

Gestion du temps & état global dans un système distribué

Babahenini Med Chaouki

WEBLIOGRAPHIE

- Durant la préparation de ce cours j'ai eu recours aux supports disponibles sur le WEB et notamment au cours des personnes suivantes:
- **Gérard FLORIN**: *professeur en informatique au CNAM de Paris.*
- **Michel RIVEILL** : *professeur d'informatique à l'Ecole Polytechnique de l'Université de Nice*
- **Yahya SLIMANI** : *professeur en informatique à la Faculté des Sciences de Tunis*
- **Sacha KRAKOWIAK**: *professeur émérite en informatique à l'université Joseph Fourier, Grenoble*
- **Eric CARIOU**: *Maître de conférences au département informatique de l'université de Pau*

État global et horloge

- Pour chaque processus du système
 - État local : valeur des variables du processus à un instant t
- État global du système
 - Valeur de toutes les variables de tous les processus du système à un instant t
 - Union à l'instant t des états locaux de tous les processus
- Problème
 - Un état est lié à un instant t mais
 - Chaque processus à une horloge physique locale
 - Pas d'horloge globale dans un système distribué
- La définition d'un état global est possible seulement si on est capable de définir un temps global

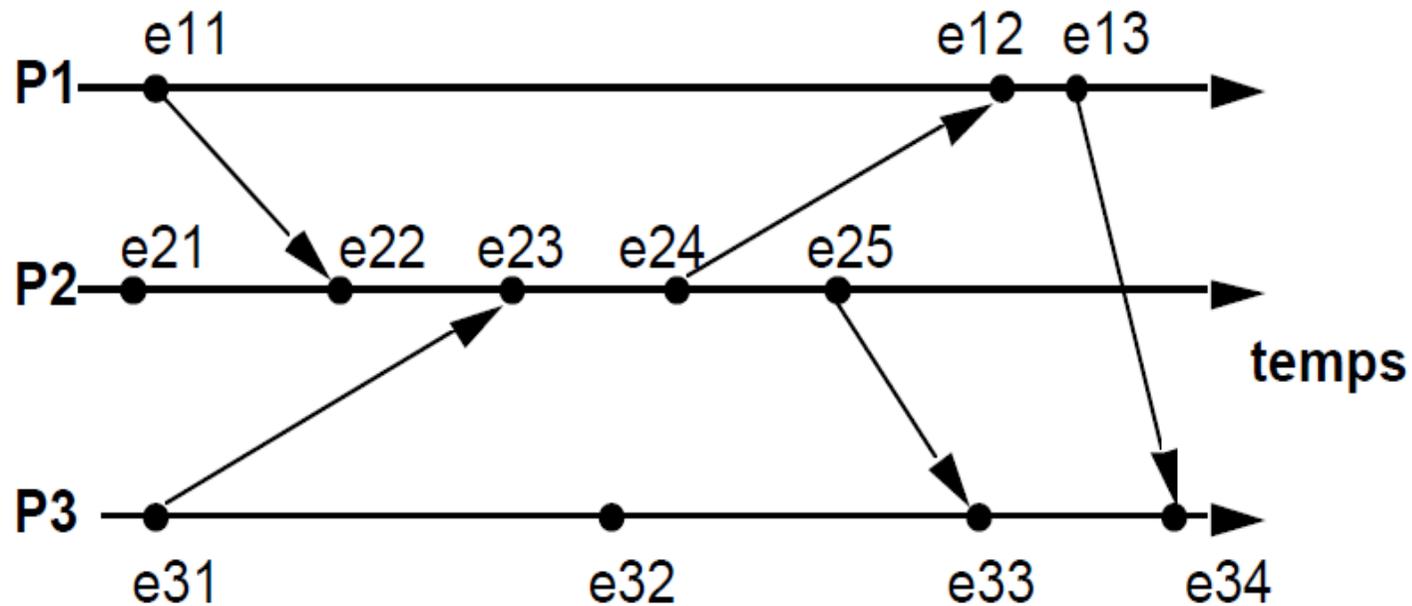
Temps

- Définir un temps global cohérent et « identique » (ou presque) pour tous les processus
 - Soit synchroniser au mieux les horloges physiques locales avec une horloge de référence ou entre elles
 - Soit créer un temps logique
- Synchronisation des horloges physiques locales
 - But est d'éviter qu'une horloge locale dérive trop par rapport à un référentiel de temps
 - La dérive est bornée en augmentation et en diminution
- Deux modes
 - Synchronisation interne
 - Synchronisation externe

Etat global et ordre des événements

- Systèmes répartis à évolution asynchrone
 - pas de mémoire commune (communication par messages)
 - pas d'horloge commune
- Propriétés du système de communication
 - temps de communication potentiellement longs par rapport aux temps de transition locaux
 - temps de transmission variables
 - (éventuellement) non-préservation de l'ordre des messages
- Conséquences
 - perception différente des mêmes événements depuis des sites distincts
 - impossibilité de définir simplement un "état global instantané"
 - conséquences sur la synchronisation, l'observation et la mise au point, tolérance aux fautes

Synchronisation dans un système réparti



Synchronisation dans un système réparti

- Système = ensemble de processus (1 par site) + canaux de communication
- Processus = séquence d'événements locaux + mémoire locale + horloge locale
- Événement = changement d'état du processus (événement interne), ou émission d'un message, ou réception d'un message
- Problèmes fondamentaux :
 - définir un état global
 - définir un temps global
 - déterminables localement par tout processus
- temps global = temps virtuel (ordre sur les événements) tout ordre doit être compatible avec la causalité

Temps logique : chronogramme

- Chronogramme: Décrit l'ordonnancement temporel des événements des processus et des échanges de messages
- Chaque processus est représenté par une ligne
- Trois types d'événements signalés sur une ligne
 - Émission d'un message à destination d'un autre processus
 - Réception d'un message venant d'un autre processus (interactions entre les sites)
 - Événement interne dans l'évolution du processus (mémoire commune)
- Les messages échangés doivent respecter la topologie de liaison des processus via les canaux

