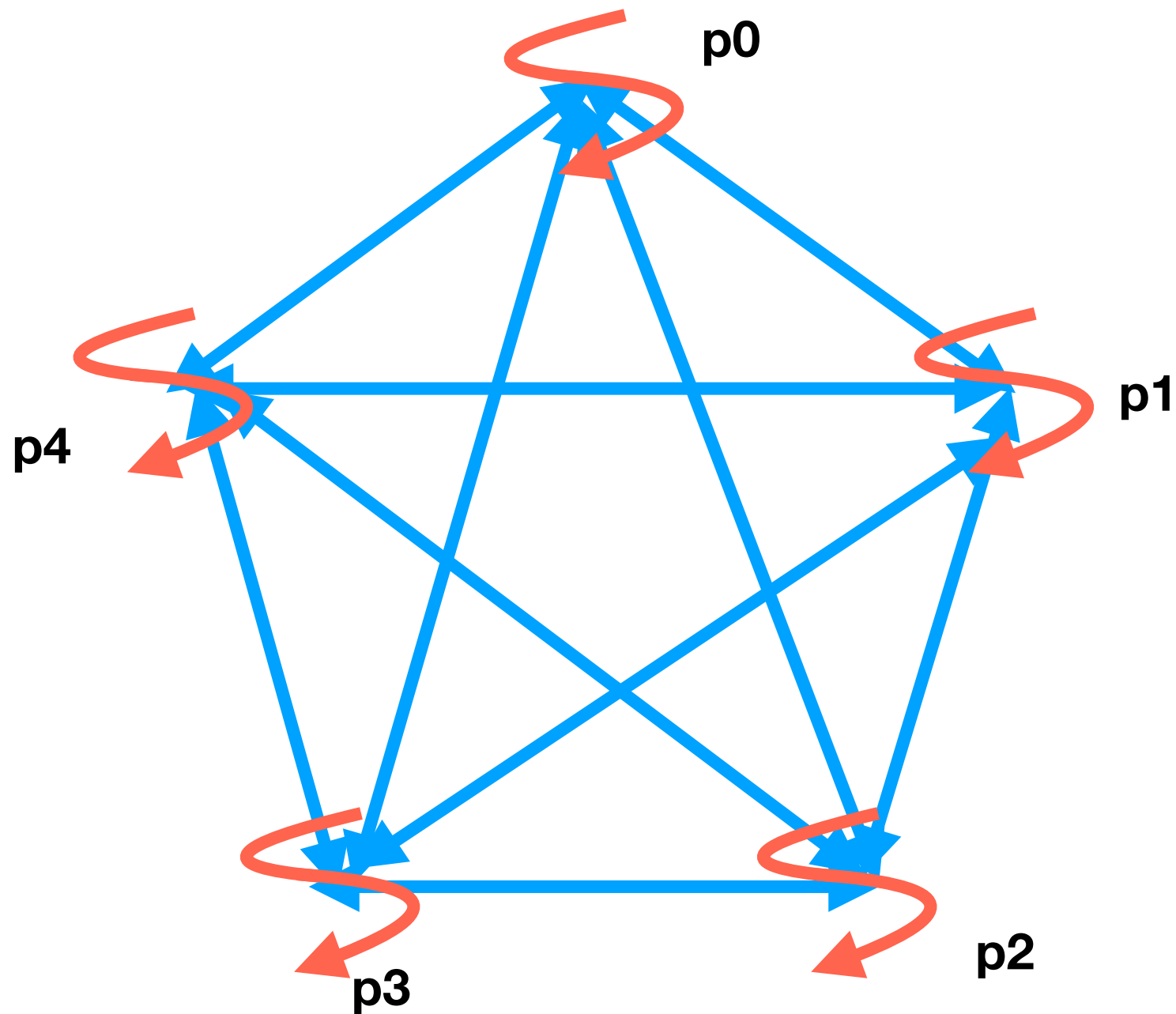
1 - Physique

Transmission binaire
numérique ou analogique

Couche transport

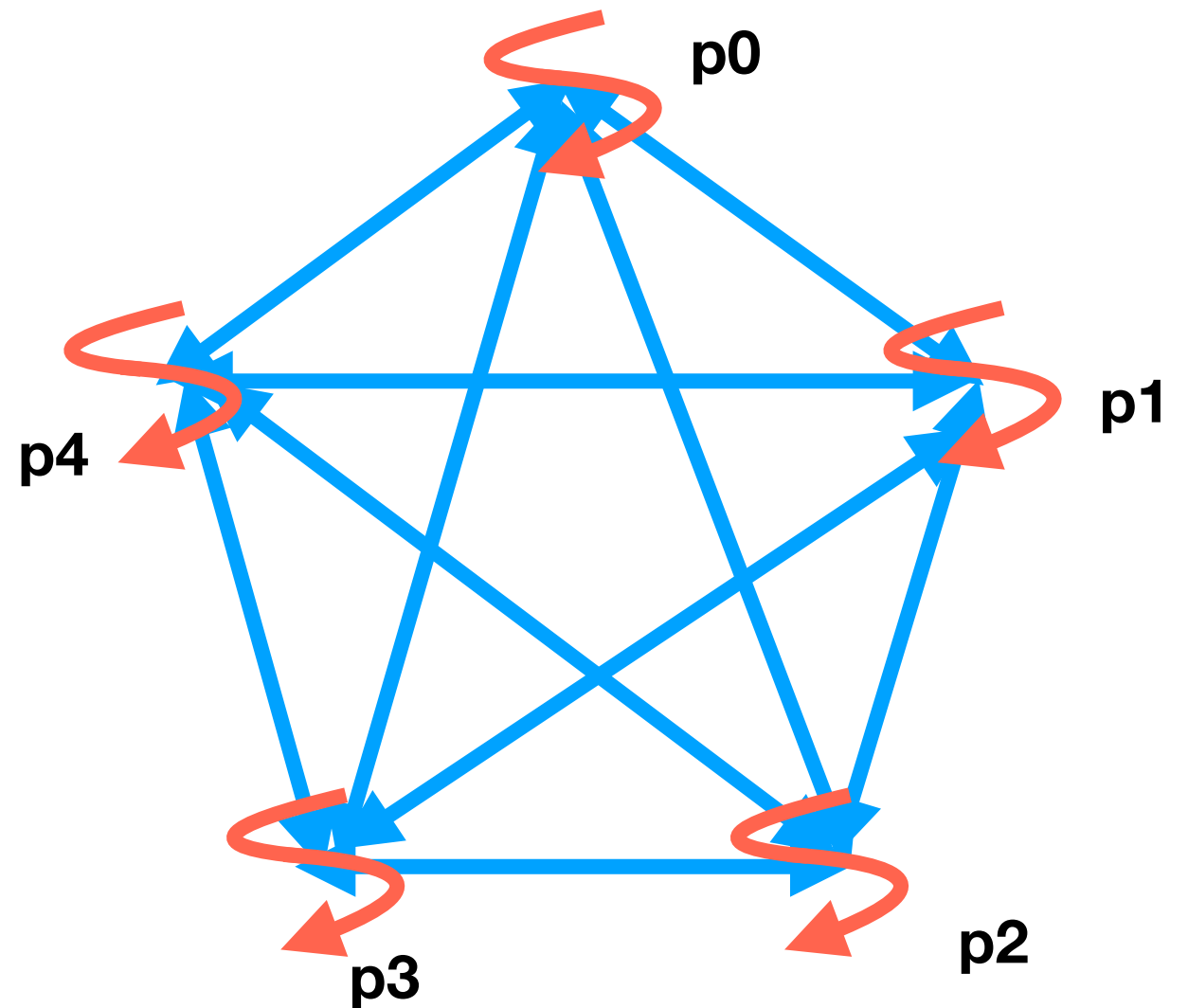
- **Canal** de communication
- Point-à-point : Entre **2 processus**
- **Fiable** (TCP) ou non (UDP)
- **Asynchrone** temps de transfert des messages non borné

Systeme distribue

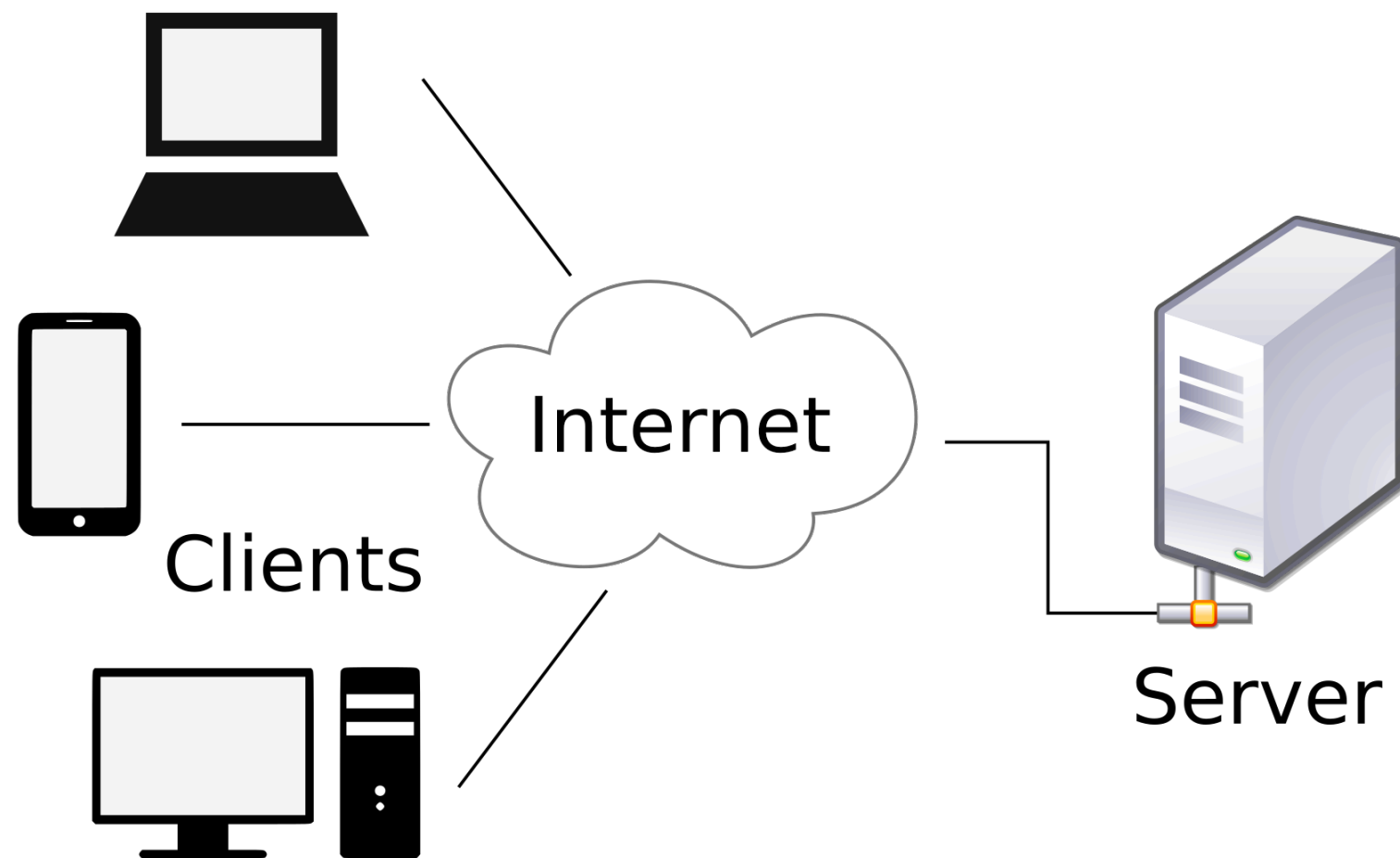


Systeme distribué

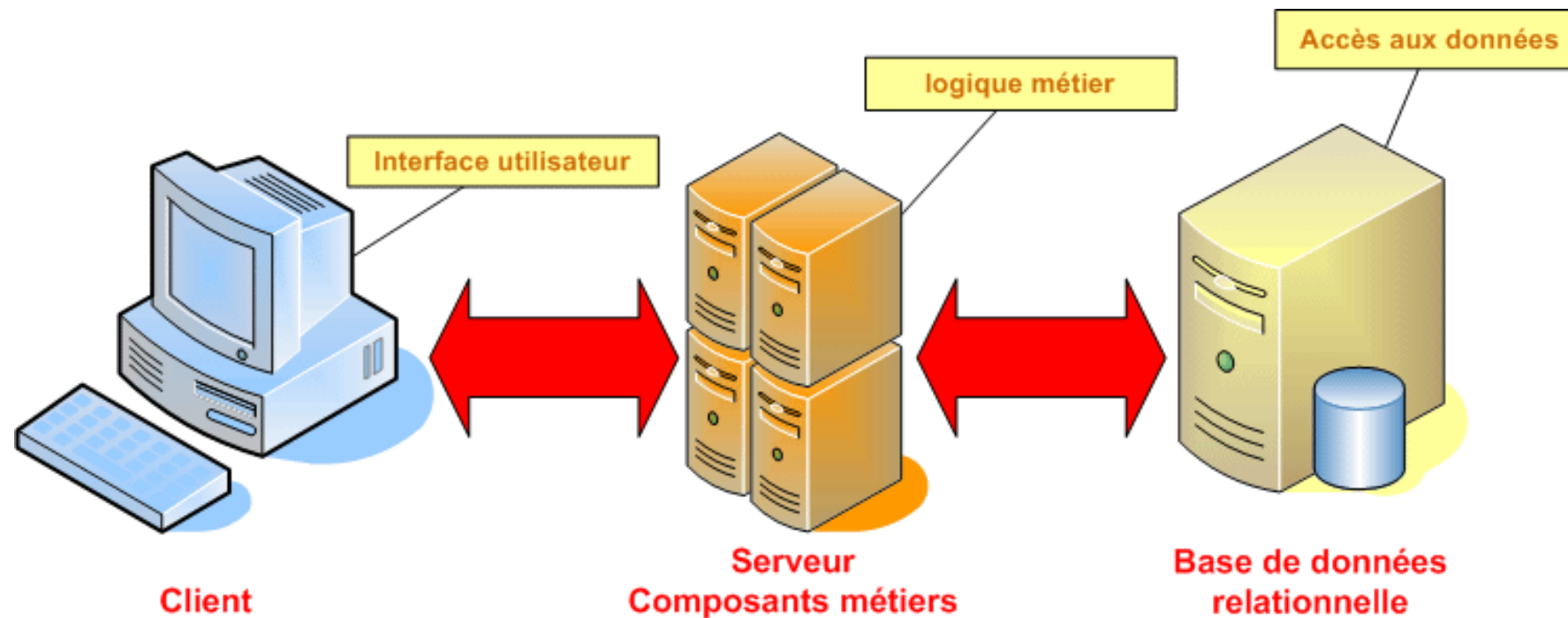
- Processus
- Canaux de communication
- Communication par **messages**



