ETAT GLOBAL

41403 – Algorithmique Répartie

UPMC - 2016-2017



Plan

- . Introduction
 - 1. Motivation
 - 2. Représentation et cohérence
- Exemple d'état global
- M. Algorithme de snapshot de Chandy et Lamport 85

Plan

- . Introduction
- Exemple d'état global
 - 1. Définition du système étudié
 - 2. Graphe des états accessibles
 - 3. Remarques
- III. Algorithme de snapshot de Chandy et Lamport 85

Plan

- . Introduction
- Exemple d'état global
- III. Algorithme de snapshot de Chandy et Lamport 85
 - 1. Principes
 - 2. Analyse
 - 3. Algorithme
 - 4. Exemples
 - 5. Conclusion

I. Introduction

- 1. Motivation
- 2. Représentation et cohérence

Calcul de l'état global d'un système réparti Dinstantané de l'état du système capture d'un moment donné de son exécution

État global d'un système réparti à un instant t
Union des états de chacun des processus du système à t
ET
Union des états de chaque canal de communication à t

Pas de solution simple

Collecte distante de données/états locaux

Pas de temps global

Intrusivité de la collecte/observation doit être minimale

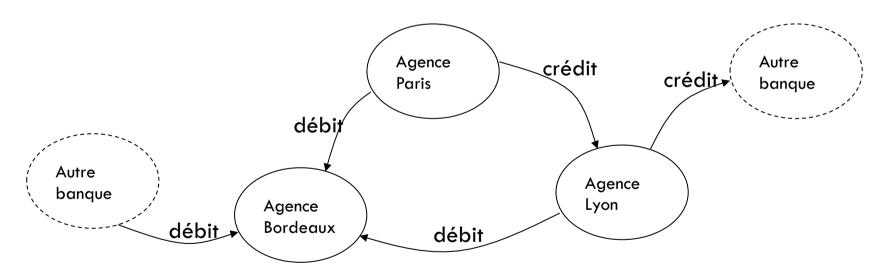
Plusieurs utilisations possibles:

- 1. Exécution d'un algorithme centralisé sur des données distantes
- 2. Détection d'états globaux stables (eg. interblocage, terminaison)
- 3. Ramasse-miette distribué
- 4. Tolérance aux fautes

Exécution d'un algorithme centralisé sur des données distantes

s'affranchir du caractère réparti du système.

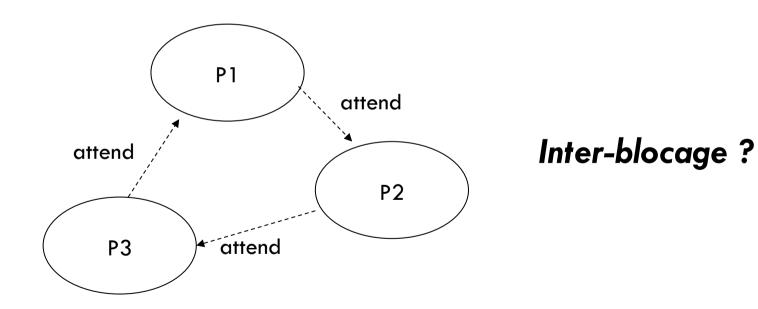
eg. calcul des avoirs d'une banque à partir de ceux de chacune de ses agences.



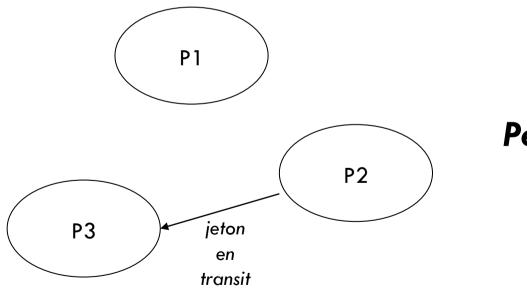
$$Avoirs_{banque} = Avoirs_{agence\ Paris} + Avoirs_{agence\ Lyon} + Avoirs_{agence\ Bordeaux}$$