Temps logique : dépendance causale

Relation de dépendance causale

- Il y a une dépendance causale entre 2 événements si un événement doit avoir lieu avant l'autre
- Notation : $e \rightarrow e'$
 - e doit se dérouler avant e'

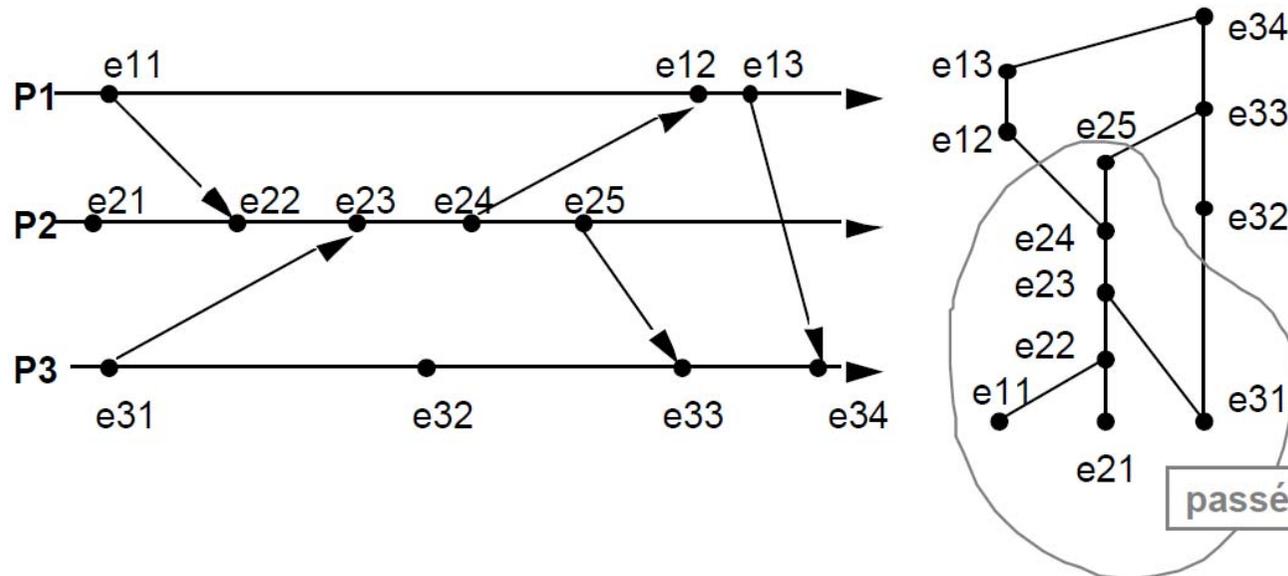
Si $e \rightarrow e'$, alors une des trois conditions suivantes doit être vérifiée pour e et e'

- Si e et e' sont des événements d'un même processus, e précède localement e'
- Si e est l'émission d'un message, e' est la réception de ce message
- Il existe un événement f tel que $e \rightarrow f$ et $f \rightarrow e'$

Relation de parallélisme : $||$

- $e || e' \leftrightarrow \neg ((e \rightarrow e') \vee (e' \rightarrow e))$
- Parallélisme logique : ne signifie pas que les 2 événements se déroulent simultanément mais qu'ils peuvent se dérouler dans n'importe quel ordre

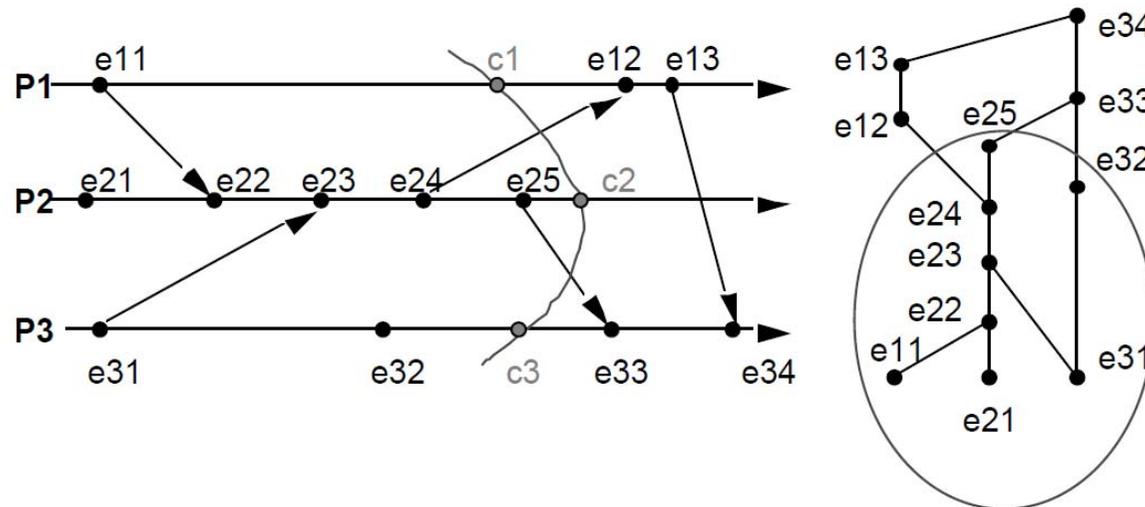
Temps logique : dépendance causale



- Passé (ou historique) d'un événement e (noté $\text{hist}(e)$) : e ensemble des événements e' tels que: $e' \rightarrow e$. Seul le passé strict de e peut influencer e .
- Chaîne causale : e_0, \dots, e_n en tels que $e_{i-1} \rightarrow e_i, i=1, \dots, n$