Client/Serveur

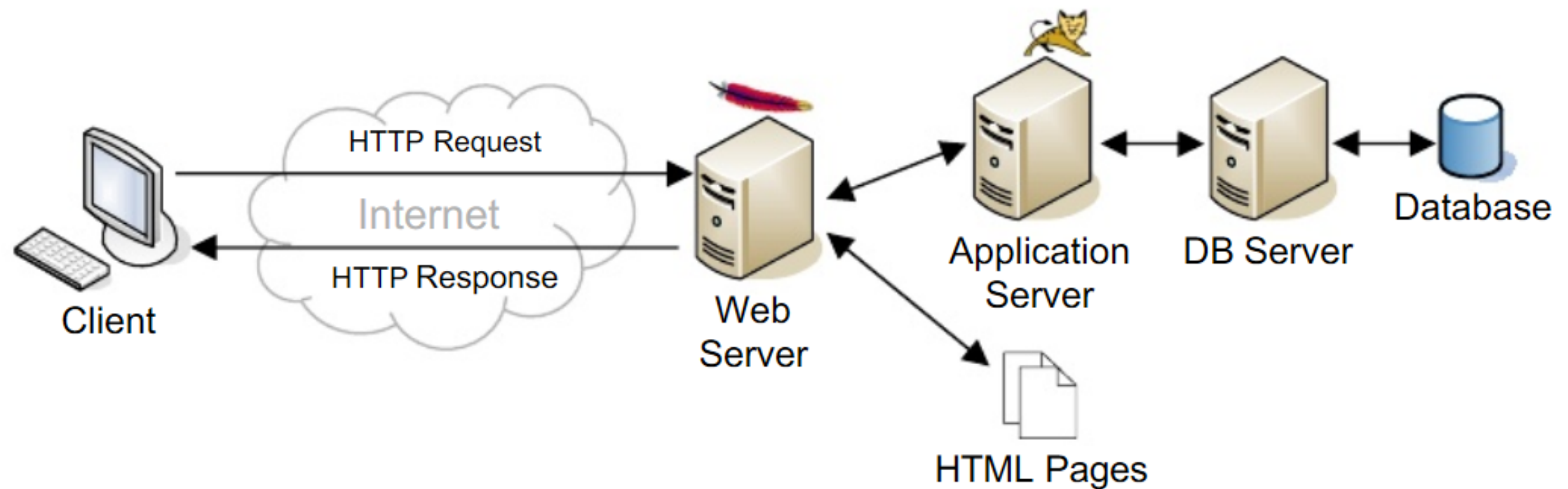


Architecture 3-tiers

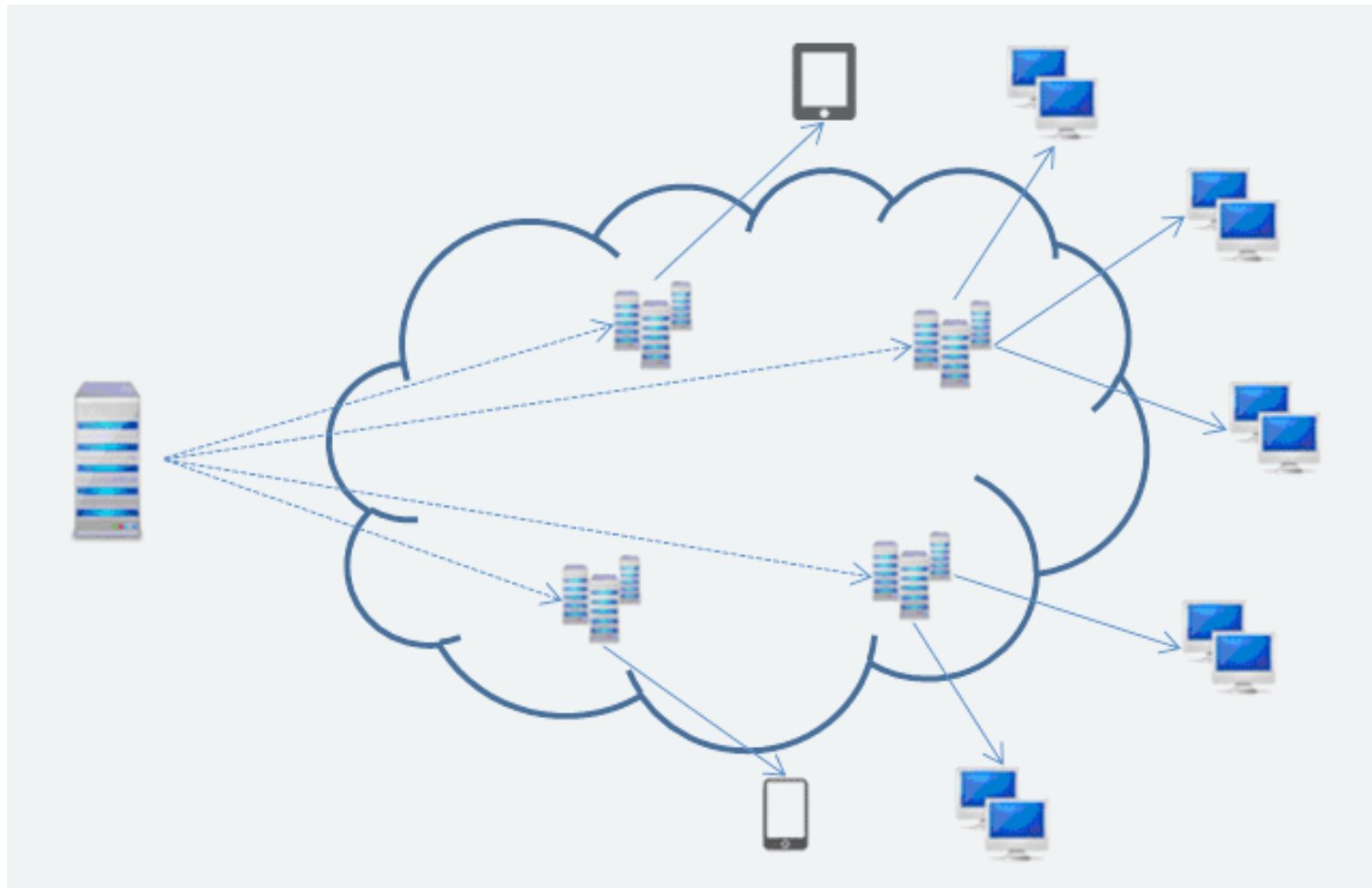


- **Présentation : le client**
- **Traitement : serveur(s)**
- **Accès/stockage des données: BDD**

Architecture N-tiers



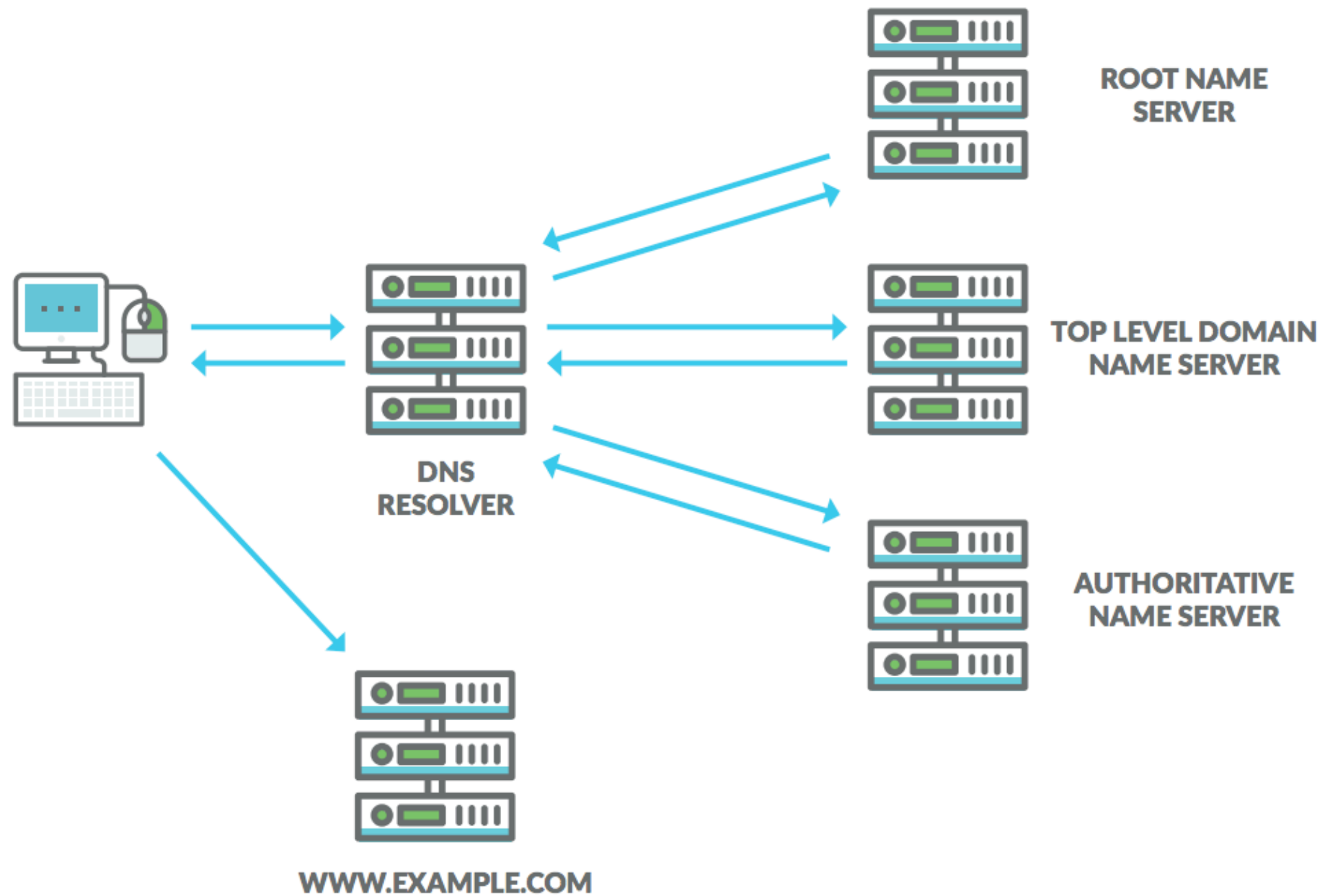
Content Delivery Network



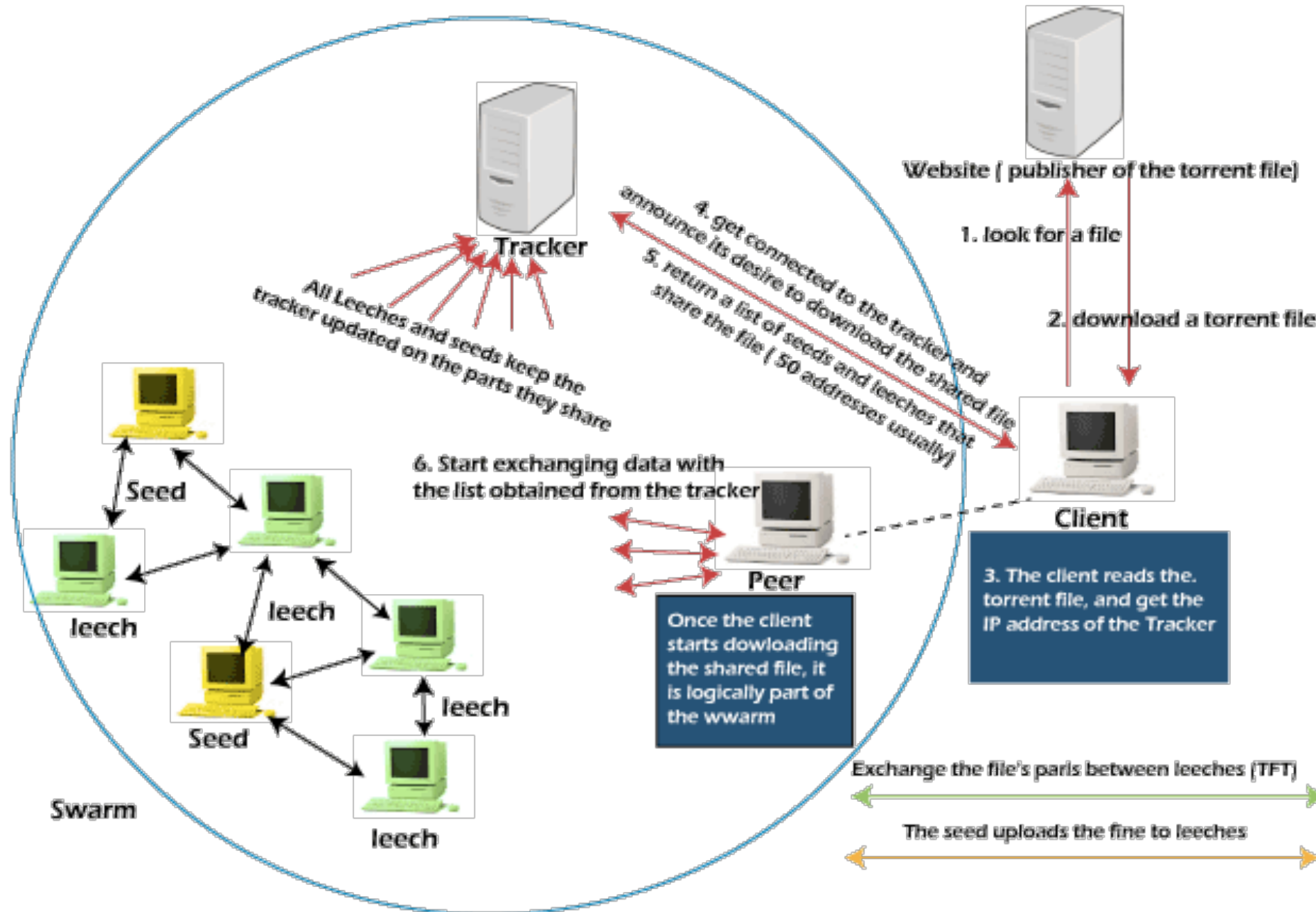
Diffusion de contenu

- **Répartition de la charge**
- **Caching**
- **Disponibilité**

DNS



Partage de fichier P2P



Systeme distribué

- Ensemble de **processus** p_0, p_1, \dots, p_N
- Qui peuvent communiquer par transmission de **messages**
- Doivent se **coordonner** pour effectuer une tâche commune/rendre un service

Systeme distribue : problematiques

- Processus peuvent mal-fonctionnés : panne, crash, bug, piratage
- Messages peuvent être perdus (UDP), modifiés, espionnés, ...
- Messages peuvent être retardés : délais de transfert **non bornés**

Exécution

p0:

```
m0 = new Message()  
send m0 to p1, p2
```

p1:

```
upon reception of m from p2  
send m to p0, p2
```

p2:

```
m2 = new Message()  
send m2 to p1
```

Exécutions

p0:

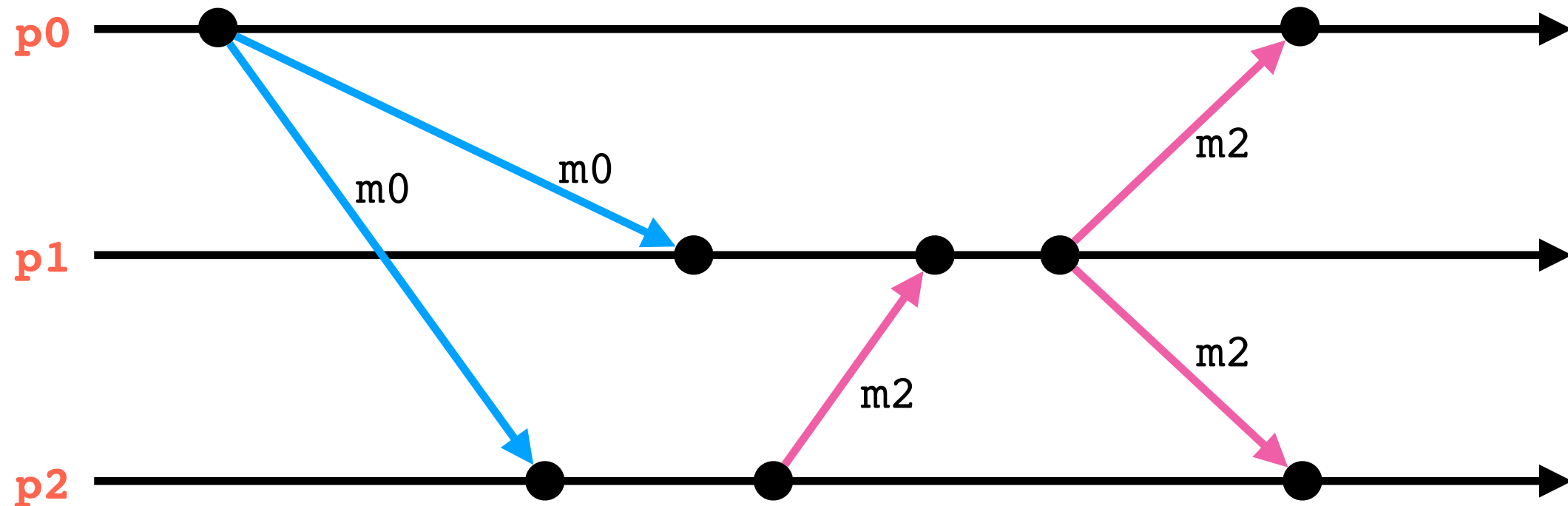
```
m0 = new Message()  
send m0 to p1, p2
```

p1:

```
upon reception of m from p2  
send m to p0, p2
```

p2:

```
m2 = new Message()  
send m2 to p1
```



Exécutions

p0:

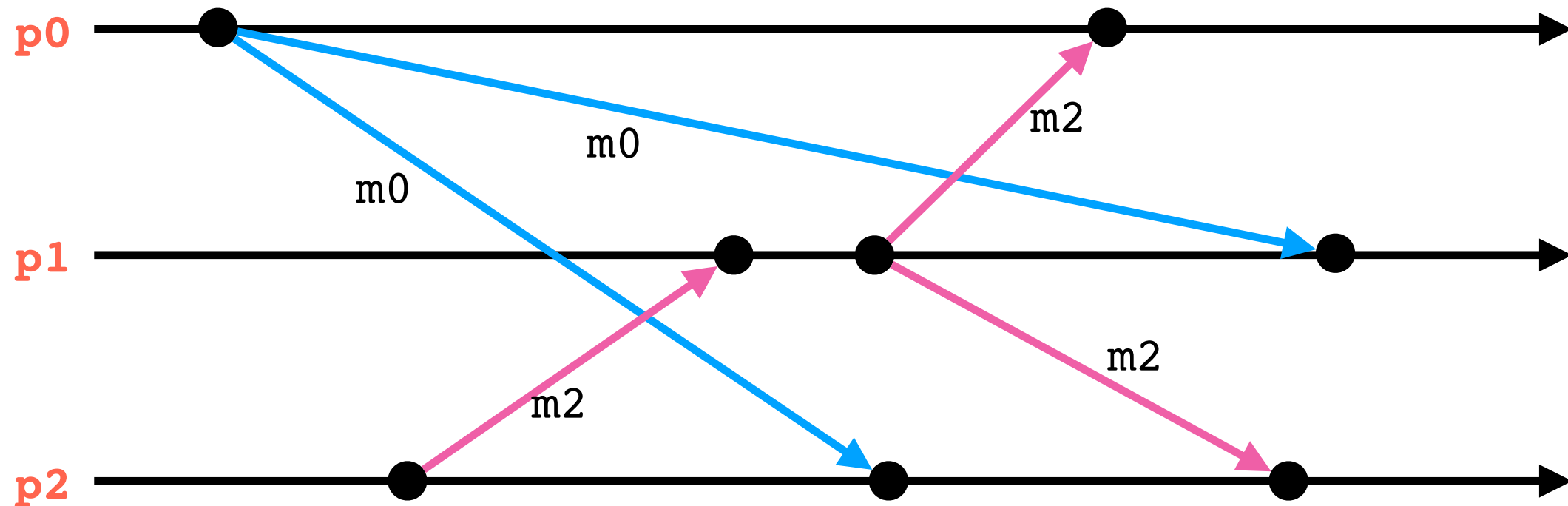
```
m0 = new Message()  
send m0 to p1, p2
```

p1:

```
upon reception of m from p2  
send m to p0, p2
```

p2:

```
m2 = new Message()  
send m2 to p1
```



Exécutions

p0:

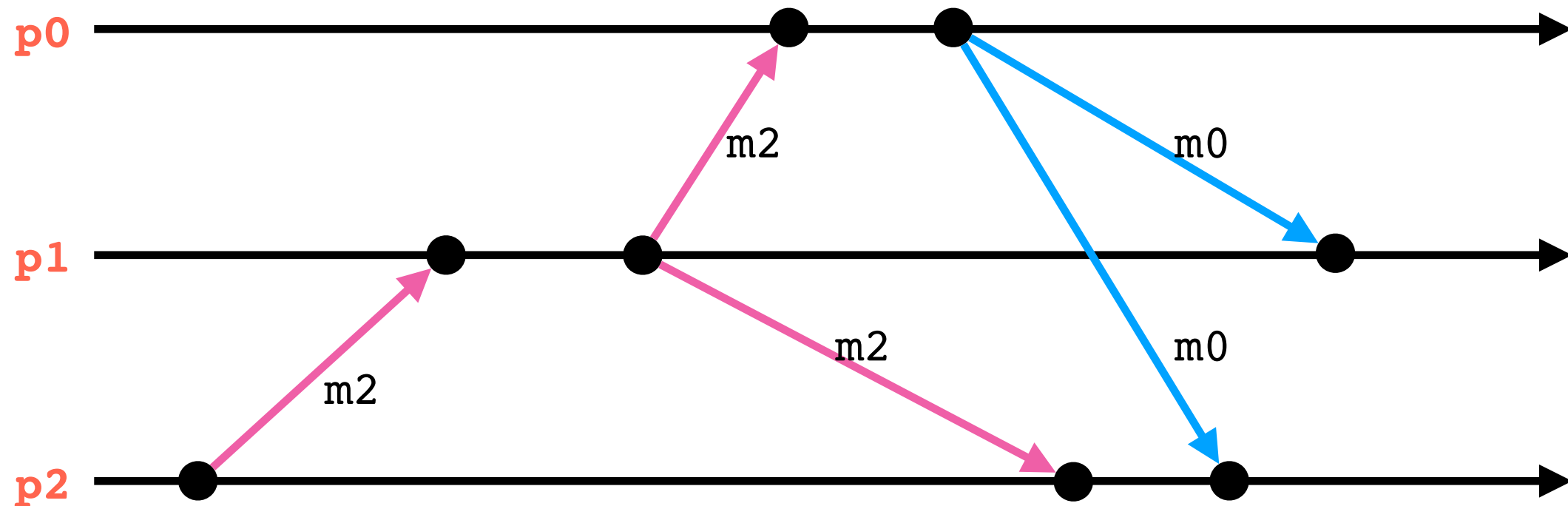
```
m0 = new Message()  
send m0 to p1, p2
```

p1:

```
upon reception of m from p2  
send m to p0, p2
```

p2:

```
m2 = new Message()  
send m2 to p1
```



Exécutions distribuées

p0:

```
m0 = new Message()  
send m0 to p1, p2
```

p1:

```
upon reception of m from p2  
send m to p0, p2
```

p2:

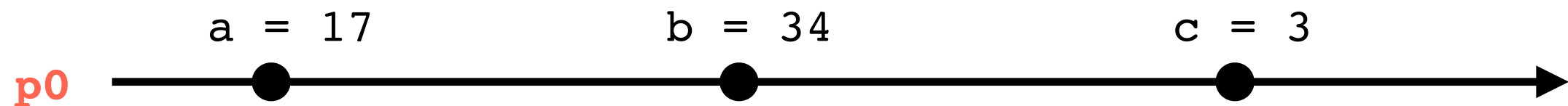
```
m2 = new Message()  
send m2 to p1
```

- Même programme
- De nombreuses exécutions possibles ! Non déterminisme !
- Ordre de réception des messages imprévisibles !
- Vitesses relatives des processus imprévisibles !

Exécution séquentielle

p0:

```
a = 17;  
b = 2*a;  
if(b == 0) { ...}  
else {c= 3}
```



- Une seule exécution ! déterministe !

Événements

- Un **événement** est l'exécution d'une instruction par un processus
- **événement** lié à un canal de communication
 - **envoi** d'un message sur ce canal
 - **réception** d'un message sur ce canal
- **événement interne** : tout autre événement
- On s'intéresse aux événements de communication et à certains événements internes (dépend de l'application)

Exécution séquentielle

- Une exécution séquentielle est une **séquence d'événements**
- Trace de l'exécution d'un programme (cf. débogage avec `printf`, fichier de log)
- Ordre **total** sur les événements dans une exécution séquentielle

Exécution distribuée

- (En général), pas d'**ordre total** sur les événements
- Certains événements sont toutefois ordonnés:
 - événements sur un même processus
 - envoi d'un message *précède* sa réception

Exécution distribuée

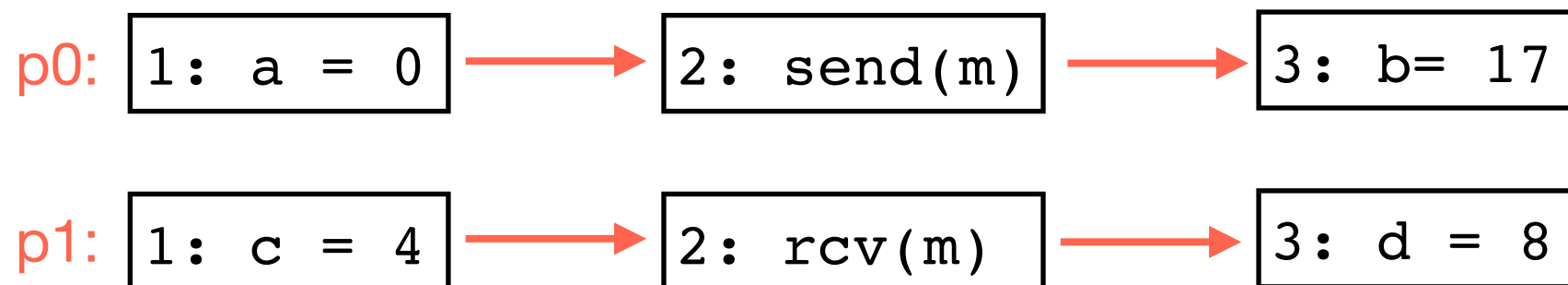
p0: trace

```
1: a = 0
2: send(m) to p1
3: b = 17
```

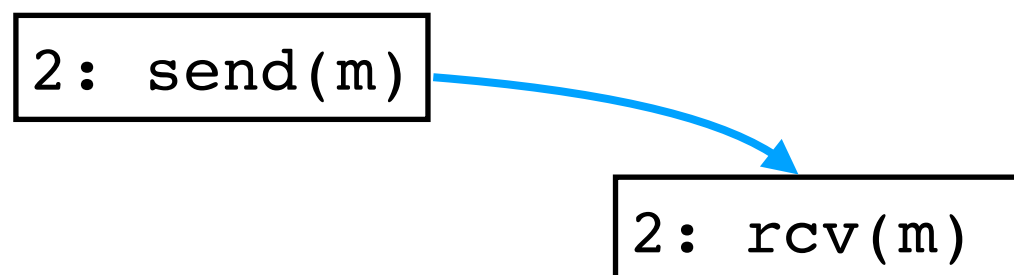
p1: trace

```
1: c = 4
2: rcv(m) from p1
3: d = 8
```

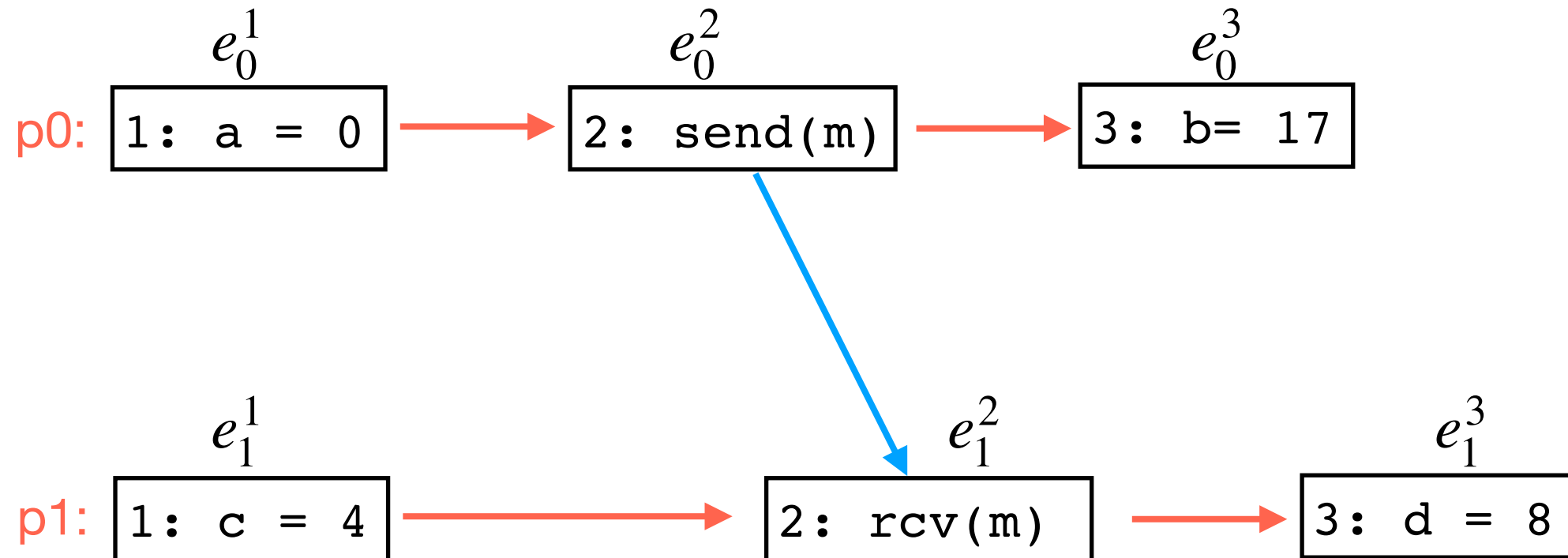
- Ordre induit par les processus :



