

Horloges logiques et datation des évènements.

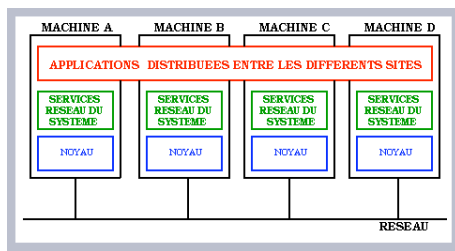
Johanne Cohen<sup>1</sup>

<sup>1</sup>LORIA/CNRS, Nancy, France.

inspiré du cours

- Jean-Marie Rifflet (<http://www.pps.jussieu.fr/~rifflet/>)
- Sacha Krakowiak (<http://proton.inrialpes.fr/~krakowia/>)

Système réparti



Être capable de raisonner sur un système réparti

- définir des prédicats (**état**)
- coordonner des instructions (**ordre**)

Système réparti

- Un système réparti est constitué de  $N$  sites (ou processus) communiquant par messages.
- Chacun de ces sites agit comme un automate : il réalise des opérations qui modifient son état.
  - Sur un même site : ordre des événements = ordre exécution des instructions
  - sur plusieurs sites : ordre des événements = ???

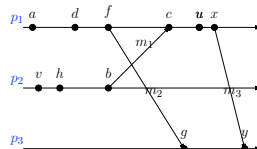
Difficile car...

- pas de mémoire commune, pas d'horloge commune
- Asynchronisme des communications et des traitements.

Modèle asynchrone

- Les processus ne communiquent que par des messages
- pas d'hypothèses
  - sur la durée de transmission des messages
  - sur la vitesse d'exécution des processus
- 3 types d'évènements

local	changement de l'état d'un processus
du à une communication	émission d'un message réception d'un message



Les communications

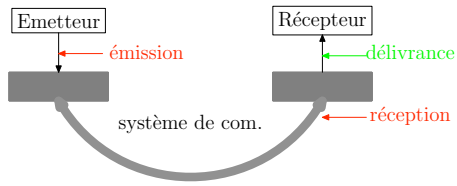
Échange d'information entre sites : envoi de messages.

Hypothèses

- Aucune perte de message (mais le temps de transmission peut être infini)
- Aucune erreur de transmission (un message arrive intacte)
- les canaux de communications peuvent être FIFO ou non.

## Délivrance

La communication entre les noeuds est en général assurée par une couche de communication spécifique utilisée par des entités de la couche supérieure pour émettre/recevoir des messages.



La **délivrance d'un message** : l'opération consistant à le rendre accessible aux applications clientes.

Par exemple : le protocole TCP dans la couche transport.

7/55

7

## Relation d'ordre sur les événements : précédence.

- Sur un même site : ordre des événements = ordre exécution des instructions
- sur plusieurs sites : **il faut définir un ordre des événements** tel que
  - il soit global
  - il doit être calculer entre deux évènements en fonction de l'information local.

Solution : Définition de la **causalité** [Lampport78].

8/55

8

## Plan

### Ordre Causal

Horloges logiques et datation des événements

- Horloges et estampilles scalaires
- Horloges et estampilles vectorielles
- Horloges et estampilles matricielles

Coupures cohérentes

- Définitions
- Protocole de Chandy-Lampport
- Coupures cohérentes et horloges vectorielles

9/55

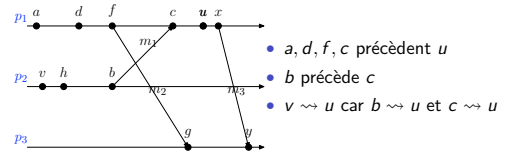
9

## Ordre causal

L'ordre causal est l'ordre partiel sur les événements :

Si  $a$  et  $b$  sont deux événements,  $a$  précède  $b$  ( $a \rightsquigarrow b$ ) si et seulement si l'une des trois conditions est vraie :

- $a$  et  $b$  ont lieu sur le même site avec  $a$  avant  $b$ .
- $a = Envoyer(< M >)$  et  $b = Recevoir(< M >)$  du même message.
- Il existe un évènements  $c$  tel que  $a \rightsquigarrow c$  et  $c \rightsquigarrow b$ .



