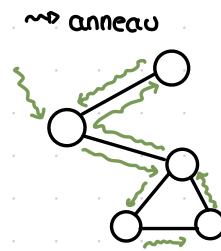


1) Communication

- diffusion
- récolte (gathering)
- construction de structures



2) Contrôle

- Election de leader
- Exclusion mutuelle
- détection de terminaison
- snapshots (instantanés)

3) Temps (séquençement)

- horloges (logiques)

4) Tolérance aux défaillances

- Autostabilisation

- Complexité:
- en messages
 ↓
 nbr de messages émis / bits émis
 - en temps

→ Algorithme diffusion (nouvelle)

type canal = entier
 voisins : ensemble de canal;
 déjà_vu: booléen initialisé à faux

un processus sur réception de Nouvelle:
 si déjà_vu = faux alors
 pour tout c dans voisins (c, Nouvelle);

Diffusion de Nouvelle sur un arbre préalablement construit (fils) avec comme racine le processus recevant Nouvelle:

fils: ensemble de canal;
 racine booléen initialisé à faux sauf pour le processus recevant Nouvelle,

le processus racine sur réception de nouvelle →
 si fils ≠ ∅ alors
 pour tout c dans fils envoyer (c, Nouvelle)

un processus sur réception de Nouvelle →
 si fils ≠ ∅ alors
 pour tout c dans fils envoyer (c, nouvelle)

* missing: cours 3

Cours 4 28/01/26

Diffusion sur un arbre avec détection de terminaison par la racine

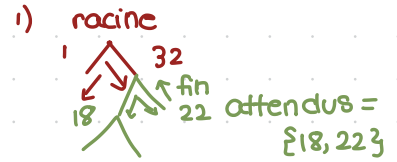
constantes: fils = ensemble de canal;
père = canal

variables: racine: booléen initialisé à faux sauf pour un processus;
attendus: ensemble de canal

la racine sur réception de nouvelle →

si fils = \emptyset alors <terminer> sinon
attendus := fils;
pour tout c dans fils envoyer(c, Nouvelle);

attendus = {1, 32}



fin = acquittement

un processus sur réception de Nouvelle →

si fils = \emptyset alors envoyer (père fin)
sinon attendus := fils; pour tout c dans attendu envoyer(c, fin(Nouvelle));

un processus sur réception de fin (Nouvelle) sur le canal c →

attendus := attendus - {c};
si attendus = \emptyset alors
si racine alors <terminer>
sinon envoyer (père, fin(Nouvelle))

nouvelle : l par arrête
acquittement : l paramètre)ⁿ ⇒ 2n

ne peut pas être exécutée avec moins de n message (sinon ∃ un processus qui n'a pas reçu de message)

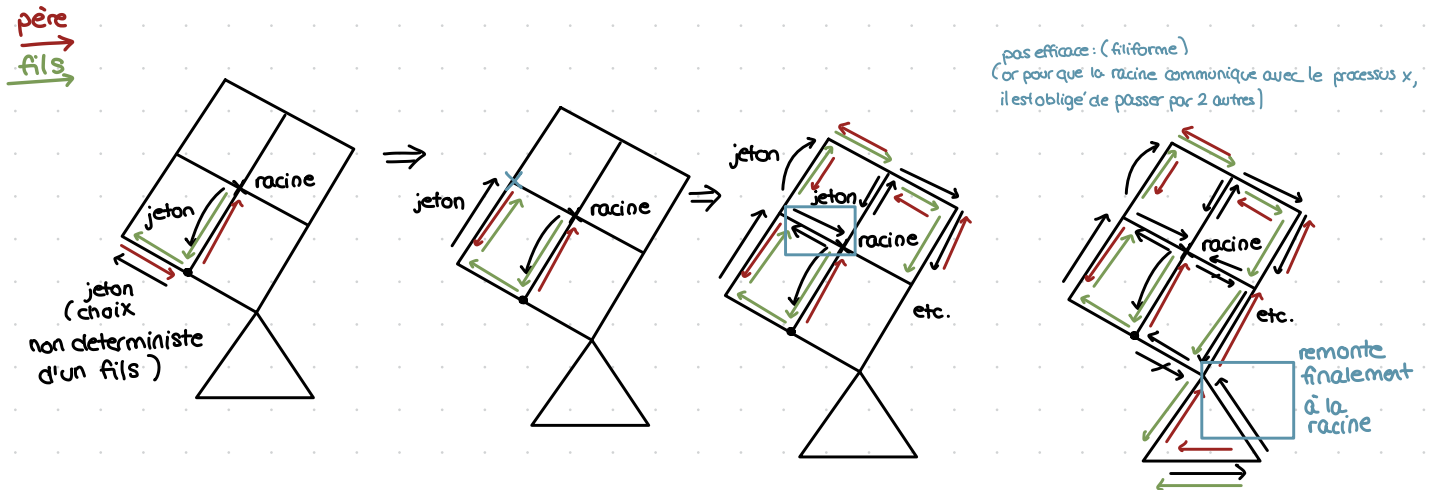
II. Comment construire un arbre recouvrant d'un graphe connexe quelconque?