41403 – Etat global

Détection d'états globaux stables

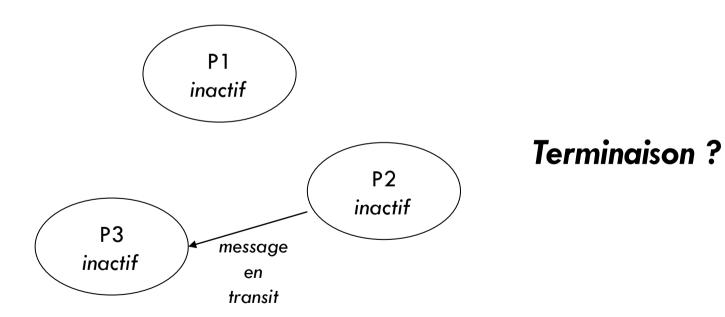


Détection d'états globaux stables



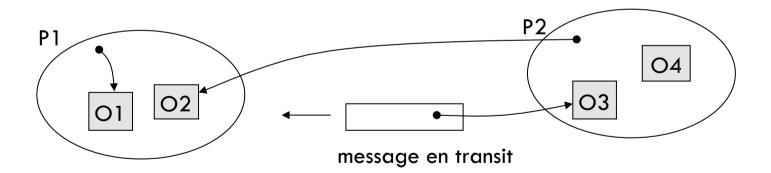
Perte du jeton?

Détection d'états globaux stables



Ramasse-miette distribué

élimination des objets inutilisables conservés en mémoire détection des références locales/distantes objet sans référence => suppression



04 n'est plus référencé dans le système réparti

03 n'est plus référencé que dans le message en transit (pour l'instant)

Tolérance aux pannes

Sauvegarde de l'état global du système => points de reprise (checkpoints) En cas de panne d'un des sites, procédure de recouvrement :

- 1. chaque site revient à son checkpoint (rollback)
- 2. reprise de l'exécution à partir de cet état.

Problème: garantir que le recouvrement est cohérent (duplication, causalité)

Etat local d'un processus Pi

Valeur des variables locales de P_i

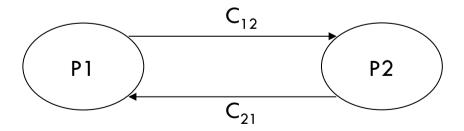
Déterminisme : état local de $P_i = f$ (état initial de P_i , succession des événements sur P_i)

Etat d'un canal de communication C_{ij} (entre P_i et P_j)

Ensemble des messages en transit entre P_i et P_i

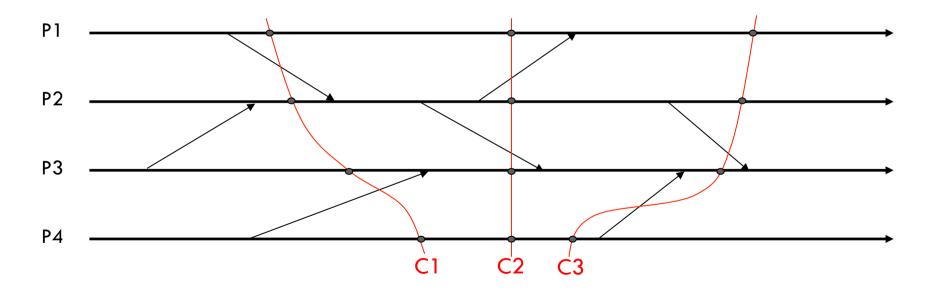
Canal FIFO => ensemble ordonné

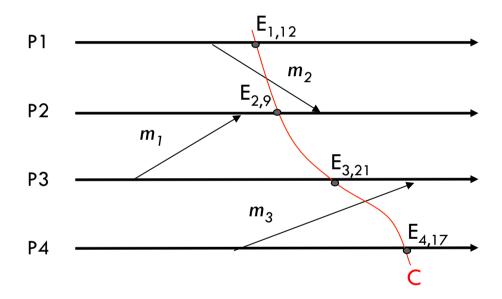
Canaux sont considérés unidirectionnels (C_{ij}, de P_i vers P_j, et C_{ji} de P_i vers P_i)