- Ordre induit par les communication :



Ordre partiel

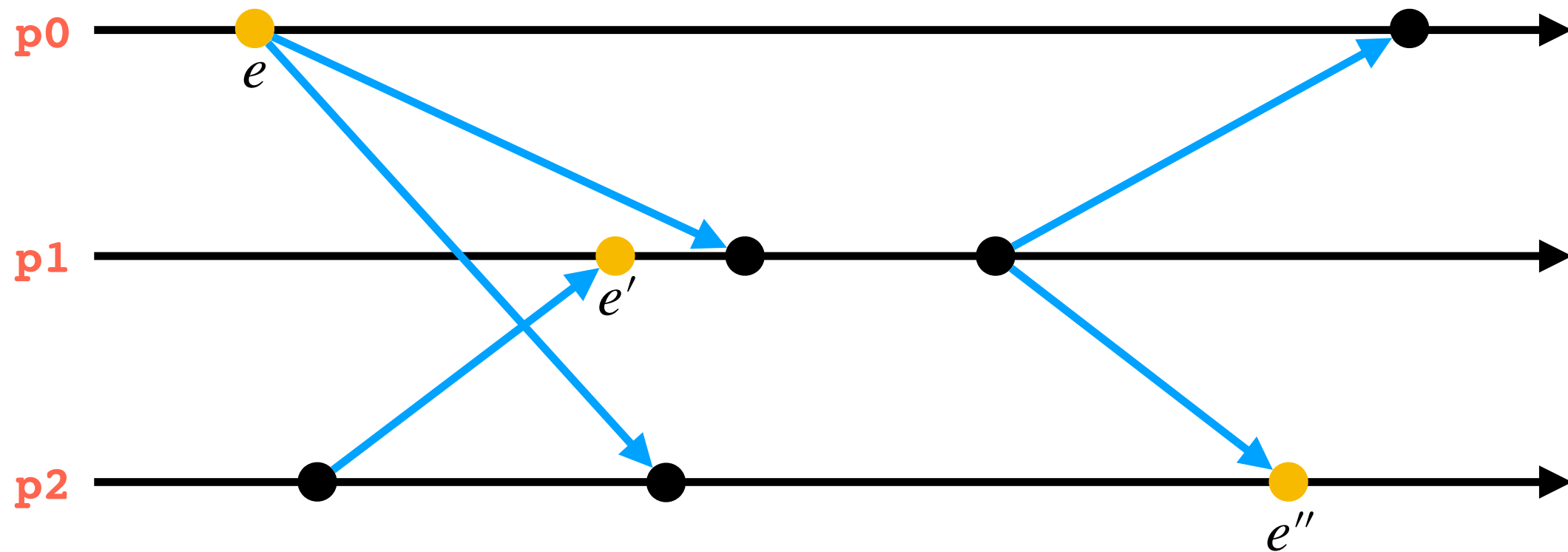


- e_0^1 précède e_0^2 ; e_0^2 précède e_1^2
- par transitivité, e_0^1 précède e_1^2
- e_0^1 et e_1^1 ne sont pas ordonnés

happens before ordre causal

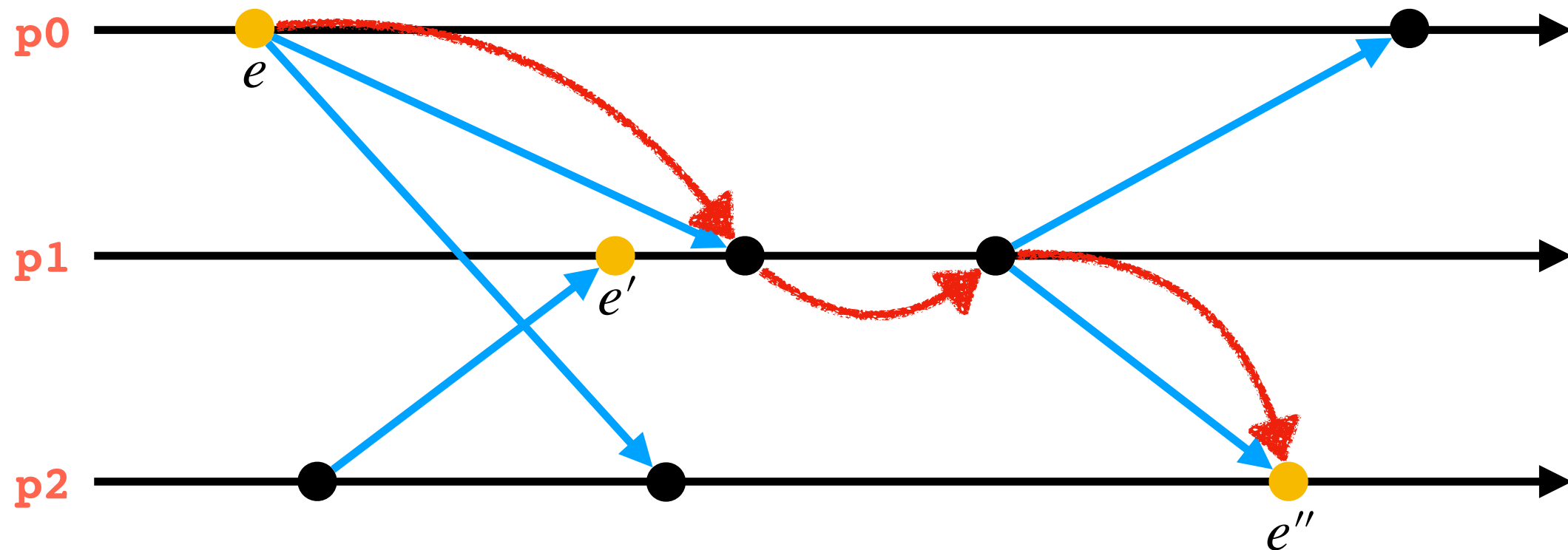
- $e \rightarrow e'$: **e a lieu avant e' ssi :**
 - e, e' ont lieu sur le même processus et e précède e' ou
 - $e = \text{send}(m)$ et $e' = \text{rcv}(m)$ ou
 - il existe e'' : $e \rightarrow e''$ et $e'' \rightarrow e'$

Ordre causal



- e et e' ?
- e et e'' ?

Ordre causal

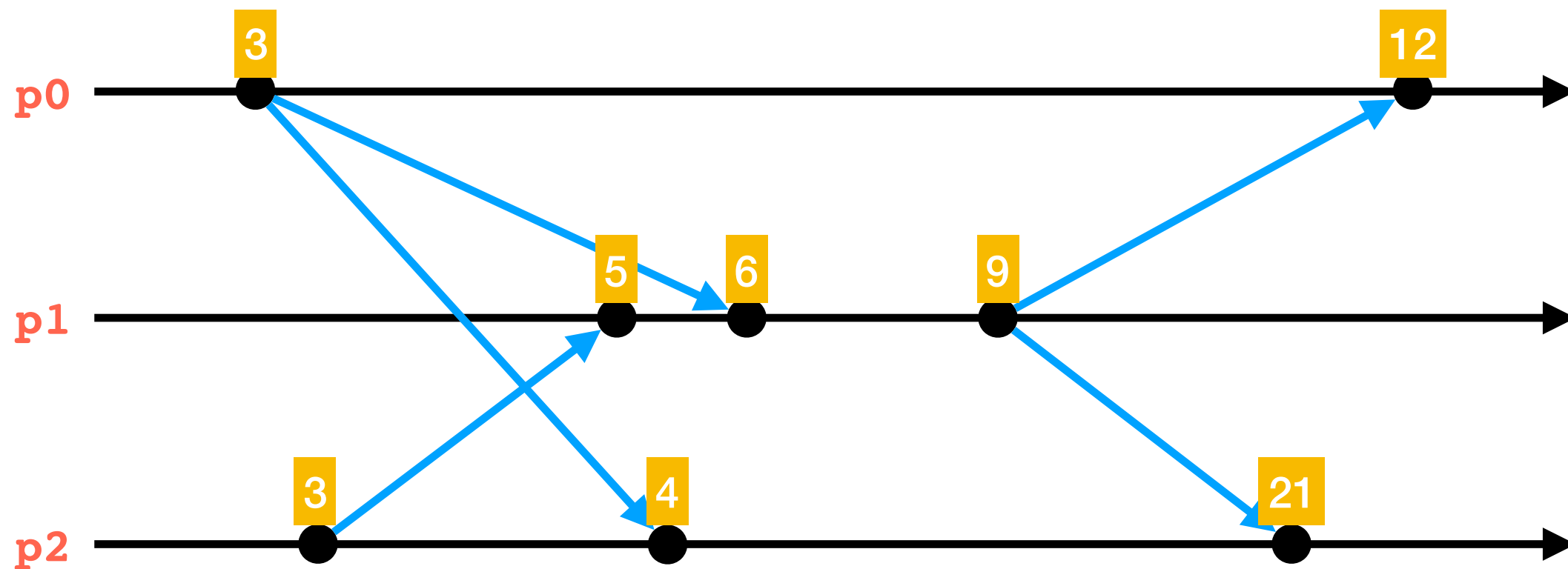


- e et e' sont indépendants noté $e \perp\!\!\!\perp e'$
- e précède causalement e'' $e \rightarrow e''$

Horloge logique

- **Objectif** : attribuer à chaque évènement une *date* telle que
 - $\forall e, e' : e < e' \implies date(e) < date(e')$
- **Moyen** : algorithme distribué
- **Remarque**: on ne demande pas que
 - $\forall e, e' : date(e) < date(e') \implies e < e'$

Exemple

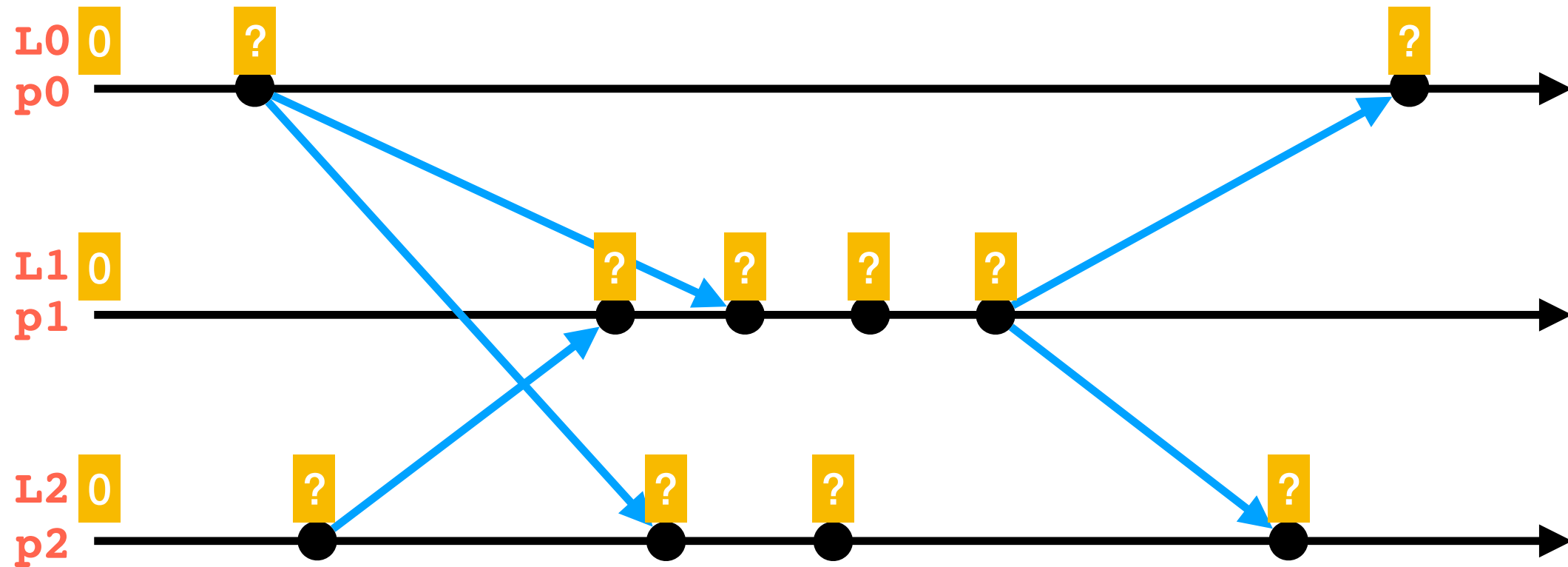


- la date est un entier
- comment faire ?