10/55

10

## Ordre causal = ordre partiel

A un événement  $e$  trois ensembles d'événements sont associés :

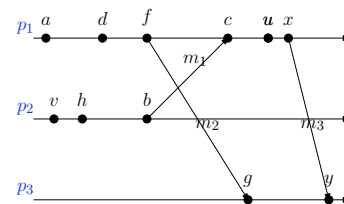
- $Passe(e)$  : ensemble des événements antérieurs à  $e$  dans l'ordre causal ( $e$  appartient à cet ensemble);
- $Futur(e)$  : ensemble des événements postérieurs à  $e$  dans l'ordre causal ( $e$  appartient à cet ensemble);
- $Concurrent(e)$  : ensemble des événements concurrents avec  $e$ .

Notation  $e \parallel e'$  : deux événements  $e$  et  $e'$  sont concurrents

11/55

11

## Exemple

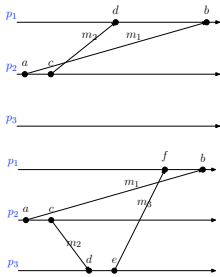


- $Passe(u) = \{a, b, c, f, v, h\}$
- $Futur(u) = \{x, y\}$
- $Concurrent(u) = \{g\}$

12/55

12

## Remarque sur certains points incohérents



Canal non FIFO :

- $c \rightsquigarrow d$
- $d \rightsquigarrow b$
- intérêt de recevoir  $m_1$  ?

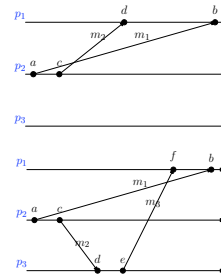
Autre scénario :

- $a \rightsquigarrow f$
- $f \rightsquigarrow b$
- intérêt de recevoir  $m_1$  ?

13/55

13

## Remarque sur certains points incohérents



Canal non FIFO :

- $c \rightsquigarrow d$
- $d \rightsquigarrow b$
- intérêt de recevoir  $m_1$  ?

Autre scénario :

- $a \rightsquigarrow f$
- $f \rightsquigarrow b$
- intérêt de recevoir  $m_1$  ?

Délivrer à l'application cliente de manière correcte

13/55

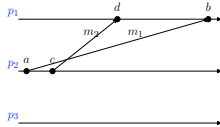
13

## Délivrance FIFO

### La délivrance FIFO

assure que si 2 messages sont envoyés successivement depuis un même site  $S_i$  vers un même destinataire  $S_j$ , alors le premier envoyé sera délivré sur le site  $S_j$  avant le second. Formellement, on a

$$send_i(m_1, j) \rightsquigarrow send_i(m_2, j) \text{ implique } del_j(m_1) \rightsquigarrow del_j(m_2)$$



14/55

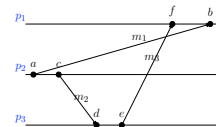
14

## Délivrance causale

### La délivrance causale

assure que si l'envoi du message  $m_1$  par le site  $S_i$  à destination du site  $S_k$  précède (causalement) l'envoi du message  $m_2$  par le site  $S_j$  à destination du site  $S_k$ , le message  $m_1$  sera délivré avant le message  $m_2$  sur le site  $S_k$ . Formellement, on a

$$send_i(m_1, k) \rightsquigarrow send_j(m_2, k) \text{ implique } del_k(m_1) \rightsquigarrow del_k(m_2)$$



15/55

15

## Plan

Ordre Causal

### Horloges logiques et datation des évènements

- Horloges et estampilles scalaires
- Horloges et estampilles vectorielles
- Horloges et estampilles matricielles

Coupsures cohérentes

- Définitions
- Protocole de Chandy-Lamport
- Coupsures cohérentes et horloges vectorielles

16/55

16

## Horloges scalaires

- Chaque site  $i$  gère un compteur (une horloge  $HL_i$ ) dont la valeur est un entier.
- chaque message  $m$  envoyé est estampillé (daté) ( $EL_m$ ) par la date de son évènement.

- Si un évènement est local, alors  $HL_i \leftarrow HL_i + 1$
- Si l'évènement est l'envoi d'un message  $m$  alors

$$\begin{cases} HL_i \leftarrow HL_i + 1 \\ EL_m \leftarrow HL_i \end{cases}$$

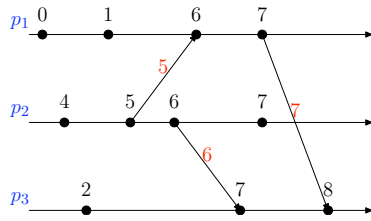
- Si l'évènement est la réception du message  $m$  alors

$$HL_i \leftarrow \max(HL_i, EL_m) + 1$$

17/55

17

### illustration



L'ordre des événements n'est pas un ordre strict : plusieurs événements peuvent porter la même valeur.