Etat global d'un système réparti

Union des états de l'ensemble des P_i et des C_{ij} Un état global peut être représenté par une courbe (*coupure*) sur le diagramme temporel

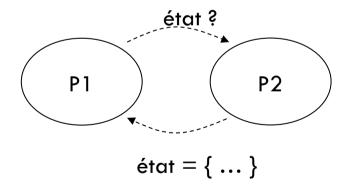


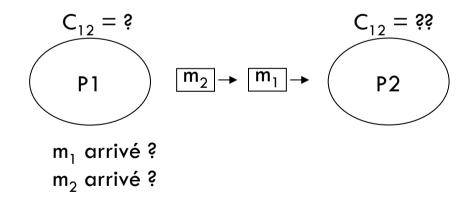


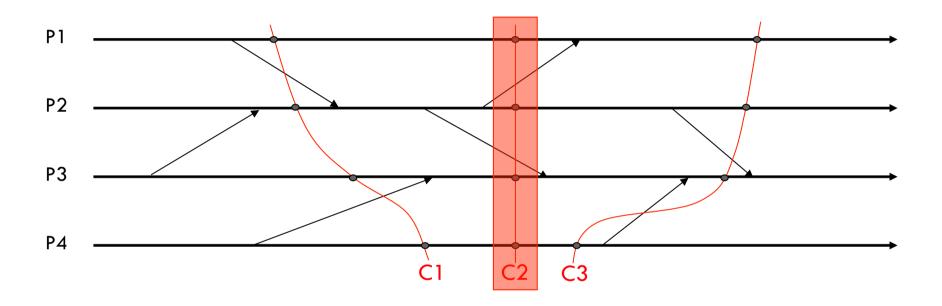
$$\mathbf{E}_{C}(SR) = \mathbf{E}(P_{1}) \cup \mathbf{E}(P_{2}) \cup \mathbf{E}(P_{3}) \cup \mathbf{E}(P_{4}) \cup \mathbf{E}(C_{12}) \cup \mathbf{E}(C_{43}) \\
= \mathbf{E}_{1,12} \cup \mathbf{E}_{2,9} \cup \mathbf{E}_{3,21} \cup \mathbf{E}_{4,17} \cup \{\mathbf{m}_{2}\} \cup \{\mathbf{m}_{3}\}$$

Problème de la collecte des états

 $E(P_i)$ n'est directement et immédiatement observable que sur P_i $E(C_{ij})$ n'est jamais directement observable, ni sur P_i , ni sur P_i







Problème de la cohérence des états collectés

Solution idéale : figer l'état de tous les P_i et C_{ij} au même instant absolu

Mais pas de temps global (ie. pas de référence temporelle commune)

Solution pragmatique : assouplir le critère de cohérence sans rendre l'état global inexploitable

=> Relation de précédence causale

Relation de précédence causale

Relation d'ordre entre événements dans le temps

e précède causalement e' est noté e → e'

2 événements d'une même exécution ne sont pas forcément liés causalement

$$\begin{array}{ccc} a) & & P_i^x & & P_i^y \\ \hline & & & & \\ \end{array}$$

devient
$$P_i^x \rightarrow P_i^y$$

$$\begin{array}{c} P_{i}^{x} \\ \hline \\ P_{j}^{x} \end{array}$$

$$\textit{devient} \quad P_i^{\ x} \ \longrightarrow \ P_j^{\ x}$$

Coupure cohérente

Définit un passé fermé pour la relation de précédence causale Pour tout événement survenu avant la coupure,

les événements qui le précèdent causalement sont également survenus avant la coupure

Formellement

C coupure cohérente $\Leftrightarrow \forall e \in Passé(C), e' \rightarrow e \Rightarrow e' \in Passé(C)$