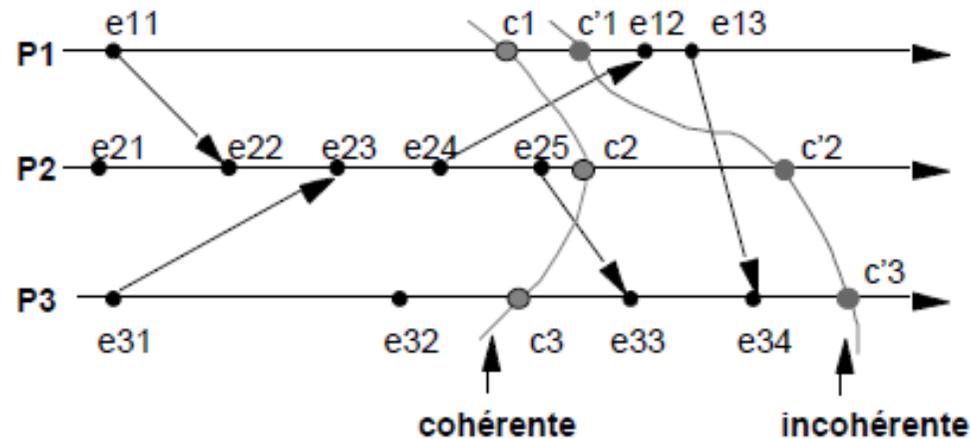
Etat global et coupures

- Calcul distribué = ensemble d'événements E
- Coupure = sous-ensemble fini C de E, tel que
 - a appartient à C et (b précède localement a) \rightarrow (b appartient à C)
- Coupure = "photographie" instantanée de l'état d'un système : ensemble d'événements (virtuels), un par processus ; permet de définir un passé et un futur (par rapport à la coupure)
- Relation d'ordre (partiel) sur les coupures, par inclusion des passés
- Etat associé à une coupure C = ensemble d'états locaux si (un par processus) tel que l'effet local de tout événement de C soit pris en compte dans s_i



Coupures cohérentes

- Cohérence = respect de la causalité (un message ne peut pas venir du futur)



- Coupure cohérente = coupure fermée par la relation de dépendance causale
 - a appartient à C et $(b \rightarrow a)$ alors (b appartient à C)
- Etat global cohérent = état associé à une coupure cohérente
- Exemple de coupure cohérente : le passé d'un événement
- Les coupures cohérentes permettent de définir un ordre de précedence entre états globaux cohérents

Calculs et observations distribués

- Observation d'un calcul réparti E = "linéarisation" de E
(définition d'un ordre total,
soit \ll , sur les événements de E)
- Exemple : vue de la séquence des événements de E par un observateur interne ou externe au système
- Observation valide = compatible avec la relation de précédence causale (pourrait être observée par un observateur externe réel)
- e, e' appartient à E : $e \rightarrow e' \Rightarrow e \ll e'$
- Exemple :
- $e_{11} e_{21} e_{31} e_{22} e_{23} e_{32} e_{24} e_{25} e_{12} e_{33} e_{13} e_{34}$ valide
- $e_{11} e_{21} e_{31} e_{22} e_{32} e_{23} e_{24} e_{25} e_{12} e_{33} e_{13} e_{34}$ valide
- $e_{11} e_{21} e_{31} e_{22} e_{23} e_{32} e_{24} e_{25} e_{12} e_{33} e_{34} e_{13}$ invalide

Temps logique : dépendance causale

- Ordonnancement des événements
- Les dépendances causales définissent des ordres partiels pour des ensembles d'événements
- But d'une horloge logique
 - Créer un ordre total global sur tous les événements de tous les processus
- Horloge logique
 - Fonction $H(e)$: associe une date à chaque événement
 - Respect des dépendances causales
 - $e \rightarrow e' \Rightarrow H(e) < H(e')$
 - $H(e) < H(e') \Rightarrow \neg (e' \rightarrow e)$
 - C'est-à-dire : soit $e \rightarrow e'$, soit $e \parallel e'$

Horloge de Lamport

- Lamport, 1978
- Horloge de Lamport : horloge logique respectant les dépendances causales
- Une date (estampille) est associée à chaque évènement : couple (s, nb)
 - s : numéro du processus
 - nb : numéro d'évènement
- Invariant sur les dates
 - Pour deux dates d'un même processus, les numéros de ces dates sont différentes
- $\forall d, d' : d.s = d'.s \Rightarrow d.nb \neq d'.nb$
- Il n'y pas deux évènements locaux ayant le même numéro

Horloge de Lamport

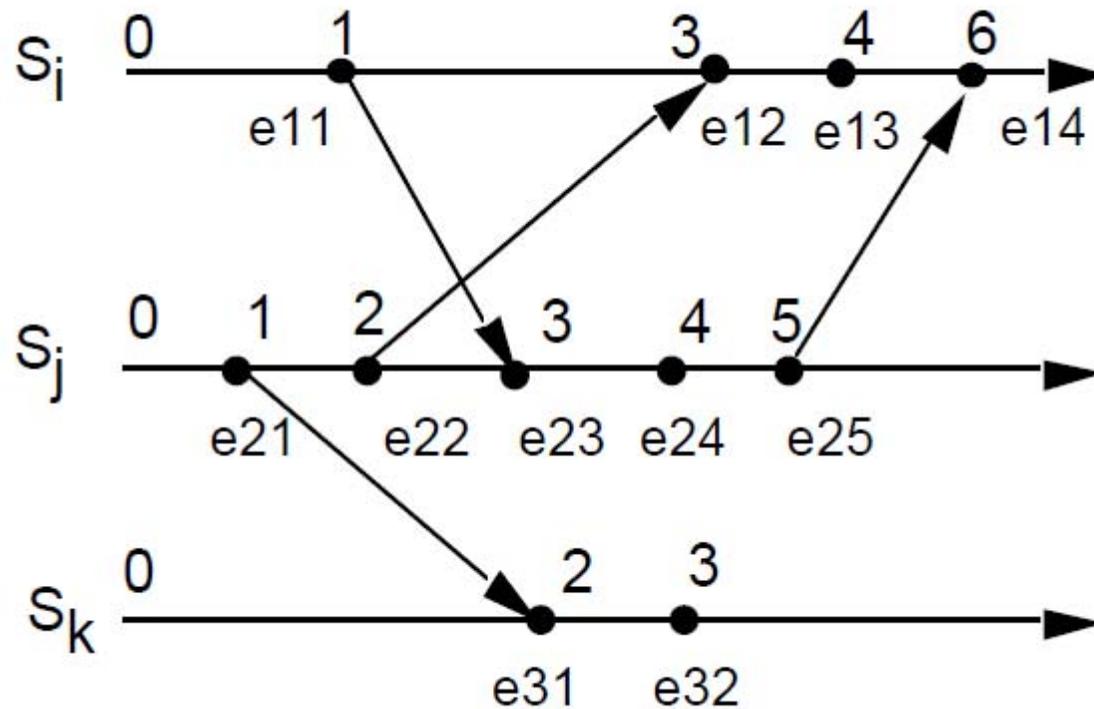
- Création du temps logique
- Localement, chaque processus P_i possède une horloge locale logique H_i , initialisée à 0
- Sert à dater les événements
- Pour chaque événement local de P_i
 - $H_i = H_i + 1$: on incrémente l'horloge locale
 - L'événement est daté localement par H_i
- Émission d'un message par P_i
 - On incrémente H_i de 1 puis on envoie le message avec (i, H_i) comme estampille
- Réception d'un message m avec estampille (s, nb)
 - $H_i = \max(H_i, nb) + 1$ et marque l'événement de réception avec H_i
 - H_i est éventuellement recalée sur l'horloge de l'autre processus avant d'être incrémentée