### Ordre strict si modification suivante

- l'estampille logique  $HL(e)$  de  $e$  sur le site  $i$  est un couple  $(HL_i, i)$ .
- L'ordre sur les estampilles est le suivant :

$$(HL_i, i) \prec (HL_j, j) \text{ si et seulement si } \begin{cases} HL_i < HL_j \\ \text{ou } HL_i < HL_j \wedge i < j \end{cases}$$

### Propriété : respect de l'ordre causal ?

#### Propriété

Si  $e \rightsquigarrow e'$ , alors  $HL(e) \prec HL(e')$

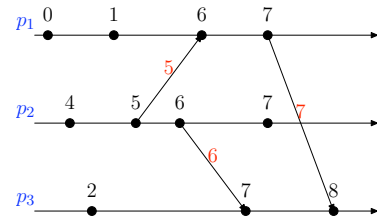
#### Démonstration.

Raisonnement par récurrence : sur la longueur  $n$  de la suite d'évènement reliant  $e$  à  $e'$

- si  $n = 0$  : alors  $e = e'$  et  $HL(e) \preceq HL(e')$
- si  $n > 1$  : soit  $e = e_0, e_1, \dots, e_n = e'$  une suite d'évènement tel que  $\forall i \in [0, \dots, n-1], e_i \rightarrow e_{i+1}$ 
  - par récurrence on a  $e \rightsquigarrow e_{n-1}$  et  $HL(e) \prec HL(e_{n-1})$
  - Comme  $e_{n-1} \rightarrow e_n$ , on a
    - Soit  $e_{n-1}, e_n$  appartiennent au même site  $HL(e_{n-1}) \prec HL(e_n)$
    - Soit il existe un message  $m$  tel que  $send(m) = e_{n-1}$  : par définition on a  $HL(e_{n-1}) \prec HL(e_n)$
  - Par transitivité de  $\preceq$ , on a  $HL(e) \prec HL(e_n)$

□

### illustration



L'ordre des événements n'est pas un ordre strict : plusieurs événements peuvent porter la même valeur.

### Ordre strict si modification suivante

- l'estampille logique  $HL(e)$  de  $e$  sur le site  $i$  est un couple  $(HL_i, i)$ .
- L'ordre sur les estampilles est le suivant :

$$(HL_i, i) \prec (HL_j, j) \text{ si et seulement si } \begin{cases} HL_i < HL_j \\ \text{ou } HL_i < HL_j \wedge i < j \end{cases}$$

### Propriété : respect de l'ordre causal ?

#### Propriété

Si  $e \rightsquigarrow e'$ , alors  $HL(e) \prec HL(e')$

#### Démonstration.

Raisonnement par récurrence : sur la longueur  $n$  de la suite d'évènement reliant  $e$  à  $e'$

- si  $n = 0$  : alors  $e = e'$  et  $HL(e) \preceq HL(e')$
- si  $n > 1$  : soit  $e = e_0, e_1, \dots, e_n = e'$  une suite d'évènement tel que  $\forall i \in [0, \dots, n-1], e_i \rightarrow e_{i+1}$ 
  - par récurrence on a  $e \rightsquigarrow e_{n-1}$  et  $HL(e) \prec HL(e_{n-1})$
  - Comme  $e_{n-1} \rightarrow e_n$ , on a
    - Soit  $e_{n-1}, e_n$  appartiennent au même site  $HL(e_{n-1}) \prec HL(e_n)$
    - Soit il existe un message  $m$  tel que  $send(m) = e_{n-1}$  : par définition on a  $HL(e_{n-1}) \prec HL(e_n)$
  - Par transitivité de  $\preceq$ , on a  $HL(e) \prec HL(e_n)$

□



## Propriété

### Propriété

- $e \rightsquigarrow e'$  si et seulement si  $EV_e \preceq EV_{e'}$
- $e \parallel e'$  si et seulement si  $EV_e \parallel EV_{e'}$

### Démonstration.

- Si  $e \rightsquigarrow e'$  alors,
  - $Passe(e) \subseteq Passe(e')$ ,
  - $\forall i, EV_e \preceq EV_{e'}$
  - $EV_e \subseteq EV_{e'}$
- Si  $e \parallel e'$  alors