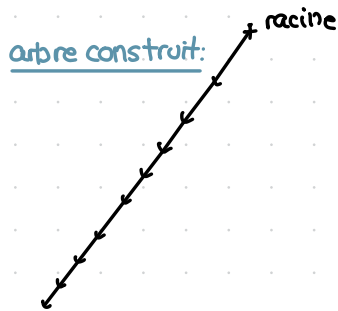
1) Construction par exploration séquentielle: (Audio 1)

Le jeton visite successivement tous les processus.

"acquittement pour signaler qu'il ne veut pas être fils de la racine"

↳ séquentiellement, il est plus probable de tomber sur un processus sans père)

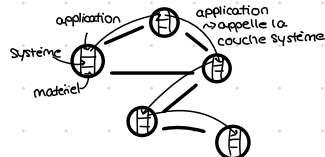




code associé:

constantes et variables

voisins = ensemble de canal;
 initiateur, participant : booléen initialisé à faux;
 encore, fils : ensemble de canal;
 père : canal;



un processus sur réception de Unit →

```

initiateur := vrai;
participant := vrai;
si voisins = ∅ alors terminer
sinon pour tout c dans voisins envoyer (c, Vu); □
  choisir c dans voisins
  envoyer (c, exploration);
  fils := {c};
  encore := encore - {fils};
  
```

un processus sur réception de Vu sur réception de Vu sur le canal c →

encore = encore - {fils};

un processus sur réception de exploration sur le canal c →

```

| si non participant alors | participant = vrai;
|                          | père = c;
|                          | pour tout c dans voisins envoyer (c, Vu);
|                          | encore := encore - {père};
|                          | si encore ≠ ∅ alors
|                          |   choisir c' dans encore;
|                          |   envoyer (c', exploration);
|                          |   fils := {c'};
|                          |   encore := encore - fils;
|                          |   sinon envoyer (père, retour);
| sinon envoyer (c, annule);
  
```

un processus sur réception de retour ou d'annule sur le canal c →

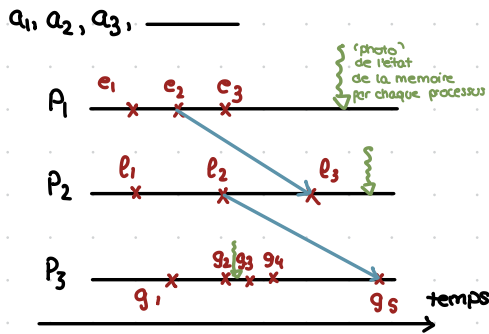
```

si annule alors fils := fils - {c}
si encore ≠ ∅ alors
  choisir c dans encore
  envoyer (c, exploration);
  fils := fils ∪ {c};
sinon si initiateur alors <terminer>
  sinon envoyer (père, retour);
  
```

Horloges logiques

horloge: application d'un ensemble d'événements dans un ensemble (partiellement) ordonné, en relation avec une relation de causalité

processus séquentiel:



événements: entiers
vecteurs d'entiers

$$a \rightarrow \theta(a)$$

Une horloge est compatible avec une relation de causalité si $a \rightarrow b \Rightarrow \theta(a) < \theta(b)$

L'horloge est caractéristique de la causalité si $a \rightarrow b \Leftrightarrow \theta(a) < \theta(b)$

Deux événements a et b sont indépendants ssi on n'a ni $a \rightarrow b$, ni $b \rightarrow a$

Horloge de Edge-Mattem

Histoire de a, $H(a)$ (K premiers événements d'un processus)

$$H(a) = \{ b \mid b \rightarrow a \} \cup \{ a \}$$

↑
passé causal de a

Projections $H_1(a), H_2(a), \dots$

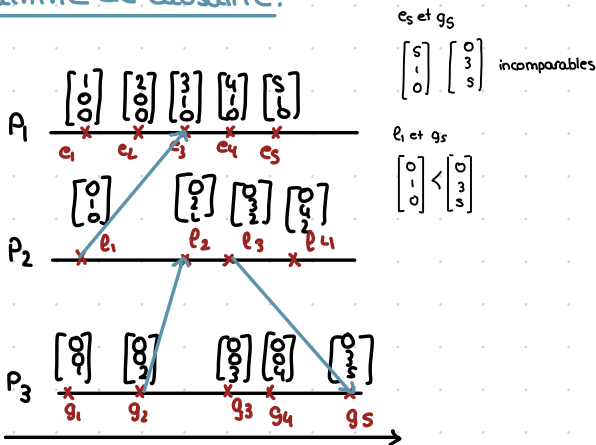
$$\theta(a) = \begin{bmatrix} \text{Card } H_1(a) \\ \text{Card } H_2(a) \\ \vdots \end{bmatrix}$$