Horloge de Lamport

- Création du temps logique
- Localement, chaque processus P_i possède une horloge locale logique H_i , initialisée à 0
 - Sert à dater les événements
- Pour chaque événement local de P_i
 - $H_i = H_i + 1$: on incrémente l'horloge locale
 - L'événement est daté localement par H_i
- Émission d'un message par P_i
 - On incrémente H_i de 1 puis on envoie le message avec (i, H_i) comme estampille
- Réception d'un message m avec estampille (s, nb)
 - $H_i = \max(H_i, nb) + 1$ et marque l'événement de réception avec H_i
 - H_i est éventuellement recalée sur l'horloge de l'autre processus avant d'être incrémentée

Horloge de Lamport

- Exemple : chronogramme avec ajouts des estampilles



Horloge de Lamport

- Réalisation d'un ordre total (\ll)
- Soit a sur S_i , b sur S_j
 - $a \ll b \Leftrightarrow (H(a) < H(b)) \text{ ou } (H(a) = H(b) \text{ et } i < j)$
 - L'ordre sur les sites est arbitraire
 - Localement (si $i = j$), H_i donne l'ordre des événements du processus
 - Les 2 horloges de 2 processus différents permettent de déterminer l'ordonnancement des événements des 2 processus
 - Si égalité de la valeur de l'horloge, le numéro du processus est utilisé pour les ordonner
- On dispose donc d'un ordre total sur les événements, détectable à partir d'informations locales uniquement
 - observation : $e_{11} e_{21} e_{22} e_{31} e_{12} e_{23} e_{32} e_{13} e_{24} e_{25} e_{14}$
- Il ne doit pas y avoir de communications "cachées" (via d'autres canaux que ceux qui estampillent les messages)

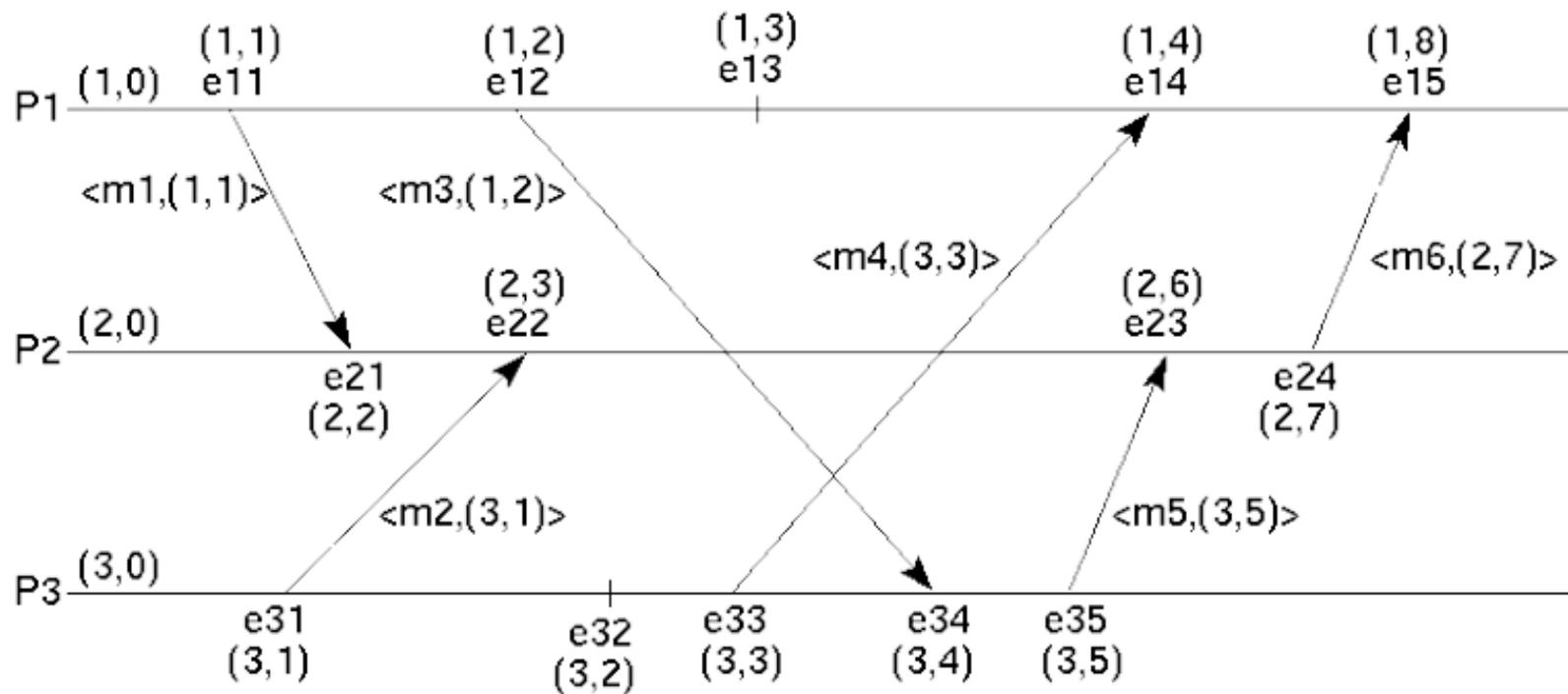
Utilisation des estampilles

- **Algorithmes utilisant une “file d’attente virtuelle” répartie**
 - exclusion mutuelle répartie
 - mise à jour de copies multiples
 - diffusion cohérente (ordre de réception uniforme)
- **Détermination de l’accès “le plus récent”**
 - gestion cohérente de caches multiples
 - fichiers répartis
 - mémoire virtuelle répartie
- **Génération de noms uniques**
 - désignation interne par noms universels
 - mécanismes d’authentification
- **Synchronisation des horloges physiques**
 - borne supérieure sur la dérive entre sites