Horloge de Lamport

- Chaque processus p_i maintient une var. locale L_i
- 3 règles :
 1. sur p_i incrémenter L_i avant chaque évènement interne:
 2. lors de l'envoi d'un message m : incrémenter L_i puis envoyer le message avec L_i . I.e $L_i = L_i + 1$; $\text{send}(m, L_i)$
 3. à la réception d'une paire (m, L_j) : m.a.j de L_i avec $\max(L_i, L_j)$ puis appliquer règle 1. $\text{rcv}(m, L_j)$; $L_i = \max(L_i, L_j)$;
 $L_i = L_i + 1$

Horloge de Lamport



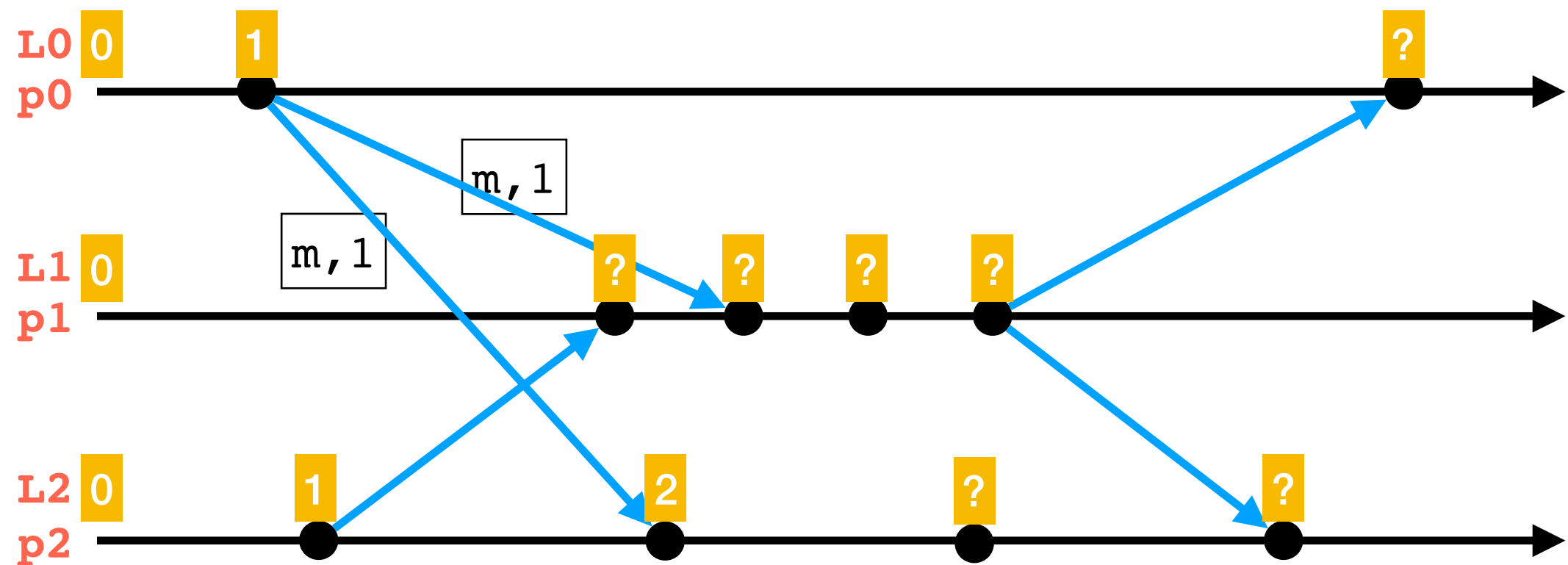
- 3 règles :

1. sur p_i incrémenter L_i avant chaque évènement interne: $L_i = L_i + 1$;

2. lors de l'envoi d'un message m : incrémenter L_i puis envoyer le message avec L_i .
 $L_i = L_i + 1$; $\text{send}(m, L_i)$

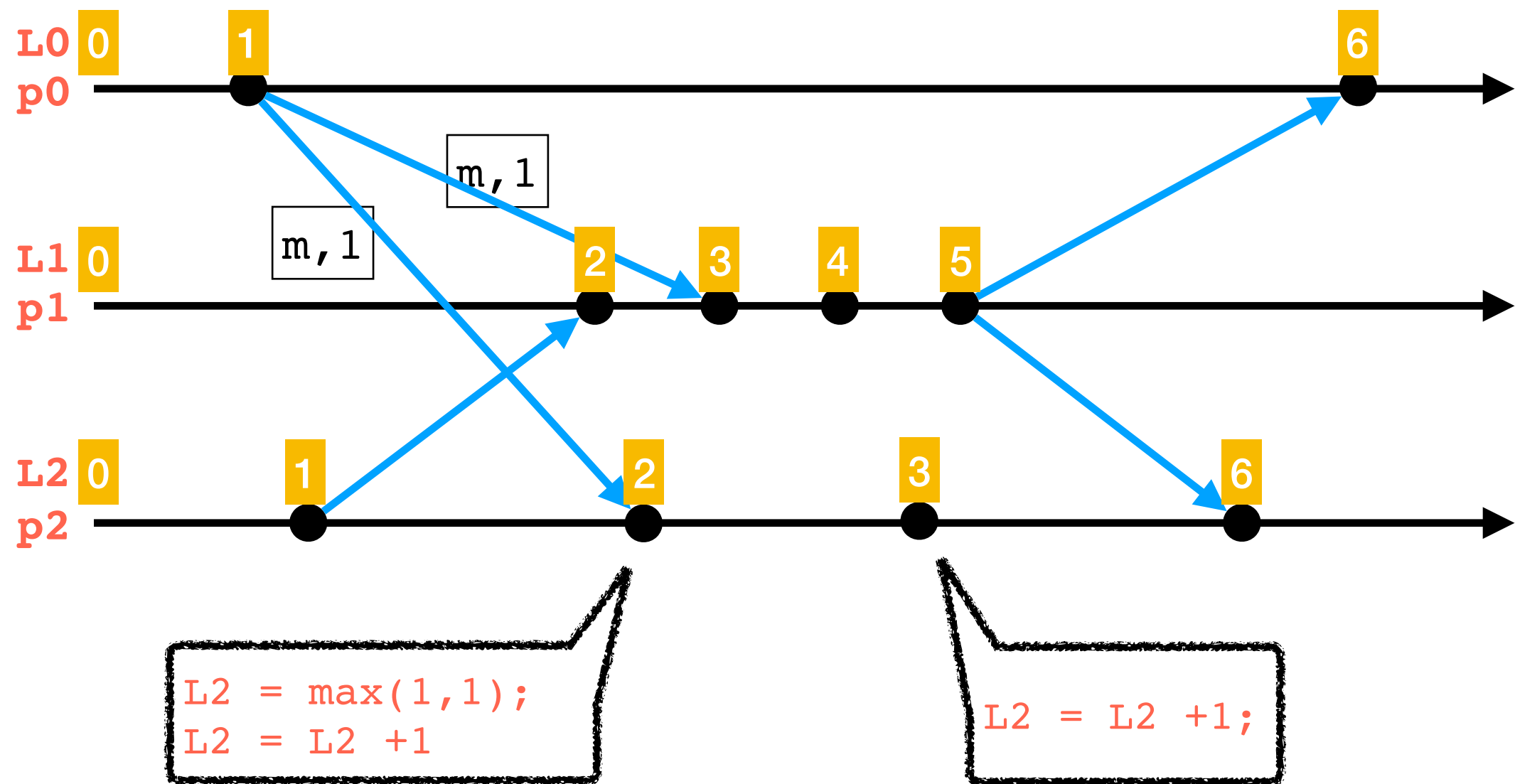
3. à la réception d'une paire (m, L_j) : maj de L_i avec $\max(L_i, L_j)$ puis appliquer règle 1.
 $\text{rcv}(m, L_j)$; $L_i = \max(L_i, L_j)$; $L_i = L_i + 1$

Horloge de Lamport

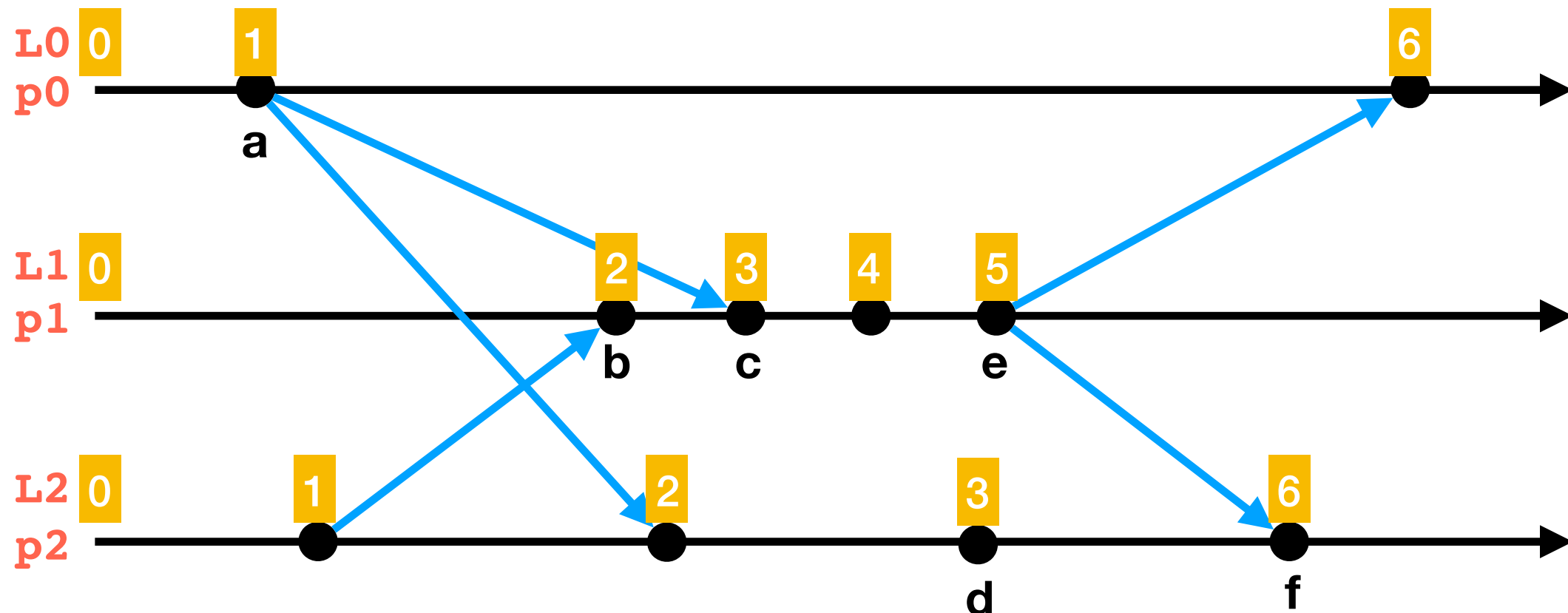


```
L2 = max(1, 1);  
L2 = L2 + 1
```

Horloge de Lamport



Horloge de Lamport

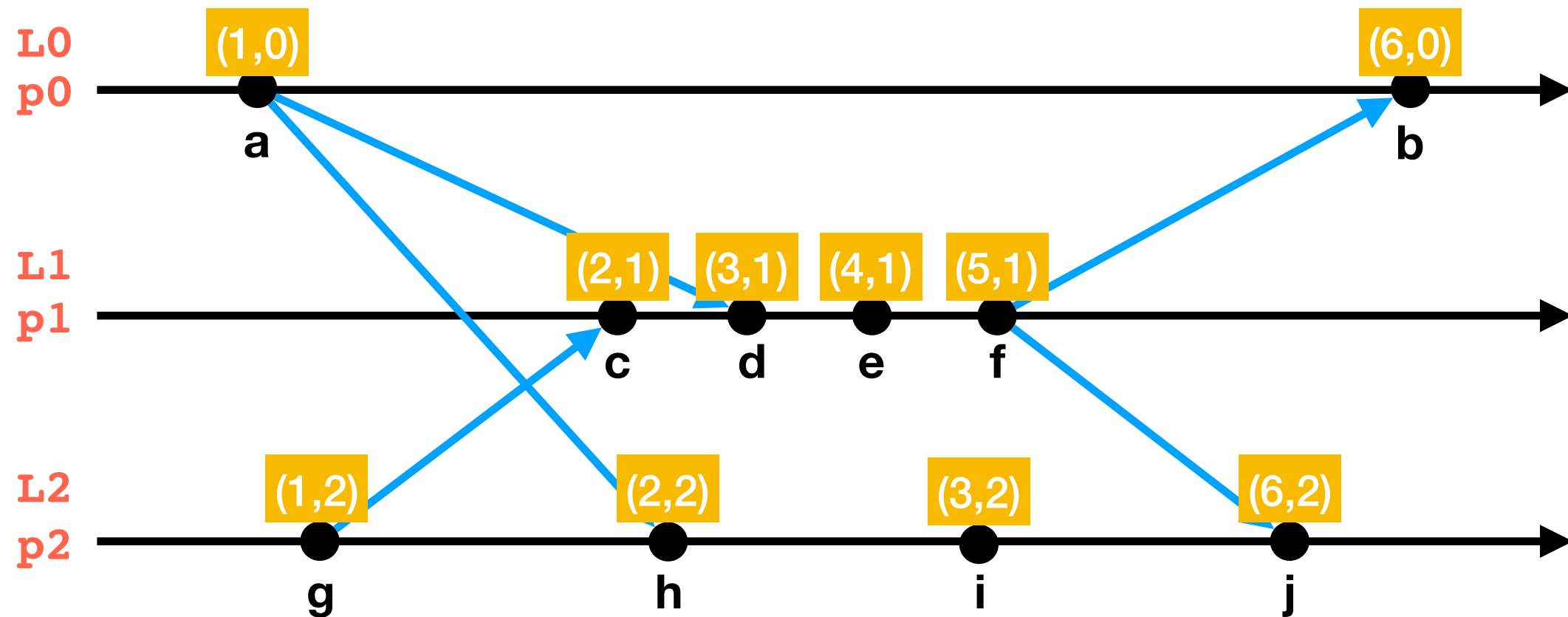


- **L** est un ordre partiel cohérent avec l'ordre causal :
- $a \rightarrow f$ on a bien $L(a) < L(f)$
- $a \rightarrow c$ on a bien $L(a) < L(c)$
- Par contre les évènements indépendants ont des dates arbitrairement ordonnées :
 - $b \parallel d$ et $L(b) < L(d)$ $c \parallel d$ et $L(c) = L(d)$ $e \parallel d$ et $L(e) < L(d)$

D'un ordre partiel à un ordre total

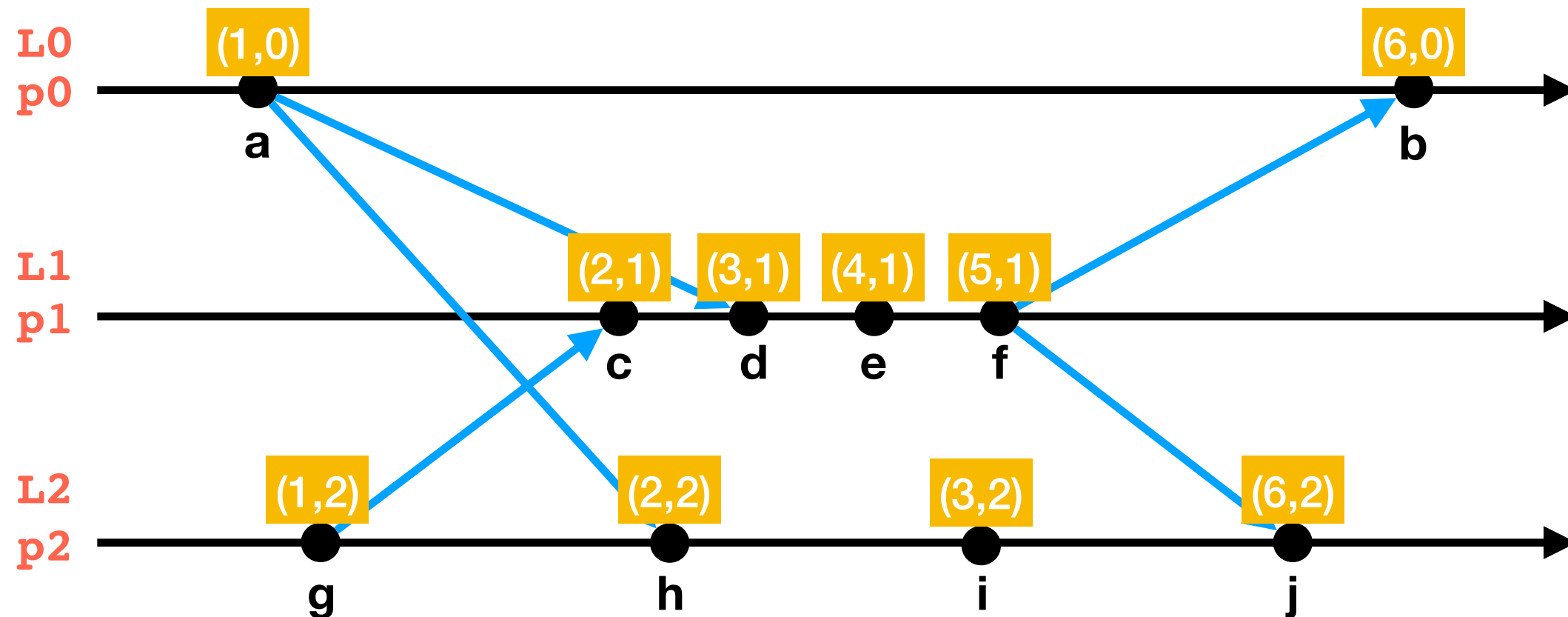
- Les horloges de Lamport ordonnent **partiellement** les évènements
 - pour $e \parallel e'$, il est possible que $L(e) == L(e')$
- Il est souvent souhaitable d'avoir un **ordre total**
- **Solution** : si e est un évènement sur p_i , lui associer le couple $(L(e), i)$
- **Ordre total cohérent avec l'ordre causal**
 $(L(e), i) < (L(e'), j) \iff L(e) < L(e') \text{ ou } (L(e) = L(e') \wedge i < j)$

Ordre total



- Ordre total ?