(* pour comparer 2 vecteurs :

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} \leq \begin{bmatrix} t \\ u \\ v \end{bmatrix} \Leftrightarrow \begin{cases} x \leq t \\ y \leq u \\ z \leq v \end{cases}$$

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} < \begin{bmatrix} t \\ u \\ v \end{bmatrix} \Leftrightarrow \begin{cases} x \leq t \\ y \leq u \\ z \leq v \end{cases} \text{ et } \neq$$

diagramme de causalité:



preuve: \Rightarrow transitivité

$$a \rightarrow b \Leftrightarrow H(a) \subseteq H(b) \text{ et } a \neq b$$

$$\Leftrightarrow H_1(a) \subseteq H_1(b)$$

$$\Leftrightarrow H_2(a) \subseteq H_2(b) \text{ et } a \neq b$$

$$\Leftrightarrow \begin{matrix} \text{Card}(H_1(a)) \leq \text{Card}(H_1(b)) \\ \text{Card}(H_2(a)) \leq \text{Card}(H_2(b)) \end{matrix} \text{ et } a \neq b$$

$$\Leftrightarrow \begin{bmatrix} H_1(a) \\ H_2(a) \\ \vdots \end{bmatrix} \not\leq \begin{bmatrix} H_1(b) \\ H_2(b) \\ \vdots \end{bmatrix}$$

Calcul de H_i par le processus P_i

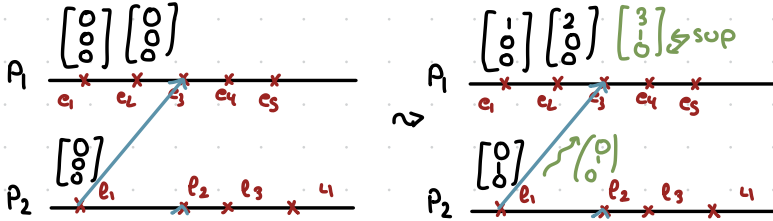
algo

Initialement $H = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

(réception locale)

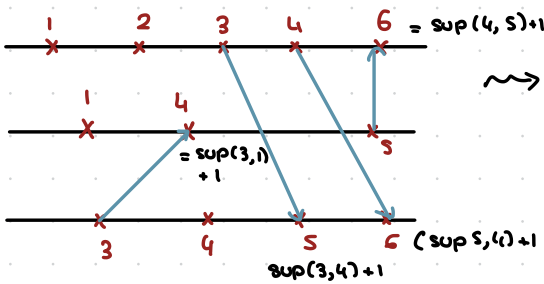
- 1) Lors d'un événement interne, incrémenter la $i^{\text{ème}}$ composante de H
- 2) Lors d'une émission d'un message m , incrémenter la $i^{\text{ème}}$ composante de b_1 et attacher H à m (piggybacking)
- 3) Lors de la réception d'un message m , incrémenter la $i^{\text{ème}}$ composante de H et prendre le sup des autres composantes.

exemple:



Horloge de Lamport:

(ajouter 1 à chaque chemin)



valeur de l'horloge de Lamport de l'événement a

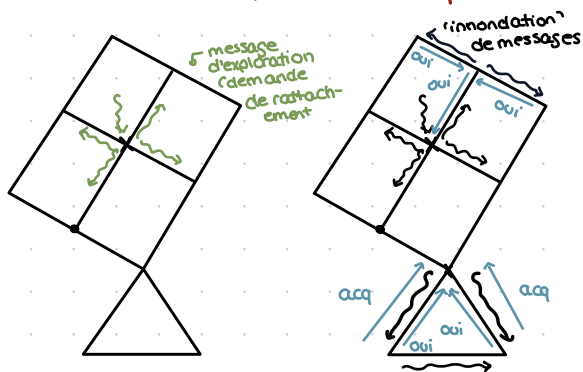
si $a \rightarrow b$ alors $L(a) \leq L(b)$

si on considère $(L(a), i)$
 $(L(a), i) \leq (L(b), j) \Leftrightarrow L(a) < L(b)$
 ou $L(a) = L(b)$ et $i < j$

permet d'avoir un ordre total sur les événements

Suite cours 4 (construction d'arbre)