Avantages

- Trois principaux avantages :
 - Première datation répartie introduite
 - Économiques : la datation est réalisée par un seul nombre et non par un vecteur
 - Causalité des messages est respectée par remise à l'heure du récepteur

Exemple



Limitations de l'ordonnancement par estampilles

- L'ordonnancement global obtenu est arbitraire et n'est pas forcément celui obtenu en pratique. => Les estampilles “effacent” artificiellement la notion de dépendance causale.
 - $H(e_{32}) = 2$ et $H(e_{22}) = 3$ pourtant e_{22} est exécuté en pratique avant e_{32}
- $E1 \rightarrow E2$ implique $H(E1) < H(E2)$ est une condition faible, car on ne peut rien affirmer si $H(E1) < H(E2)$
 - soit les évènements $E1$ et $E2$ sont concurrents
 - soit les évènements $E1$ et $E2$ sont ordonnés
- Ordre total global obtenu pour l'exemple
 - $e_{11} \ll e_{31} \ll e_{12} \ll e_{21} \ll e_{32} \ll e_{13} \ll e_{22} \ll e_{33} \ll e_{14} \ll e_{34} \ll e_{35} \ll e_{23} \ll e_{24} \ll e_{15}$
 - D'autres sont valides ...

Limitations de l'ordonnement par estampilles

- Les estampilles ne sont pas “denses” => Horloge de Lamport respecte les dépendances causales mais avec e et e' tel que $H(e) < H(e')$ on ne peut rien dire sur la dépendance entre e et e'
 - Dépendance causale directe ou transitive entre e et e' ?
 - Exemple : $H(e_{32}) = 2$ et $H(e_{13}) = 3$ mais on a pas $e_{32} \rightarrow e_{13}$

soient $E1$ et $E2$ tels que $H(E1) < H(E2)$, on ne peut pas savoir, s'il existe $E3$ tel que $E1 \rightarrow E3$ et/ou $E3 \rightarrow E2$

Horloges vectorielles (de Mattern)

- Les estampilles ne conservent pas la trace des chemins de causalité : réciproque de la dépendance causale
- Non constantes = pas de lien avec la précédence
- Nécessité de conserver l'état de chaque estampille pour déterminer localement la précédence entre deux événements
- Vecteur d'horloges
- Horloge de Mattern & Fidge, 1989-91
 - Horloge qui assure la réciproque de la dépendance causale
 - $H(e) < H(e') \Rightarrow e \rightarrow e'$
 - Permet également de savoir si 2 événements sont parallèles (non dépendants causalement)
 - Ne définit par contre pas un ordre total global