Ordre total (L,ID)



- **Ordre total (a,g,c,h,d,i,e,f,b,j)**
- **Cohérent avec l'ordre causal**

Horloge vectorielle

- Horloge de Lamport : on ne peut pas détecter les évènements indépendants. Pour $e \parallel e'$, on peut avoir $L(e) = L(e')$, $L(e) < L(e')$ ou $L(e) > L(e')$
- **Objectif** attribuer à chaque événement une date qui capture l'ordre causal :

$$\forall e, e', e \rightarrow e' \iff date(e) < date(e')$$

Horloge vectorielle

- Chaque processus p_i maintient un **vecteur** local
 $V_i = [c_0, \dots, c_{n-1}]$
 - (n est le nombre de processus)
- Chaque évènement e est associé à un vecteur
 $V(e) = [c_0, \dots, c_{n-1}]$
- c_j est le nombre d'évènements qui précèdent causalement e sur le processus p_j

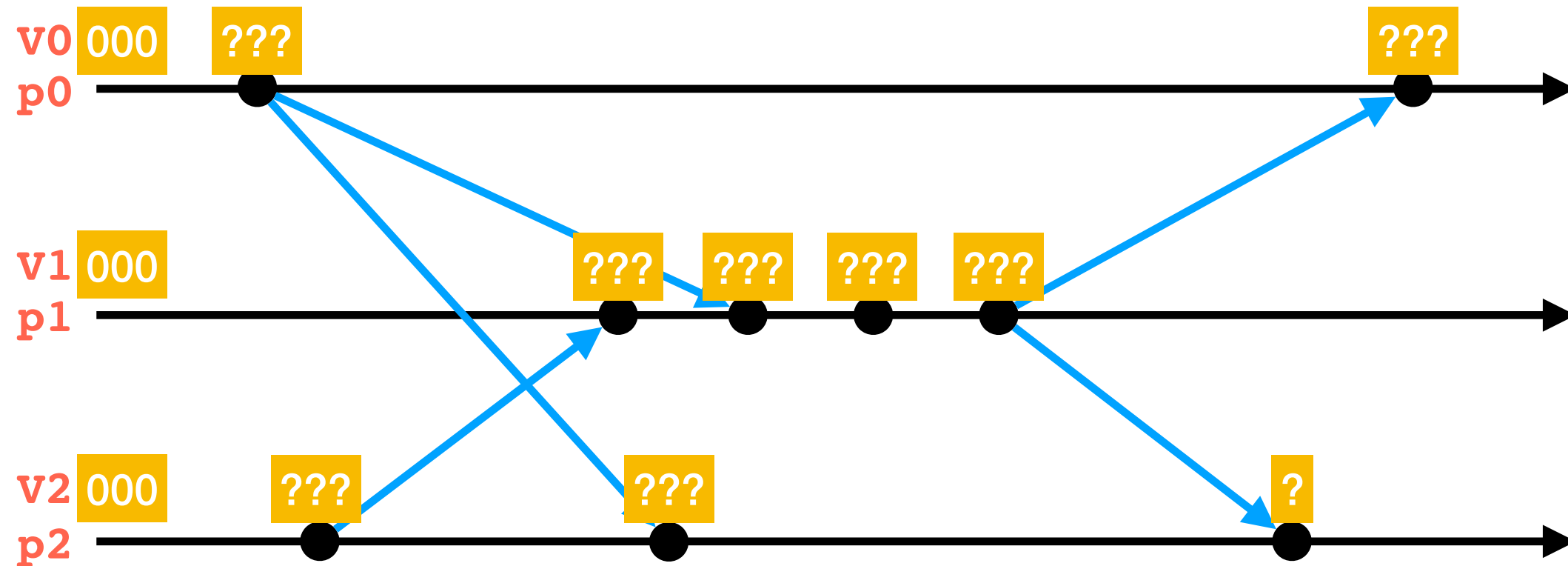
Mise à jour des vecteurs

- Initialement, pour chaque processus p_i , $V_i = [0, \dots, 0]$
 1. à chaque évènement **par** p_i , incrémenter $V_i[i]$
 $V_i[i] = V_i[i] + 1$
 2. **envoi** d'un message **m** incrémenter $V_i[i]$, puis envoi de **m** et V_i $V_i[i] = V_i[i] + 1$; **send**(**m**, V_i)
 3. **réception** de (**m**, V) : m.a.j V_i avec max. par composante de V_i et V , incrémenter $V_i[i]$
rcv(**m**, V) $V_i = \max(V_i, V)$; $V_i[i] = V_i[i] + 1$

Opérations sur les vecteurs

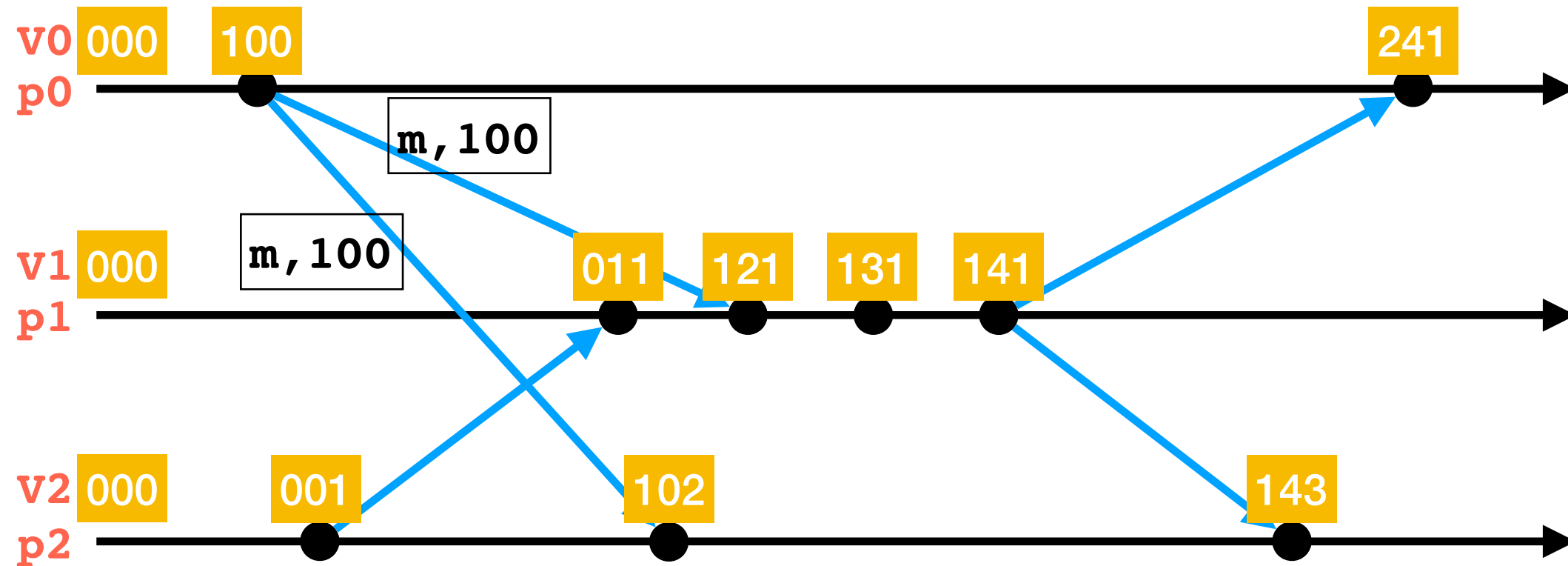
- $U = \max(V, W)$: pour chaque $i, 0 \leq i \leq n - 1$:
 $U[i] = \max(V[i], W[i])$
- $U \leq V$ si et seulement si pour chaque $i, 0 \leq i \leq n - 1$:
 $U[i] \leq V[i]$
- $[2, 2, 3] = \max([1, 2, 3], [2, 1, 3])$
- $[1, 2, 3] < [2, 3, 4]$, $[1, 2, 3]$ et $[2, 1, 3]$ sont incomparables

Horloge vectorielle



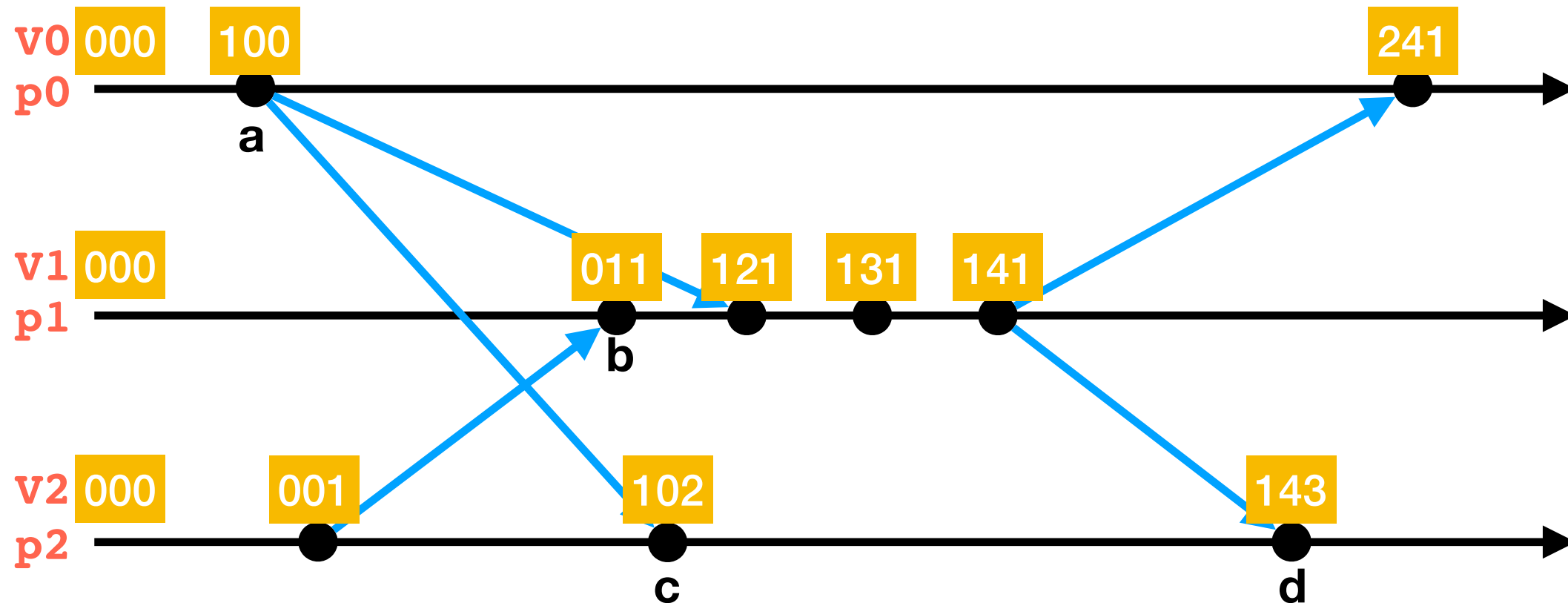
1. à chaque évènement interne par p_i : $V_i[i] = V_i[i] + 1$
2. envoi d'un message m : $V_i[i] = V_i[i] + 1$; $\text{send}(m, V_i)$
3. réception de (m, V) : $\text{rcv}(m, V)$ $V_i = \max(V_i, V)$; $V_i[i] = V_i[i] + 1$

Horloge vectorielle



1. à chaque évènement interne par p_i : $V_i[i] = V_i[i] + 1$
2. envoi d'un message m : $V_i[i] = V_i[i] + 1$; $\text{send}(m, V_i)$
3. réception de (m, V) : $\text{rcv}(m, V)$ $V_i = \max(V_i, V)$; $V_i[i] = V_i[i] + 1$