□

30/55

30

## Propriété

### Propriété

- $e \rightsquigarrow e'$  si et seulement si  $EV_e \preceq EV_{e'}$
- $e \parallel e'$  si et seulement si  $EV_e \parallel EV_{e'}$

### Démonstration.

- Si  $e \rightsquigarrow e'$  alors,  $EV_e \preceq EV_{e'}$ 
  - $Passe(e) \subseteq Passe(e')$ ,
  - $\forall i, EV_e \preceq EV_{e'}$
  - $EV_e \subseteq EV_{e'}$
- Si  $e \parallel e'$  alors

□

30/55

30

## Suite de la preuve : si $e \parallel e'$ , alors

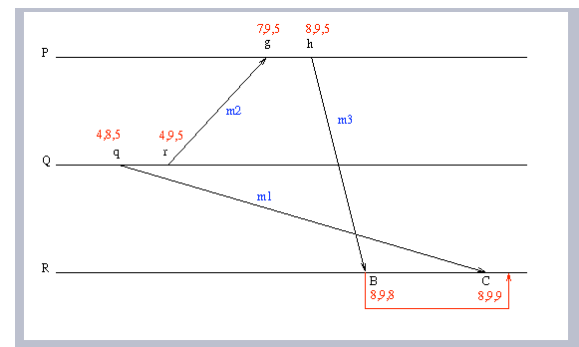
Soit  $S_j$  et  $S_j$  les sites appartenant à  $e$  et  $e'$

- si  $j = i$  alors,  $e$  et  $e'$  ne peuvent pas être concurrents.
- Donc  $j \neq i$
- Soit  $x(e, j)$  le plus grand élément du  $Passe_j(e)$
- on a  $x(e, j) \rightarrow^+ e'$  (car si  $e' \rightsquigarrow x(e, j)$ , alors  $e' \in Passe(e)$ )
- Donc :  $Passe_j(e) = Passe_j(x(e, j)) \subset Passe_j(e')$
- Par conséquent :  $EV_e[j] < EV_{e'}[j]$
- Utilisant le même raisonnement : on a  $EV_{e'}[i] < EV_e[i]$

31/55

31

## Cet ordre permet-il la délivrance causale ? NON



32/55

32

## Horloges et estampilles matricielles

- Chaque site  $i$  gère une horloge (notée  $HM_i$ ) matricielle ( $n \times n$ ) avec  $n$  le nombre de sites
- chaque message  $m$  envoyé est estampillé (daté) ( $EM_m$ ) par la valeur courante de  $HM_i$ .
- Que signifie  $HM_i[j, k]$ ?
  - $HM_i[j, k]$  = nbre de messages issus de  $p_j$  vers  $p_k$  dont  $p_i$  a connaissance.
  - $HM_i[j, j]$  : correspond au nbre d'évènements locaux du site  $j$ .

33/55

33

## Comment $HM_i$ est-elle modifiée ? (1/2)

- Évènement local à  $p_i$  :  $HM_i[i, i] \leftarrow HM_i[i, i] + 1$
- Émission de  $m$  vers  $p_j$  :
 
$$\begin{cases} HM_i[i, i] \leftarrow HM_i[i, i] + 1 \\ HM_i[i, j] \leftarrow HM_i[i, j] + 1 \\ m \text{ est estampillé : } EM_m = HM_i \end{cases}$$
- Réception de  $(m, EM_m)$  en provenance  $p_j$  :

34/55

34

### Comment $HM_i$ est-elle modifiée ? (1/2)

- Évènement local à  $p_i$  :  $HM_i[i, i] \leftarrow HM_i[i, i] + 1$
- Émission de  $m$  vers  $p_j$  :  $\begin{cases} HM_i[i, i] \leftarrow HM_i[i, i] + 1 \\ HM_i[i, j] \leftarrow HM_i[i, j] + 1 \\ m \text{ est estampillé : } EM_m = HM_i \end{cases}$
- Réception de  $(m, EM_m)$  en provenance  $p_j$  :  
Toutes les informations sont connues pour assurer la délivrance causale

34/55

34

### Comment $HM_i$ est-elle modifiée ? (2/2)

- Réception de  $(m, EM_m)$  en provenance  $p_j$  :
  - Quand un message peut-il être délivré ?
- Lors de la délivrance du message  $(m, EM_m)$  en provenance de  $p_j$ , mettre à jour l'horloge