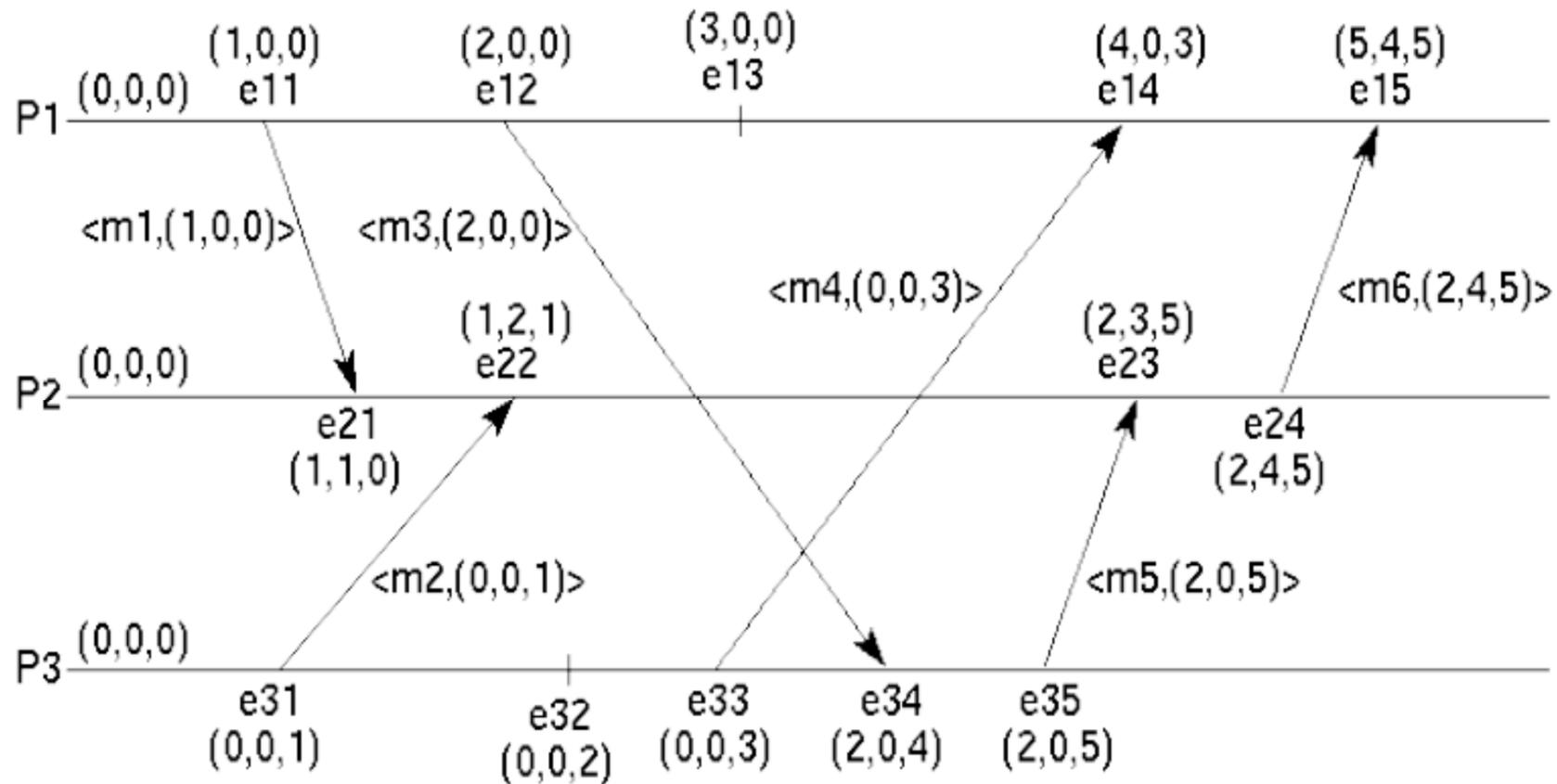
Principe des horloges vectorielles

- Principe
 - Utilisation d'estampilles sur les messages
 - Vecteur V de taille égale au nombre de processus
 - Localement, chaque processus P_i a un vecteur V_i
 - Pour chaque processus P_i , chaque case $V_i[j]$ du vecteur contiendra des valeurs de l'horloge du processus P_j

Horloges vectorielles (de Mattern)

- On associe un vecteur V_i à chaque processus P_i (ou site S_i), $i = 1, \dots, n$
- Initialement, $V_i = (0, \dots, 0)$
- A chaque événement local à P_i ,
 - $V_i[i] := V_i[i] + 1$
- Chaque message m porte une estampille V_m
 - ($V_m = V_i$ de l'émetteur)
- A la réception de (m, V_m) par un processus P_i :
- $V_i[j] := \max(V_i[j], V_m[j])$ pour $j = 1, \dots, n, j \neq i$
 - Mémorise le nombre d'événements sur P_j qui sont sur P_j dépendants causalement par rapport à l'émission du message
 - La réception du message est donc aussi dépendante causalement de ces événements sur P_j
- $V_i[i] := V_i[i] + 1$

Exemple : chronogramme d'application horloge de Mattern



Caractéristiques des HV

- Relation d'ordre partiel sur les dates
 - $V \leq V'$ défini par $\forall i : V[i] \leq V'[i]$
 - $V < V'$ défini par $V \leq V'$ et $\exists j$ tel que $V[j] < V'[j]$
 - $V \parallel V'$ défini par $\neg(V < V') \wedge \neg(V' < V)$
- Dépendance et indépendance causales
 - Horloge de Mattern assure les propriétés suivantes, avec e et e' deux événements et $V(e)$ et $V(e')$ leurs datations
 - $V(e) < V(e') \Rightarrow e \rightarrow e'$
 - Si deux dates sont ordonnées, on a forcément dépendance causale entre les événements datés
 - $V(e) \parallel V(e') \Rightarrow e \parallel e'$
 - Si il n'y a aucun ordre entre les 2 dates, les 2 événements sont indépendants causalement