Horloge vectorielle

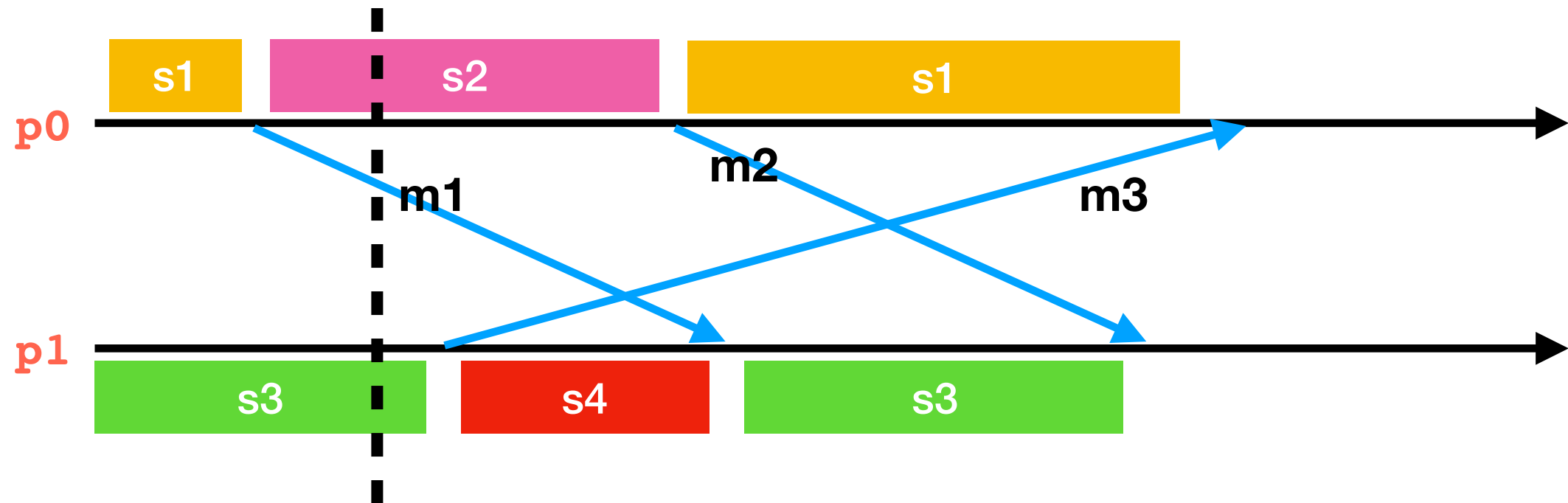


- événements indépendants ont des vecteurs incomparables
 - $\mathbf{b} \parallel \mathbf{c} : V(b) \not\leq V(c) \text{ et } V(c) \not\leq V(b)$
- événements causalement ordonnés ont des vecteurs ordonnés
 - $a \rightarrow d : V(a) < V(d)$

Etat global

- **Etat global** : états locaux des processus + états des canaux (messages en transit)
- **Cohérent** : ont pu se produire ensemble à un moment de l'exécution

état global



- état global
 - état p_0 : s_2
 - état p_1 : s_3
 - canal c_{01} : m_1
 - canal c_{10} : vide

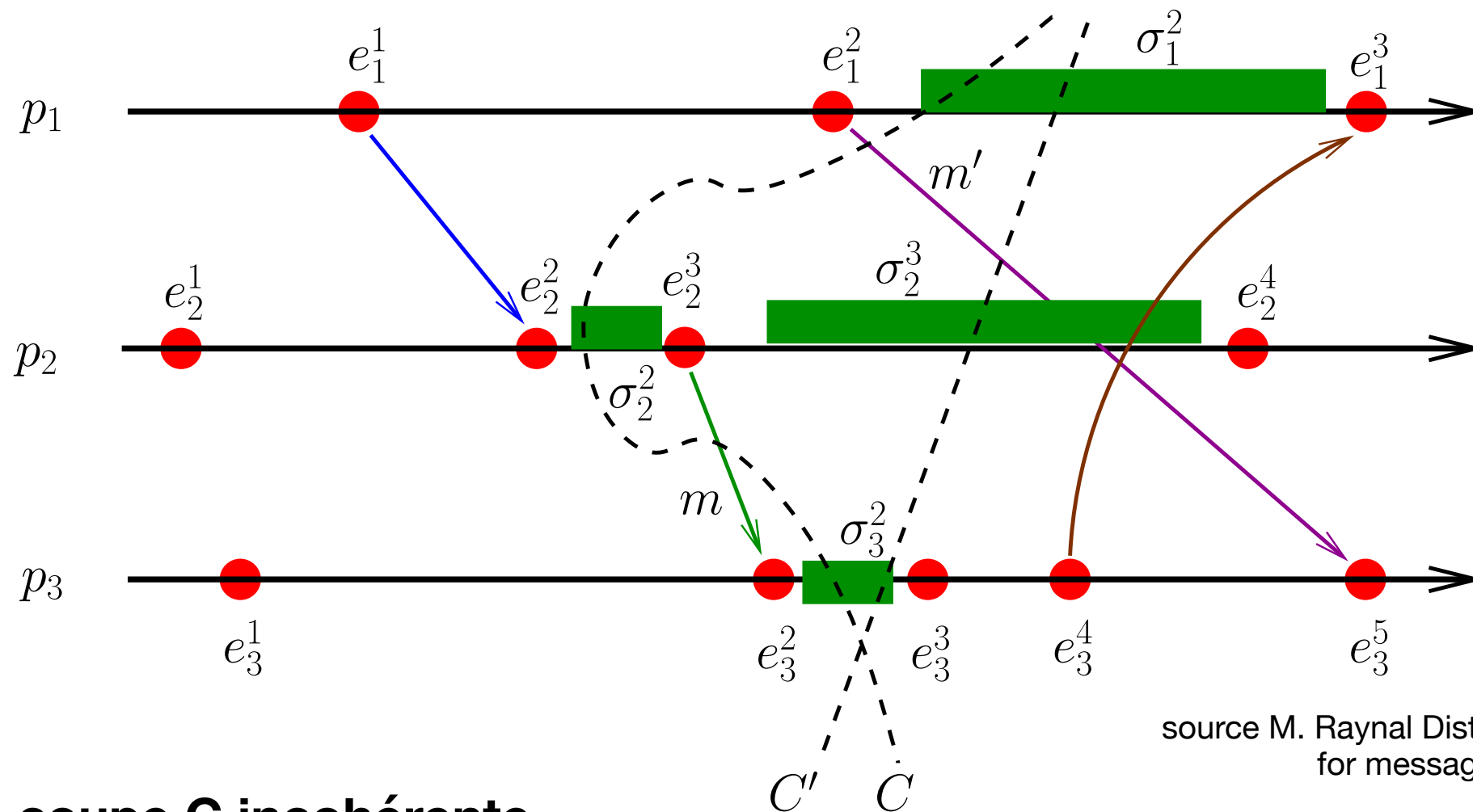
Application

- Point de reprise après panne (checkpointing)
- Garbage collection
- Débogage
- Détection de terminaison
- Vérification de propriétés globales (absence de deadlock)

Coupe cohérente

- **Coupe** : préfixes des événements de processus
- Coupe C **cohérente** si et seulement si :
$$\forall e, e' : e' \in C \wedge e \rightarrow e' \implies e' \in C$$
- I.e., pour chaque événement e dans C, le passé causal de e est aussi dans C

Coupes (in)cohérentes



source M. Raynal Distributed Algorithms
for message-passing systems

- coupe **C** incohérente
- coupe **C'** cohérente

Modèle

- **n** processus
- **canaux FIFO** first-in first-out : messages de p à q sont reçus dans l'ordre de leur envoi
- **asynchrone** ni perte, ni altération des messages mais délai non bornés

Algorithme de Chandy-Lamport

La machine initiatrice :

- sauvegarde son état
- envoie un **jeton** à tous les autres
- sauvegarde les messages reçus après l'envoi du jeton

Quand une machine reçoit un jeton pour la première fois :

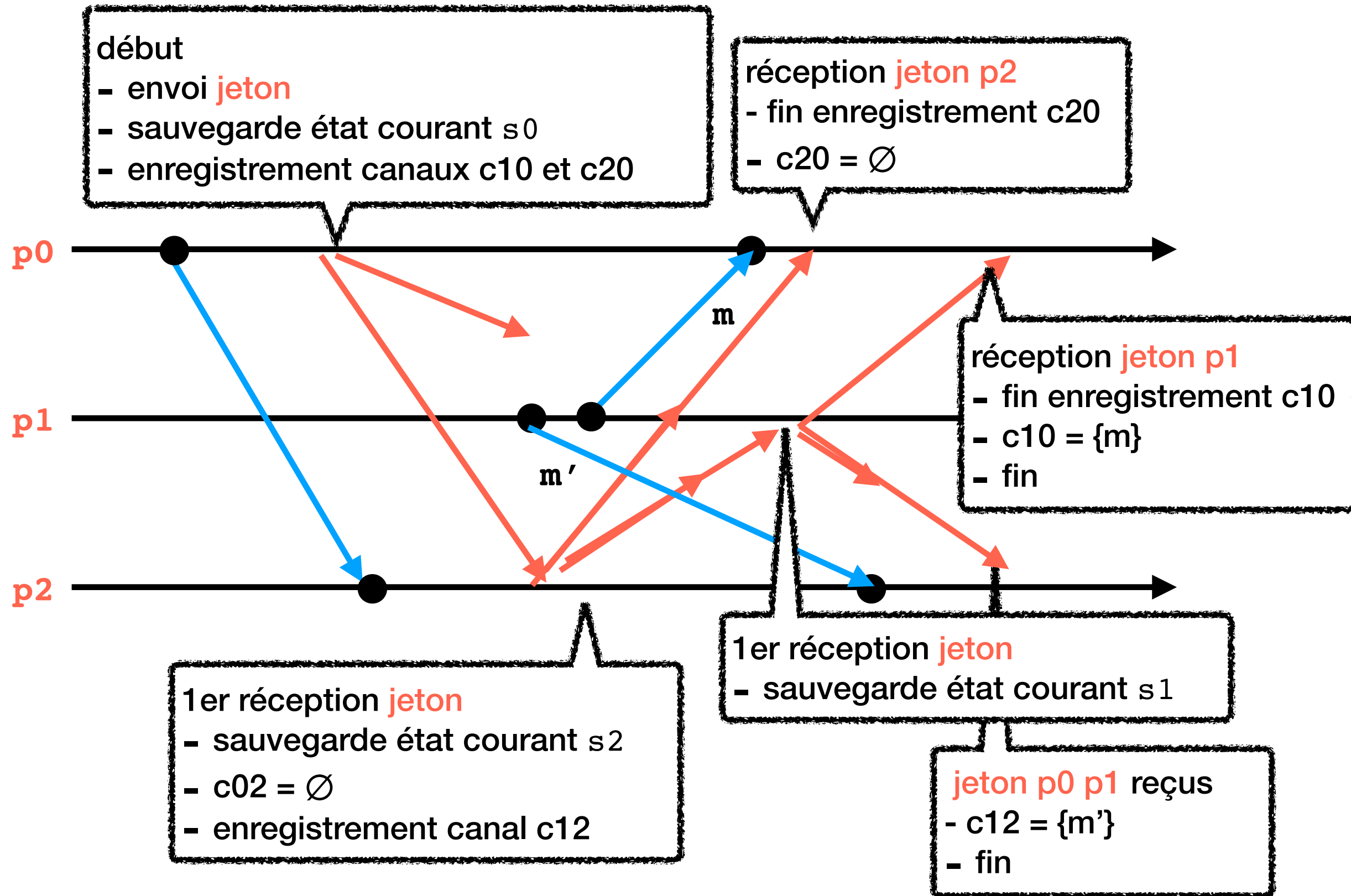
- sauvegarde son état
- envoie un **jeton** à tous les autres
- sauvegarde les messages reçus après l'envoi du jeton **sauf** ceux venant de machines dont elle a déjà reçu un jeton

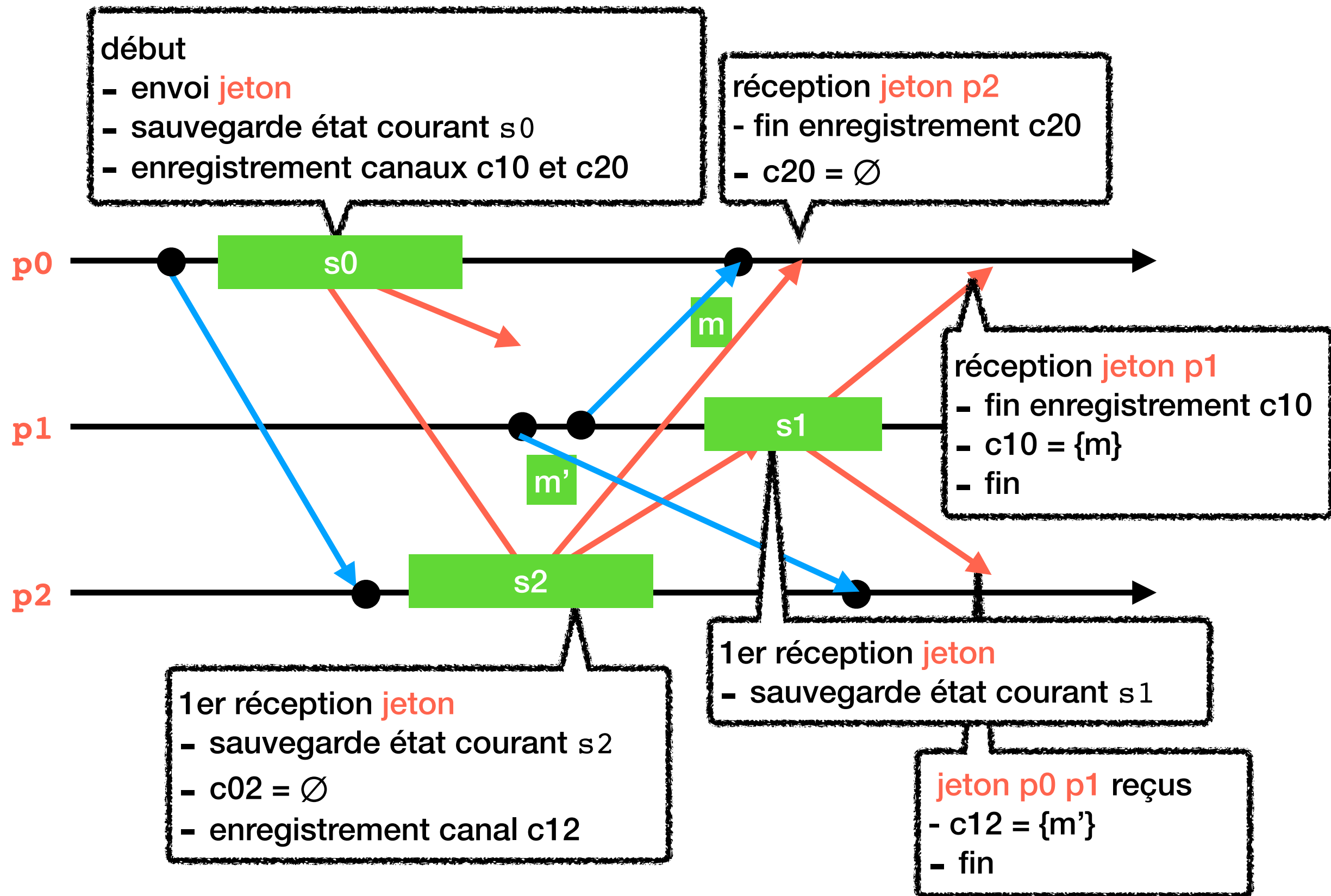
Quand une machine a reçu un jeton de tous les membres :