35/55

35

### Comment $HM_i$ est-elle modifiée ? (2/2)

- Réception de  $(m, EM_m)$  en provenance  $p_j$  :
  - Quand un message peut-il être délivré ? quand tous les messages antérieurs à lui ont été délivrés, i.e :
    - Respect de l'ordre FIFO sur le canal  $j \rightarrow i$  :  $EM_m[j, i] = HM_i[j, i]$
    - Respect de l'ordre de réception :  $EM_m[k, i] = HM_i[k, i], \forall k \neq i$
  - Lors de la délivrance du message  $(m, EM_m)$  en provenance de  $p_j$ , mettre à jour l'horloge

35/55

35

### Comment $HM_i$ est-elle modifiée ? (2/2)

- Réception de  $(m, EM_m)$  en provenance  $p_j$  :
  - Quand un message peut-il être délivré ? quand tous les messages antérieurs à lui ont été délivrés, i.e :
    - Respect de l'ordre FIFO sur le canal  $j \rightarrow i$  :  $EM_m[j, i] = HM_i[j, i]$
    - Respect de l'ordre de réception :  $EM_m[k, i] = HM_i[k, i], \forall k \neq i$
  - Lors de la délivrance du message  $(m, EM_m)$  en provenance de  $p_j$ , mettre à jour l'horloge
    - $HM_i[i, i] \leftarrow HM_i[i, i] + 1$
    - $HM_i[i, j] \leftarrow HM_i[i, j] + 1$
    - pour tout  $k \neq i$ , pour tout  $\ell \neq i : HM_i[k, \ell] \leftarrow \max(HM_i[k, \ell], EM_m[k, \ell])$

35/55

35

### illustration (début)

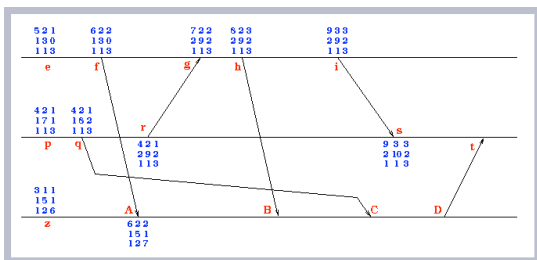


FIG.: datation avec l'horloge matricielle

36/55

36

### illustration (attention à la délivrance des messages)

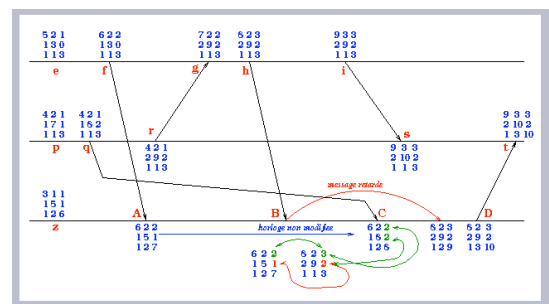


FIG.: datation avec l'horloge matricielle

37/55

37

## En résumé

- **Horloge scalaire :**  
 $HL_i$  : ce que  $p_i$  connaît du système (nbre d'évènements)
- **Horloge vectorielle :**  
 $HV_i[j]$  : ce que  $p_i$  connaît du site  $p_j$
- **Horloge matricielle :**  
 $HM_i[j, k]$  : ce que  $p_i$  connaît de la connaissance du site  $p_j$  sur le site  $p_k$

38/55

38

## Plan

Ordre Causal

Horloges logiques et datation des évènements  
 Horloges et estampilles scalaires  
 Horloges et estampilles vectorielles  
 Horloges et estampilles matricielles

Coupures cohérentes

Définitions  
 Protocole de Chandy-Lamport  
 Coupures cohérentes et horloges vectorielles

39/55

39

## État d'un système réparti

- Notion non-triviale : observateur universel ayant un accès instantané à l'ensemble de systèmes  
**Irréalisable**
- notion de **coupure** : capture l'état du dernier évènement avant "la photo" sur chaque site.