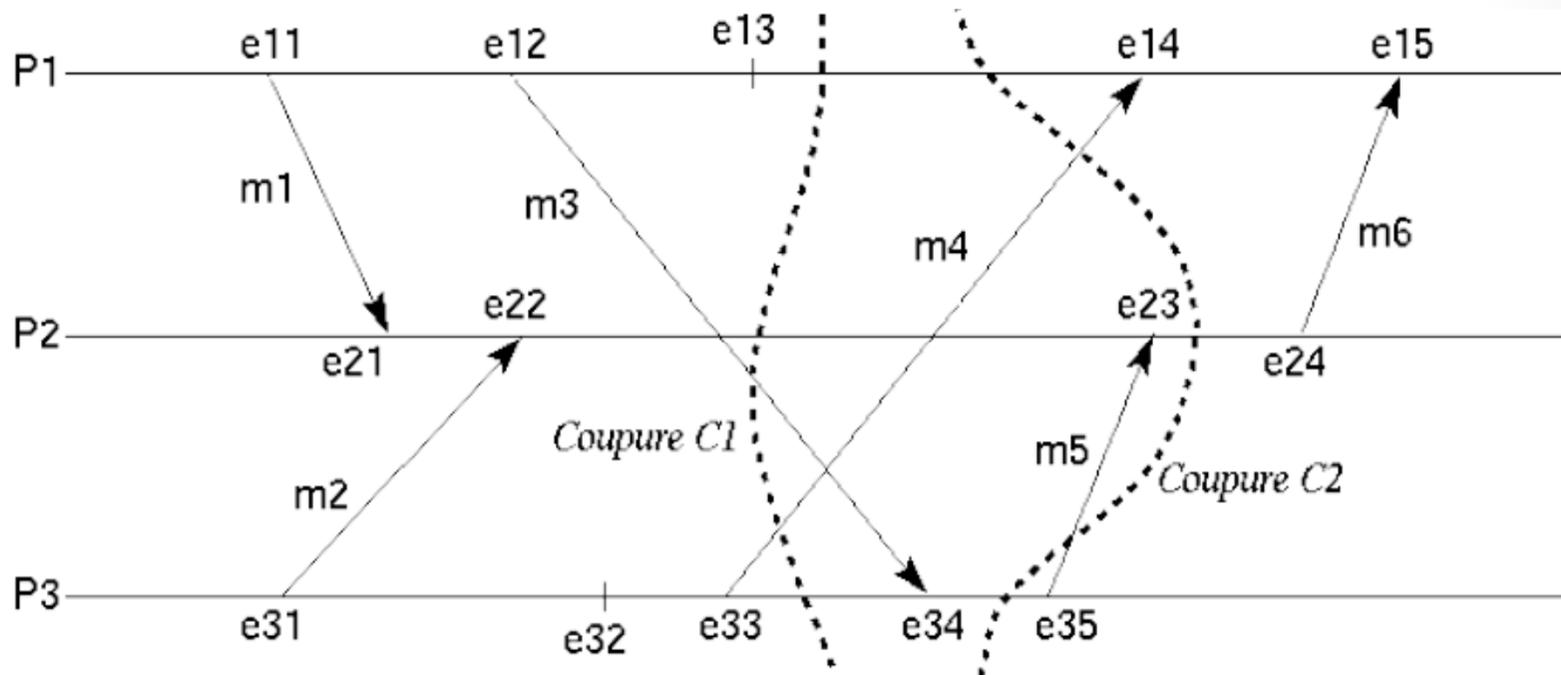
Caractéristiques des HV

- $V(e_{13}) = (3,0,0)$, $V(e_{14}) = (4,0,3)$, $V(e_{15}) = (5,4,5)$
 - $V(e_{13}) < V(e_{14})$ donc $e_{13} \rightarrow e_{14}$
 - $V(e_{14}) < V(e_{15})$ donc $e_{14} \rightarrow e_{15}$
- $V(e_{35}) = (2,0,5)$ et $V(e_{23}) = (2,3,5)$
 - $V(e_{35}) < v(e_{23})$ donc $e_{35} \rightarrow e_{23}$
- L'horloge de Mattern respecte les dépendances causales des événements
 - Horloge de Lamport respecte cela également
- $V(e_{32}) = (0,0,2)$ et $V(e_{13}) = (3, 0, 0)$
 - On a ni $V(e_{32}) < V(e_{13})$ ni $V(e_{13}) < V(e_{32})$ donc $e_{32} \parallel e_{13}$
 - L'horloge de Mattern respecte les indépendances causales
 - L'horloge de Lamport impose un ordre arbitraire entre les événements indépendants causalement.
- Limite de l'horloge de Mattern
- Ne permet pas de définir un ordre global total
- En cas de nombreux processus, la taille du vecteur peut-être importante et donc des données à transmettre relativement importante

HV et Coupures cohérentes

- Horloge de Mattern permet de dater la coupure
 - Soit N processus, C la coupure, e_i l'événement le plus récent pour le processus P_i , $V(e_i)$ la datation de e_i et $V(C)$ la datation de la coupure
 - $V(C) = \max (V(e_1), \dots , V(e_N)) :$
 $\forall i : V(C)[i] = \max (V(e_1)[i] , \dots , V(e_N)[i])$
 - Pour chaque valeur du vecteur, on prend le maximum des valeurs de tous les vecteurs des N événements pour le même indice
- Permet également de déterminer si la coupure est cohérente
 - Cohérent si $V(C) = (V(e_1)[1] , \dots , V(e_i)[i] , \dots , V(e_N) [N])$
 - Pour un processus P_i , si l'événement e_i est le plus récent c'est lui qui a la date la plus récente pour C : sinon un événement e_j d'un processus P_j ($i \neq j$) s'est déroulé après un événement $e_{i'}$ de P_i avec $e_{i'}$ plus récent que e_i
 - $e_i \rightarrow e_{i'}$ et $e_{i'} \rightarrow e_j$ avec $e_i \in C$, $e_j \in C$ et $e_{i'} \notin C$

Coupure cohérente: Exemple



- Coupure C1 : cohérente
- Coupure C2 : non cohérente car $e35 \rightarrow e23$ mais $e35 \notin C2$
 - La réception de $m5$ est dans la coupure mais pas son émission
 - $m5$ vient du futur par rapport à la coupure

Datation des coupures de l'exemple

- Coupure C1 : état = (e13, e22, e33)
 - $V(e13) = (3,0,0)$, $V(e22) = (1,2,1)$, $V(e33) = (0,0,3)$
 - $V(C) = (\max(3,1,0), \max(0,2,0), \max(0,1,3)) = (3,2,3)$
 - Coupure cohérente car $V(C)[1] = V(e13)[1]$, $V(C)[2] = V(e22)[2]$, $V(C)[3] = V(e33)[3]$
- Coupure C2 : état = (e13, e23, e34)
 - $V(e13) = (3,0,0)$, $V(e23) = (2,3,5)$, $V(e34) = (2,0,4)$
 - $V(C) = (\max(3,2,2), \max(0,3,0), \max(0,5,4))$
 - Non cohérent car $V(C)[3] \neq V(e34)[3]$
 - D'après la date de e23, e23 doit se dérouler après 5 événements de P3 or e34 n'est que le quatrième événement de P3
 - Un événement de P3 dont e23 dépend causalement n'est donc pas dans la coupure (il s'agit de e35 se déroulant dans le futur)

Utilisation des horloges vectorielles

- Datation
- Diffusion fiable avec ordonnancement causal
- Simulation répartie
- Calcul d'état global
- Mise au point

Séquencement des messages selon l'ordre causal

- Définition

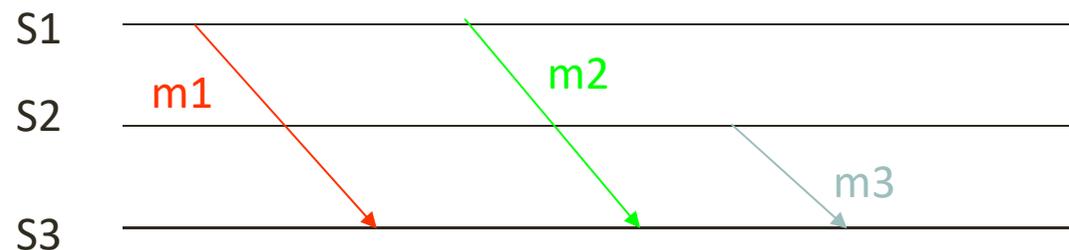
Soient les deux événements $E1 = \text{émission}(m1)$

et $E2 = \text{émission}(m2)$,

L'ordre causal est respecté si la propriété suivante est respectée :