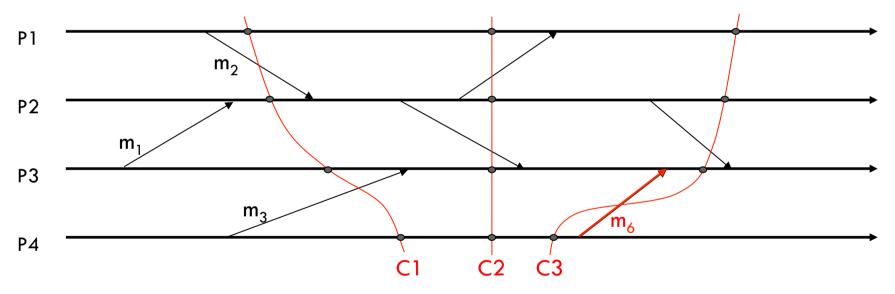
Intuitivement

La cause d'un événement ne peut être dans son futur (de la coupure).

Par définition

un état global cohérent est un état global obtenu selon une coupure cohérente

Pour déterminer si une coupure est cohérente, il suffit d'examiner les messages en transit. Aucun message ne doit aller du futur vers le passé.



C1 et C2 sont cohérentes

C3 **n'est pas** cohérente : réception $(m_6) \in Passé(C3)$ et émission $(m_6) \notin Passé(C3)$

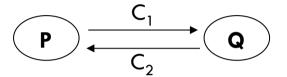
II. Exemple d'état global

- Définition du système étudié
- 2. Graphe des états accessible
- 3. Remarques

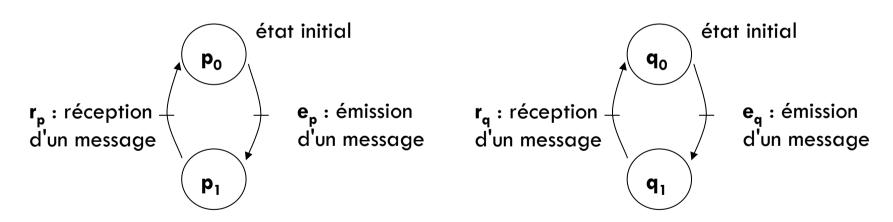
II.1 Définition du système étudié

Définissons le système réparti SR :

- 2 processus P et Q
- 2 canaux unidirectionnels C_1 (P vers Q) et C_2 (Q vers P)



Chacun des processus P et Q envoie puis recoit un message, alternativement. Automates des processus :



Etat global du système : quadruplet <état local P, état C1, état C2, état local Q>

II.2 Graphe des états accessibles

Définition

Ensemble des états que peut prendre le système réparti

Chaque nœud du graphe correspond à un état global du système réparti

Chaque arête correspond à une transition d'un état à un autre suite à un événement

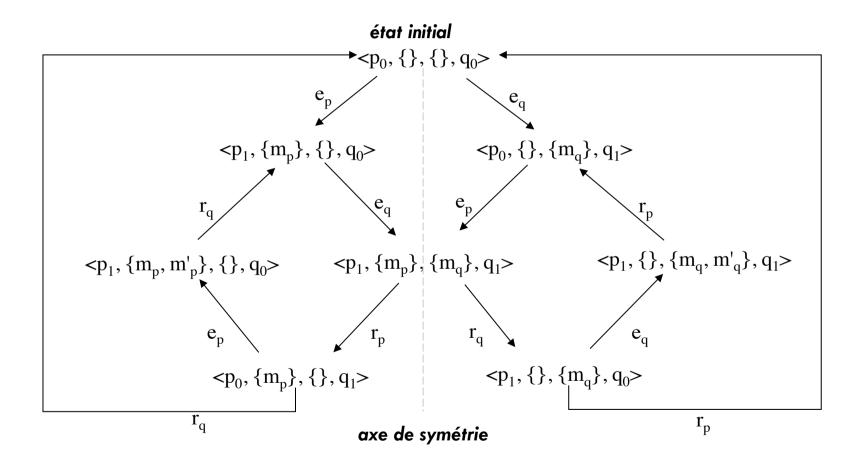
Une exécution particulière du système

Un cheminement dans ce graphe à partir de l'état initial.