- Elle envoie son état et ses messages sauvegardés à l'initiateur

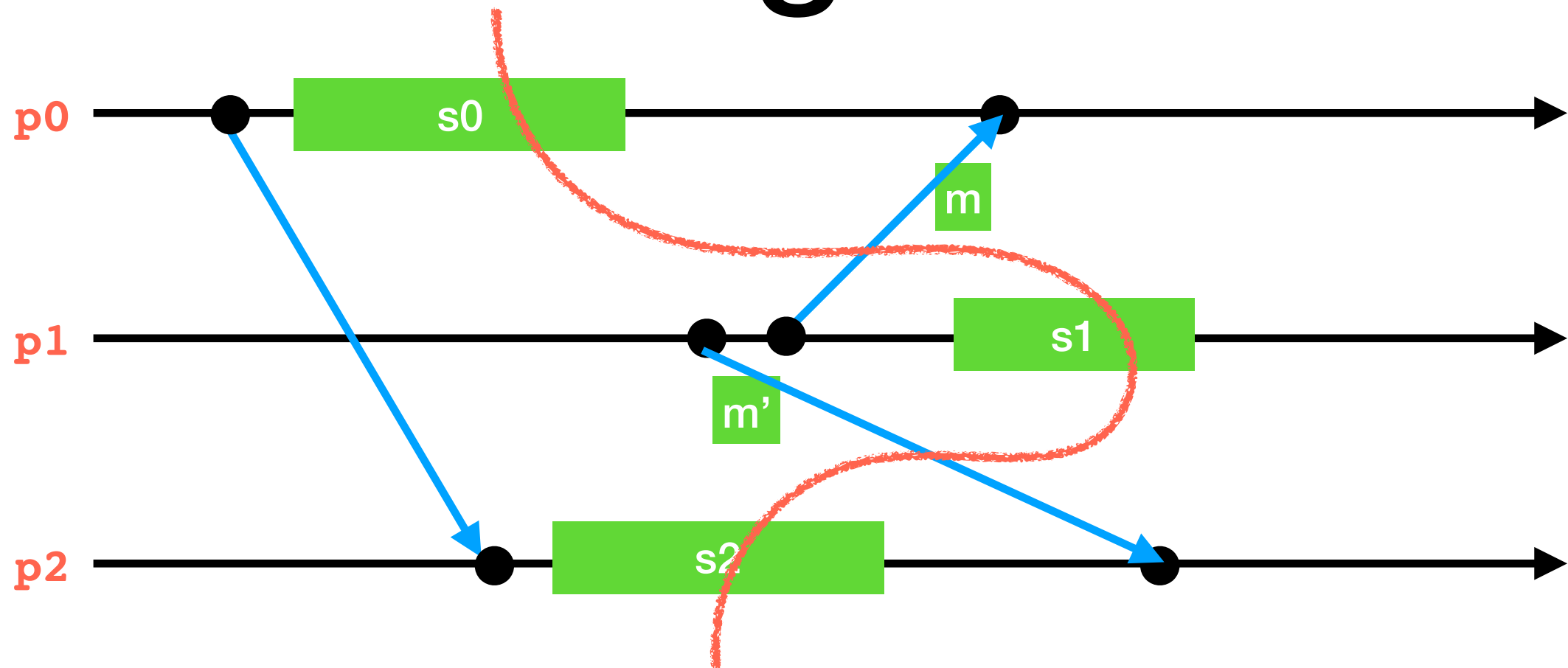
⇒ État global cohérent !

(en incluant les messages sauvegardés)



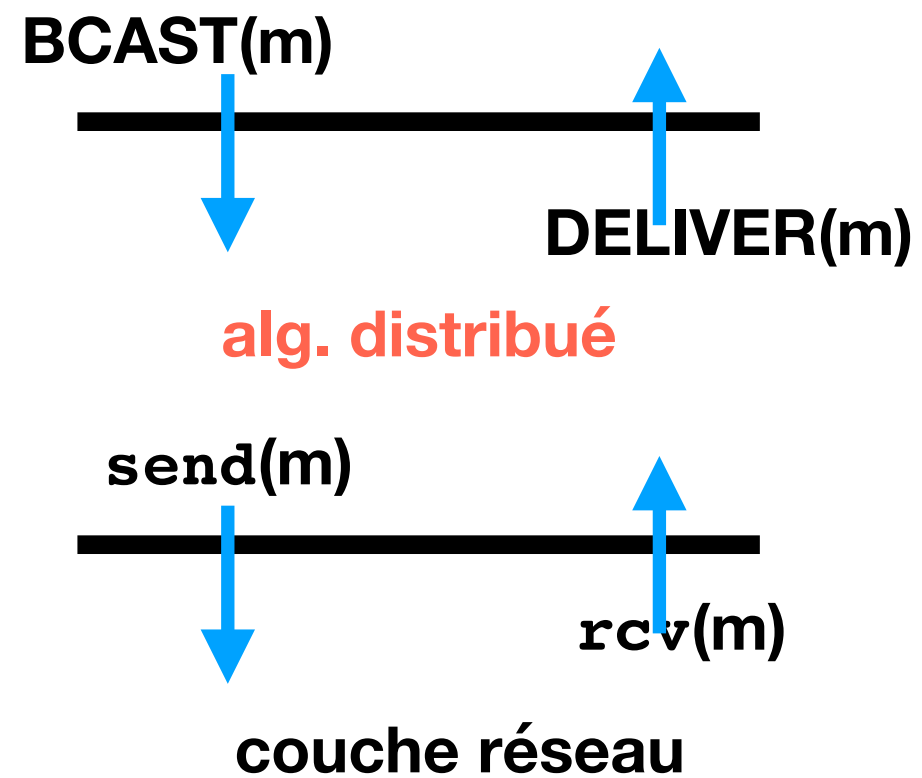


état global



- états processus $p_0 : s_0$, $p_1 : s_1$, $p_2 : s_2$
- états canaux : c_{01} :vide, c_{02} : vide, c_{10} : $\{m\}$, c_{12} : $\{m'\}$, c_{20} : vide, c_{21} :vide
- -> coupe cohérente

Diffusion (broadcast)



Total Order Broadcast

- Tout message diffusé (**BCAST(m)**) doit être livré (**DELIVER(m)**) par tous les processus
- Pour un processus et un message **m**, **m** est livré au plus une fois
- Les messages sont livrés dans **le même ordre** sur chaque processus

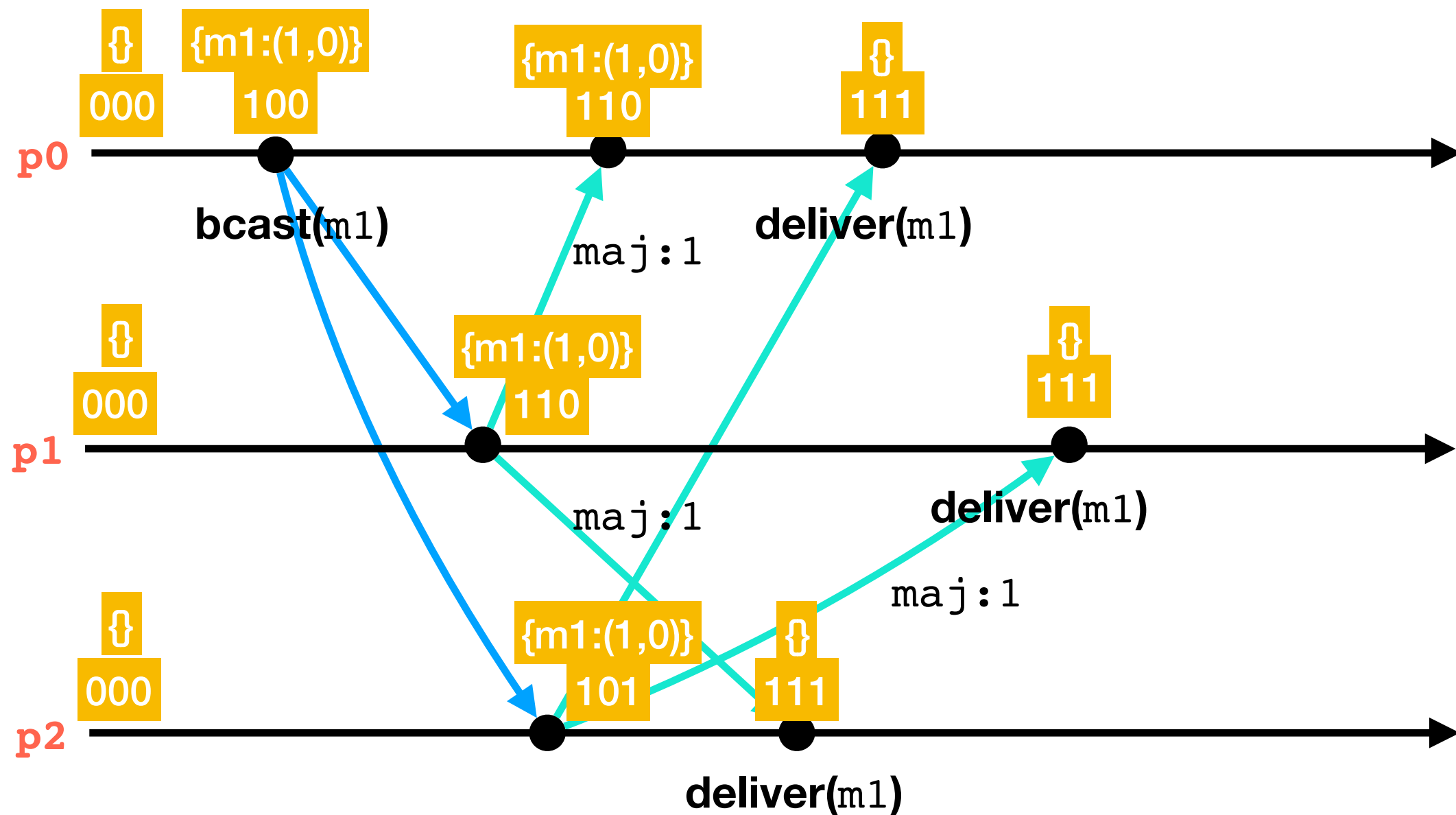
Modèle

- canaux asynchrones mais fiables
- chaque canal est FIFO
- n processus

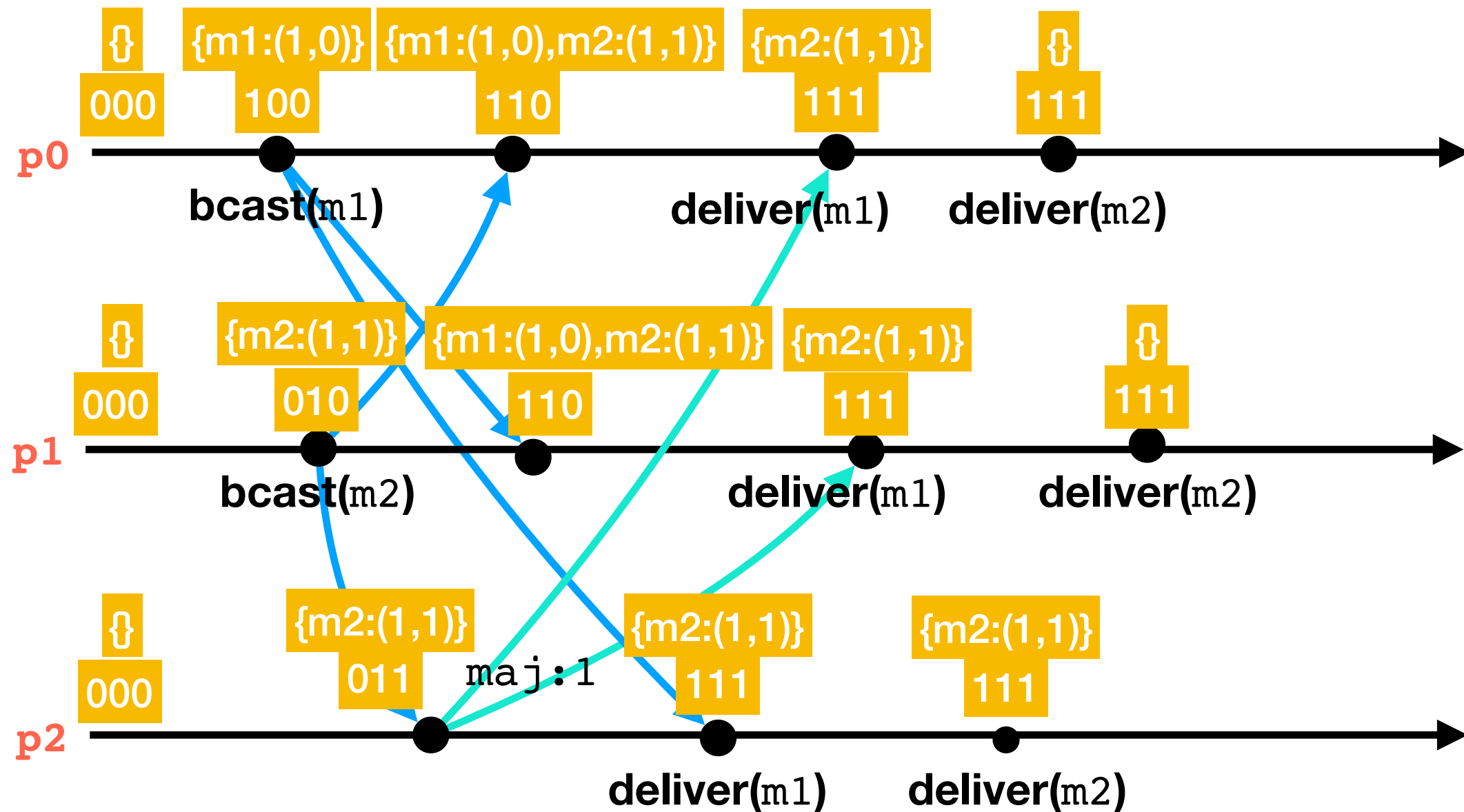
Algorithme

- **Variables locales** chaque processus p_i : $V[0..n-1]$ et *attente*
 - initialement $V = [0, \dots, 0]$ et *attente* = { }
- **BCAST(m) :**
 - $V[i]++$; ajouter $(m, V[i], i)$ à *attente* ; *send*($m, V[i], i$) à chaque processus $p_j, j \neq i$
- **réception de (m, v, j) :**
 - $V[j] = v$; ajouter (m, v, j) à *attente*.
 - si $v > V[i]$ alors $V[i] = v$; *send* (MAJ, $V[i], i$) à chaque processus $p_j, j \neq i$
- **réception de (MAJ, v, j) :** $V[j] = v$
- **DELIVER(m,j)** quand :
 - (m, v, j) a le plus petit couple (v, j) dans *attente* ET $v \leq V[k]$ pour tout k
 - enlever (m, v, j) de *attente*

Example



ordre total



Total order broadcast

- Ordre causal + ordre total
 - ordre causal : si **deliver(m)** \rightarrow **bcast(m')** alors **deliver(m)** \rightarrow **deliver(m')** pour tous les processus
- Canaux asynchrone, mais fiable et FIFO
- Non-tolérant aux pannes

Référence