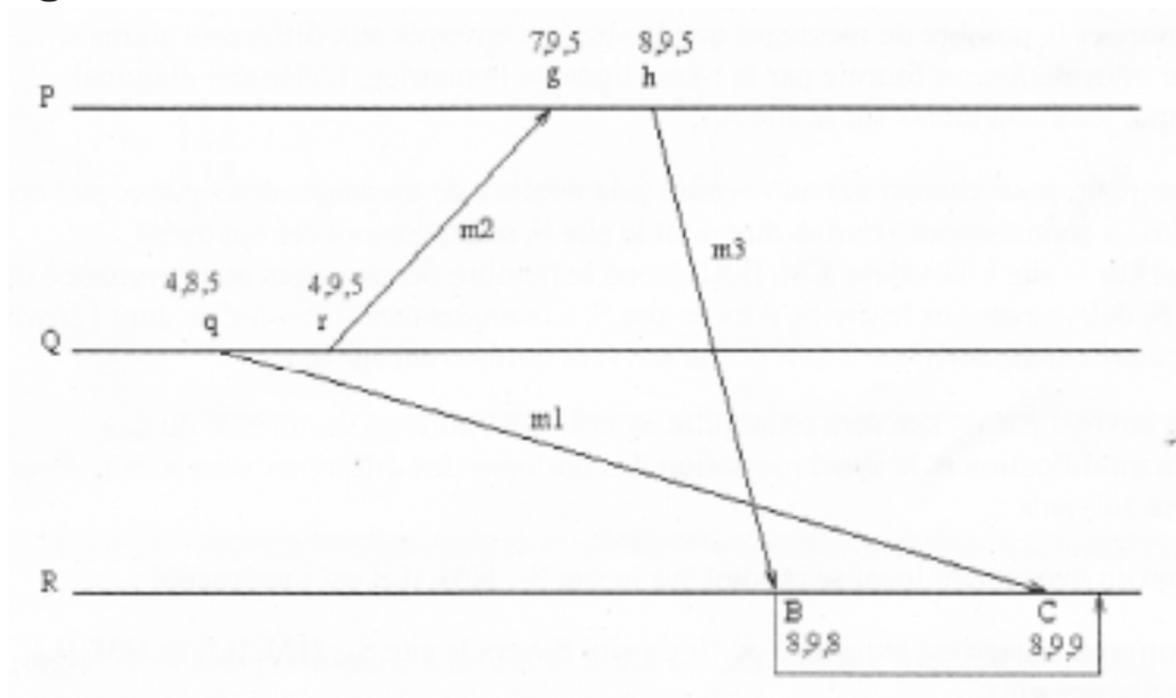
si $(E1 \rightarrow E2)$ et $m1$ et $m2$ ont même destinataire alors
 $\text{réception}(m1) \rightarrow \text{réception}(m2)$

- Exemple



Inconvénient des estampilles vectorielles:

- ne permettent pas de garantir une délivrance causale des messages.



- Dans l'exemple ci-dessus le message m3 est reçu trop tôt: pour respecter l'ordre causal de la délivrance il ne doit être délivré qu'après le message m1.
- Les estampilles permettent de détecter (a posteriori !) la violation de l'ordre causal dans la réception des messages

Horloges et estampilles matricielles

Définition et principe des estampilles matricielles

- Inconvénients des estampilles vectorielles :
 - ⇒ Ne permettent pas de corriger la défaillance vis-à-vis de la délivrance causale des messages
- Principe des estampilles matricielles :
 - Dans un système de n sites, les horloges d'un site i et les estampilles des événements (et des messages) sont des matrices carrées d'ordre n .
 - H_{Mi} désigne l'horloge matricielle du site S_i ,
 - EM_m désigne l'estampille matricielle du message m ,

Fonctionnement des horloges matricielles

Sur le site S_i , la matrice H_{Mi} va :

- mémoriser le nombre de messages que le site S_i a envoyé aux différents autres sites : cette information est fournie par la $i^{\text{ème}}$ ligne de la matrice.
- L'élément diagonal représente la connaissance qu'a S_i du nombre d'événements locaux qui se sont déjà produits sur les différents sites S_j : correspond à l'estampille vectorielle.
- mémoriser, pour chacun des autres sites j , le nombre de messages émis par ce site dont le site i a connaissance (i.e dont le site S_i sait qu'ils ont été envoyés).
- Ainsi sur le site i , la valeur $E_{Mi} [j,k]$ donne le nombre de messages en provenance du site S_j délivrables sur le site S_k dont le site S_i a connaissance (i.e dont l'envoi est causalement antérieur à tout ce qui arrivera dorénavant sur S_i).

La modification synchronisation des horloges des différents sites

- lorsqu'un événement local se produit sur le site S_i : $H_{M_i} [i,i]$ est incrémenté;
- lorsqu'un message est expédié à partir du site S_i vers le site S_j : $H_{M_i} [i,i]$ et $H_{M_i} [i,j]$ sont incrémentés ;
- lorsqu'un message m en provenance du site S_j est reçu sur le site S_i , il faut s'assurer que tous les messages envoyés antérieurement au site S_i y sont effectivement arrivés. Cela suppose donc que S_i ait reçu :
 - d'une part tous les messages précédents de S_j
 - d'autre part tous ceux envoyés plus tôt causalement depuis d'autres sites.
- Cela correspond aux conditions suivantes (à vérifier dans l'ordre):
 1. $H_{M_m} [j,i] = H_{M_i} [j,i] + 1$ (ordre FIFO sur le canal (j,i))
 2. pour tout $k \neq i$ et j , $H_{M_m} [k,i] \leq H_{M_i} [k,i]$ (tous les messages en provenance des sites différents de S_j ont été reçus).

La modification synchronisation des horloges des différents sites

- Si toutes ces conditions sont vérifiées, le message est délivrable et l'horloge du site S_i est mise à jour :
 1. $H_{Mi} [i,i] ++$ (incrémententation);
 2. $H_{Mi} [j,i] ++$ (incrémententation);
 3. pour le reste de la matrice : $H_{Mi} [k,l] = \max (H_{Mi} [k,l], E_{Mm} [k,l])$
- *Si les conditions ne sont pas toutes vérifiées*, la délivrance du message est différée jusqu'à ce qu'elles le deviennent et l'horloge n'est pas mise à jour.
- La délivrance d'un message pourra ainsi provoquer celle des messages arrivés prématurément.

Exemple d'horloge matricielle