Dans SR

un nœud est un quadruplet <état local P, état C1, état C2, état local Q> les arêtes sont e_p , r_p , e_q ou r_q

II.2 Graphe des états accessibles



II.3 Remarques

Le graphe des états accessibles permet :

- de détecter les deadlocks état puits : état sans arête sortante
- 2. de dimensionner les canaux visibilité du nombre maximal de messages dans chaque canal eg. sur SR, au plus deux messages dans les canaux C1 et C2

Chacun des états du graphe est cohérent : il correspond à un état réellement accessible.

Si le système réparti est symétrique, son graphe l'est également.

Pb : la taille du graphe croît **très** rapidement eg. bien que SR soit simple, son graphe comprend déjà 8 états

III. Snapshot algorithm (Chandy & Lamport 85)

- 1. Principes
- 2. Analyse
- 3. Algorithme
- 4. Exemples
- 5. Conclusion

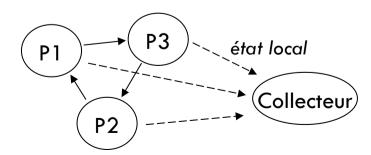
III.1 Principes

Objectif

Calculer **dynamiquement** un **état global cohérent** d'un système réparti en cours d'exécution Minimiser la complexité et le caractère intrusif de la solution

Modèle système

N processus répartis sur un **réseau fortement connexe**Canaux de communication **fiables** et **FIFO**Un processus (collecteur) désigné à l'avance et connu de tous les autres **Aucune panne** de processus au cours du déroulement de l'algorithme

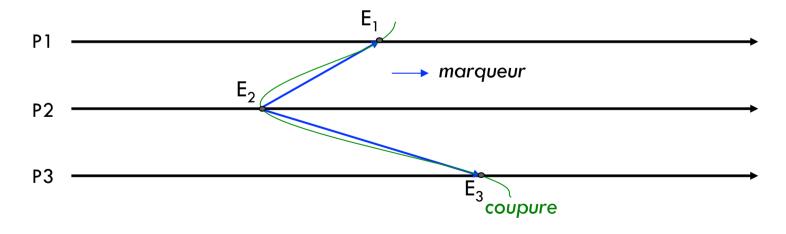


On part de l'algorithme le plus simple pour lister les problèmes à résoudre. Un processus déclenche le calcul d'un état global

il sauvegarde sont état local puis demande aux autres pcs d'en faire autant par un message spécial (marqueur)

2 propriétés à assurer :

- 1. la cohérence de la coupure,
- 2. et la complétude l'état sauvegardé pour les canaux de communication.

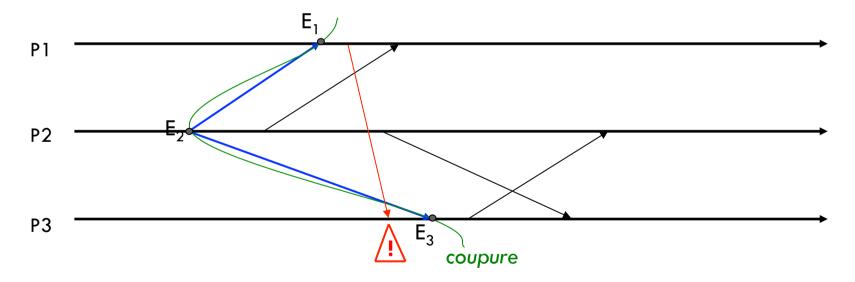


41403 – Etat global

Cohérence

Un message reçu avant la coupure doit être émis avant la coupure

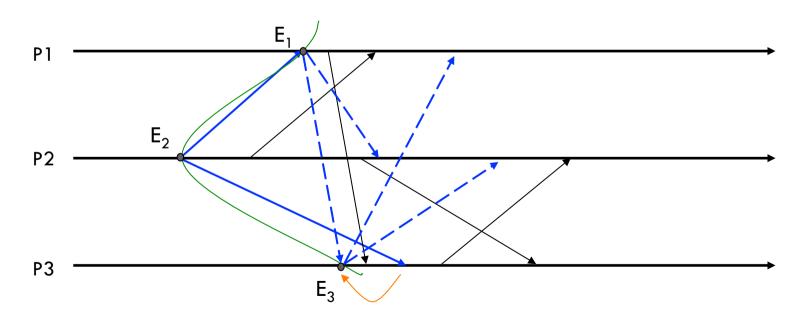
- de P₂ vers P₁ et P₃: pas de problème puisque les canaux sont FIFO
- de P₃ ou P₁ vers P₂ : pas de problème en raison de la causalité
- aucune garantie pour les autres couples de processus