40/55

40

## Coupure

Définition

Une **coupure**  $C$  est l'ensemble des évènements tels que

Si  $e \in C$  et ( $e'$  précède localement  $e$ ), alors  $e' \in C$

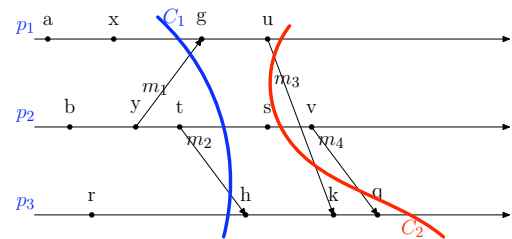


FIG.: la coupure  $C_1$  contient les évènements  $\{a, x, b, y, t, r\}$

41/55

41

## Coupures cohérentes

Définition

Une coupure **cohérente**  $C$  est une coupure cohérente si :

$$(e \in C) \text{ et } (e' \rightsquigarrow e) \implies e' \in C$$

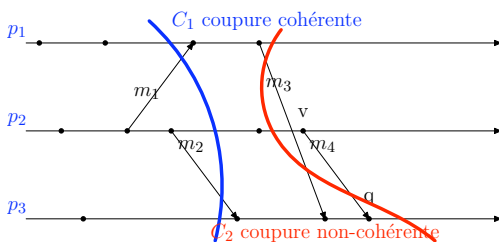


FIG.:  $C_2$  non-cohérente car  $v \notin C_2$  et  $q \in C_2$

42/55

42

## Propriétés sur les coupures cohérentes

Définition

Une coupure **cohérente**  $C$  est une coupure cohérente si :

$$(e \in C) \text{ et } (e' \rightsquigarrow e) \iff e' \in C$$

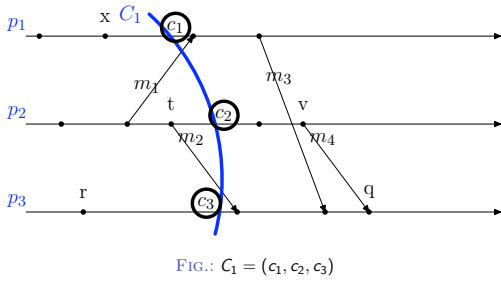
43/55

43



### Caractérisation des coupures cohérentes

Soit  $c_i$  les événements définissant la coupure :  $C = (c_1, \dots, c_n)$ .  
**Propriété :**  $(\forall i, j, c_j \parallel c_i) \Leftrightarrow C$  est une coupure cohérente



### Caractérisation des coupures cohérentes

Soit  $c_i$  les événements définissant la coupure :  $C = (c_1, \dots, c_n)$ .

**Propriété :**  $(\forall i, j, c_j \parallel c_i) \Leftrightarrow C$  est une coupure cohérente

**Démonstration.**

- Soit  $C$  est une coupure cohérente. Par contradiction, supposons  $\exists i \neq j, c_i \rightsquigarrow c_j$ .
  - $c_i \rightsquigarrow a_i \rightsquigarrow a_j \rightsquigarrow c_j$  (message entre les deux sites)
  - Comme  $a_j \rightsquigarrow c_j$ , on a  $a_j \in C$ .
  - Ceci est contradictoire avec le fait que  $c_j \rightsquigarrow a_j$
- Soit  $C = (c_1, \dots, c_n)$  avec  $\forall i, j, c_j \parallel c_i$ .
  - Soit  $a_i \rightsquigarrow c_i$ , et  $a_j \rightsquigarrow a_i$
  - Alors  $a_j \rightsquigarrow c_i$ , et  $a_j \in C$
  - Supposons que  $c_j \rightsquigarrow a_j$  : on aurait  $c_j \rightsquigarrow a_j \rightsquigarrow a_i \rightsquigarrow c_i$  (Contradictoire avec le fait que  $c_j \parallel c_i$ ).
  - Donc  $C$  est une coupure cohérente.

□

### Enregistrement l'état d'un système

- **Motivation :**
  - Observation, détection de propriétés.
  - reprise en cas de panne
- **Contraintes :**
  - L' état enregistré doit être cohérent
  - Le coût de l'enregistrement doit être raisonnable
- **Difficulté :**
  - opérations locales /propriété globale

### Protocole de Chandy-Lamport [1985]

**Hypothèses :**

- les canaux de communication entre processus sont FIFO
- Un seul processus décide de lancer la procédure d'enregistrement

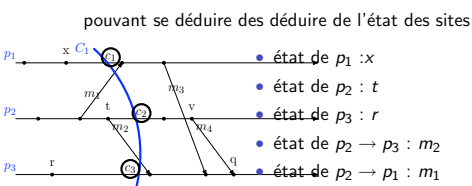
**Objectif :**

enregistrement d'un état cohérent en un temps fini sous quelle forme ?