Solution pour la cohérence : propager le marqueur.

Chaque sauvegarde provoque la diffusion du marqueur.

Le premier marqueur reçu provoque une sauvegarde, mais pas sur les suivants.



41403 – Etat global

Complétude

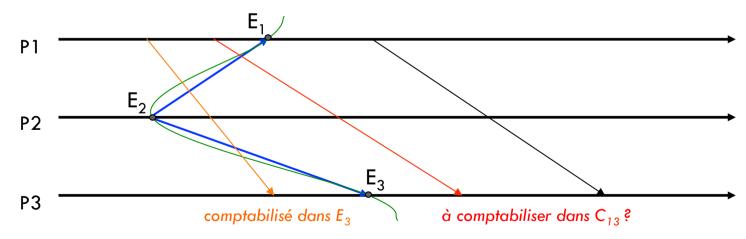
Tout message émis **avant** une sauvegarde doit être comptabilisé

Soit dans la sauvegarde du destinataire en étant reçu **avant** cette sauvegarde

Soit dans l'état du canal parce qu'il a été reçu après cette sauvegarde

Problème

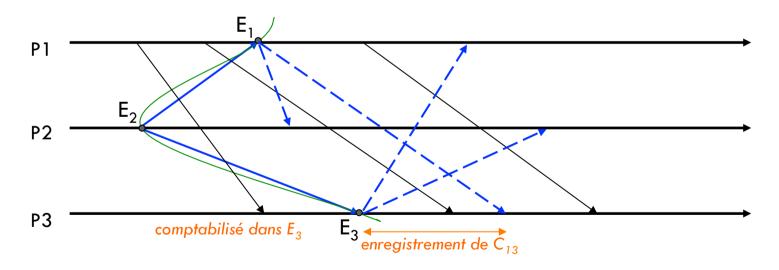
quand débuter et arrêter l'enregistrement des messages en transit ? Quels sont les messages émis qui ont été reçus avant la sauvegarde du destinataire ?



Solution pour la complétude : partir des messages reçus

Utiliser les marqueurs propagés

Message en transit : reçu pendant la phase de sauvegarde (ie. entre le premier marqueur et celui associé au canal)



Note: C₂₃ est par définition vide car P₃ reçoit son premier marqueur du collecteur (P₂)

Résumé

Utilisation de messages marqueurs

messages de contrôle, orthogonaux aux messages applicatifs (coloration)

Mise en œuvre de 2 règles

- 1. Règle d'envoi du marqueur sauvegarde de l'état local $P_i =>$ envoi du marqueur sur tous les canaux sortants C_{ix}
- 2. Règle de réception du marqueur sur réception du premier marqueur, P_i
 - a. sauvegarde son état
 - b. débute l'enregistrement des messages reçus sur chacun de ses canaux
 - c. clôt l'enregistrement sur réception du marqueur associé à chaque canal

Le processus initiateur agit comme sur réception d'un marqueur (fictif).