### Définition d'un état.

L'état d'un système de réparti est

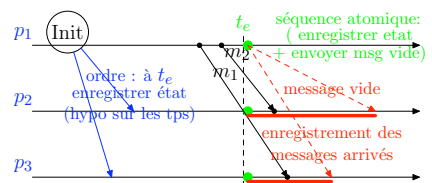
- l'ensemble des états de chacun de ses sites
- et l'ensemble des états de chacun de ses canaux de communications



### Première version : (1/2)

**Hypothèses :**

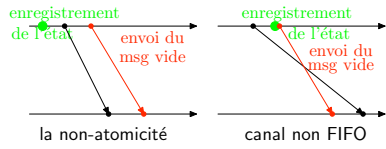
- horloge globale  $HR$  (temps réel).
- délai de transmission et rapport des vitesses des sites bornés.



### Première version :(2/2)

Propriétés : la coupure  $C$  défini au tps  $t_e$  est une coupure cohérente

- si  $e \in C$  et si  $e' \rightsquigarrow e$  (localement) alors  $HR(e') < HR(e)$  et  $e' \in C$
- L'état du canal est correct du
  - au fait de l'atomicité des instructions enregistrement état ; envoi message vide
  - et du fait que le canal de communication est FIFO

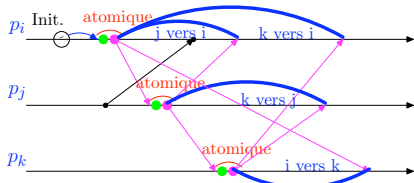


### Deuxième version :(1/3)

Algorithme sur le site  $p_i$  :

1. lière réception du marqueur depuis  $p_j$ 
  - enregistrer état de  $p_i$
  - envoyer le marqueur à tous ses voisins atomicité
  - état du canal  $j \rightarrow i := \emptyset$
  - enregistrer les msg sur les canaux entrants
2. Réception suivante du marqueur depuis  $p_k$  :
  - état du canal  $k \rightarrow i :=$  msg reçus depuis l'enregistrement.
  - ne plus enregistrer les msg provenant de  $p_k$

### Deuxième version :(2/3)



### Deuxième version :(3/3)

Propriétés : la coupure  $C$  ainsi calculée est une coupure cohérente  
Démonstration.

En exercice ?

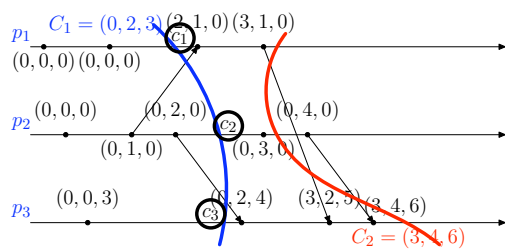
□

### Coupures cohérentes et horloges vectorielles

- La date de la coupure  $C = (c_1, \dots, c_n)$  est définie par

$$VH(C) = \sup(VH(c_1), \dots, VH(c_n))$$

- Illustration

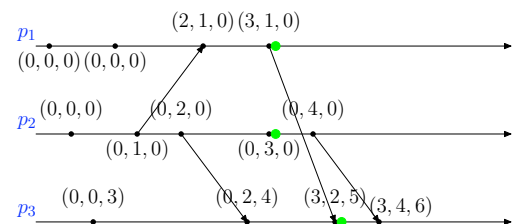


### Condition pour les coupures soient cohérentes.

Propriété :  $C$  cohérente  $\Leftrightarrow VH(C) = (VH(c_1)[1], \dots, VH(c_n)[n])$

Démonstration.

- Soit  $C$  cohérente.  $\forall i \neq j, VH(c_i)[i] \geq VH(c_j)[i]$ .
  - $VH(c_i)[i]$  modifié par un évènement local ou d'un msg passé



## Suite de la démonstration

- Soit  $C = (c_1, \dots, c_n)$  non-cohérente.
  - $\exists i, j$  tels qu'un message  $m$  de  $p_i$  émis après  $C$  a été reçu par  $p_j$  avant  $C$ . Soit  $EV(m)$  son estampille
  - Alors  $VH(c_i)[i] < EV(m)[i] \leq VH(c_j)[j]$
  - Donc  $(VH(c_1)[1], \dots, VH(c_n)[n]) < VH(C)$