III.3 Algorithme

```
Variables locales :  enreg_i = faux; \\ etatLocal_i = \varnothing; \\ etatCanal[N]_i = \{\varnothing, \varnothing, \ldots\}; \\ marqRecu[N]_i = \{faux, faux, \ldots\};
```

```
\begin{split} & \textbf{sauvegarder()}: \\ & \textbf{enreg}_i = \textbf{vrai}; \\ & \textbf{etatLocal}_i = \textbf{etatLocal()}; \\ & \textbf{POUR} \ \textbf{j} = 1 \ .. \ \textbf{N} \ \textbf{SAUF} \ \textbf{i} \\ & \textbf{envoyer(} \ \textbf{MARQ,} \ \textbf{P}_i \textbf{)} \\ & \textbf{etatCanal[i]}_i = \varnothing; \\ & \textbf{marqRecu[i]}_i = \textbf{faux;} \\ & \textbf{FPOUR} \end{split}
```

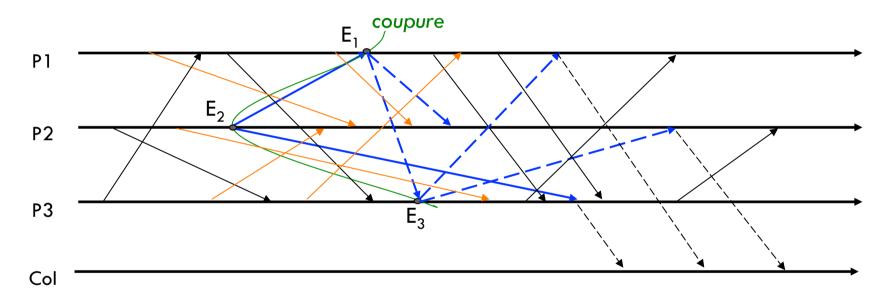
```
recevoir( m, P<sub>i</sub>):
SIm == MARQ ALORS
   SI! enreg; ALORS
      sauvegarder()
  FSI
  marqRecu[i]_i = vrai;
  SI toutRecu() ALORS
      envoyer( <etatLocal, etatCanal[], Collecteur );</pre>
     enreg_i = faux;
     marqRecu[]_i = \{faux, faux, ...\};
  FSI
   RETOURNER
FSI
SI enreg; ET! marqRecu[j]; ALORS
  etatCanal[i]: \cup = \{ m \};
FSI
```

III.3 Algorithme (2)

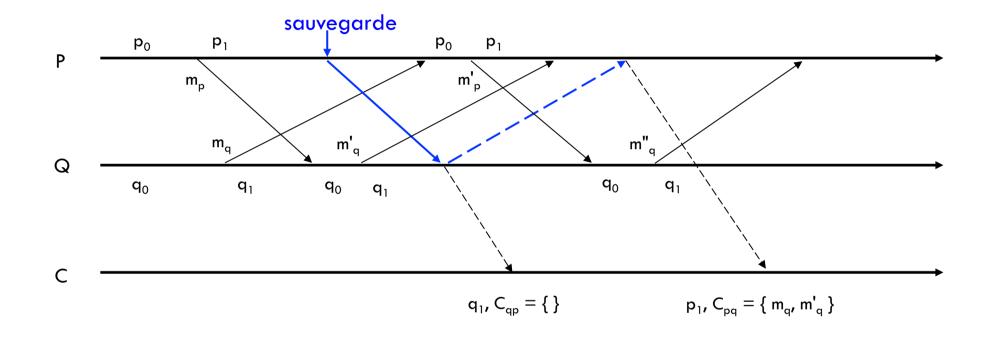
etatLocal(): RETOURNER variables locales du processus

```
toutRecu():
POUR j = 1 .. N SAUF i
   SI ! marqRecu[j]; ALORS
     RETOURNER faux;
FSI
FPOUR
RETOURNER vrai;
```

III.4 Exemples



III.4 Exemples



Etat global collecté :
$$< p_1$$
, $\{ \}$, $\{ m_q, m_q' \}$, $q_1 >$

III.5 Conclusion

Propriétés

L'algorithme se termine.

L'état global enregistré est complet.

L'état global enregistré est cohérent.

Remarque

L'état enregistré peut ne pas correspondre à un état global effectivement atteint par le système au cours de la même exécution. Mais on peut montrer qu'il correspond à un état global que le système aurait atteint en réordonnant les événements de façon compatible avec la causalité (i.e. en réordonnant les seuls événements concurrents)

III.5 Conclusion

Extension à des canaux non FIFO [Lai-Yang 87]