Construction d'arbre par exploration parallèle



code:

variables / constantes:

- voisins : ensemble de canal
- parti a part : booléen initialisé à faux
- attendus, fils : ensemble de canal
- père : canal
- racine : booléen initialisé à faux

un processus sur réception de Init \rightarrow

participant := vrai; racine := vrai; fils = \emptyset ;

attendus := voisins;

on supprime graphes non-dégénérés

pour tout c dans voisins envoyer (c , explorer);

un processus sur réception de Explorer sur le canal $c \rightarrow$

si non participant alors participant := vrai;

père := c ;

fils := \emptyset ;

attendus := voisins - $\{c\}$;

si attendus $\neq \emptyset$ alors pour-tout c' dans attendus
envoyer (c' , explorer)

sinon envoyer (père, oui)

sinon envoyer (c , Non);

un processus sur réception de oui sur le canal $c \rightarrow$

fils := fils $\cup \{c\}$;

attendus := attendus - $\{c\}$;

si attendus = \emptyset alors

si racine alors <terminer>

sinon envoyer (père, oui);

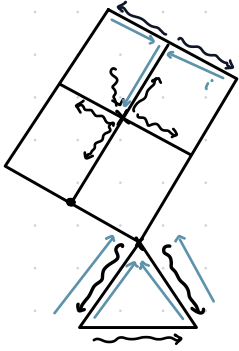
un processus sur réception de Non sur le canal $c \rightarrow$

attendus := attendus - $\{c\}$;

si attendus = \emptyset alors

si racine alors <terminer>

sinon envoyer (père, oui);



attendus: acquittements de vagues
encore: acquittements de terminaux

Bellman-Ford

Algorithme de construction d'arbre recouvrant par vagues synchronisées et détection de terminaison

- variables / constantes:
- voisins;
- encore, attendus, fils : ensemble de canal;
- père : canal;
- niveau : entier;
- racine : booléen initialisé à faux;

Initialement, l'environnement contacte un processus.

Un processus sur communication avec l'environnement →

participant := vrai; racine := vrai;
niveau := 0; encore := voisins; ^{pour avoir la terminaison}
si encore ≠ ∅, alors pour tout c dans voisins envoyer(c, aller(c)) ^{vague à distance de la racine}
attendus := voisins;

un processus sur réception de aller(k) sur le canal c →

si non participant alors participant = vrai; père := c; niveau := k; fils = ∅;
encore := voisins - {père};
si encore = ∅ alors envoyer(père, retour(fin(k)))
sinon envoyer(père, retour(again(k)));

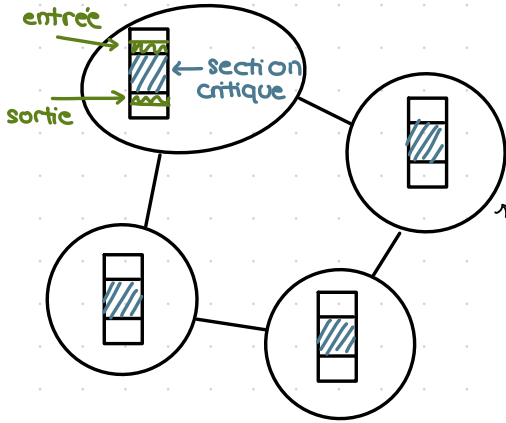
sinon // si participant

si père = c alors père := voisins - {père};
pour tout c dans encore
envoyer(c, aller(k));
attendus := encore
sinon envoyer(c, retour(non(k)));

un processus sur réception (retour(reponse(k)) sur le canal c →

attendus := attendus - {c};
si (reponse = non) ou (reponse = fin) alors encore := encore - {c};
si (reponse = again) ou (reponse = fin) alors fils := fils ∪ {c};
si encore ≠ ∅ et attendus = ∅ alors
si (¬ racine) alors envoyer(père, retour(again(k)))
sinon pour tout c dans encore envoyer(aller(k+n));
attendus := encore;
si encore = ∅ alors si (non racine) alors envoyer(père, retour(fin(k)));
sinon < terminée >

Exclusion mutuelle des sections critiques.



un processus est dans sa section critique si son compteur de programme pointe vers une instruction de la section critique.

↳ on veut que l'accès à la section critique soit séquentiel

↳ à chaque instant, au plus un processus exécute une instruction dans sa section critique (→)

↳ causalement indépendants

Spécification précise d'un problème:

≡ propriété d'exécution

≡ conditions sur les exécutions

Safety (sûreté) = satisfaisabilité à chaque partie d'exécution (→ ne pas construire un arbre en cycle par exemple)

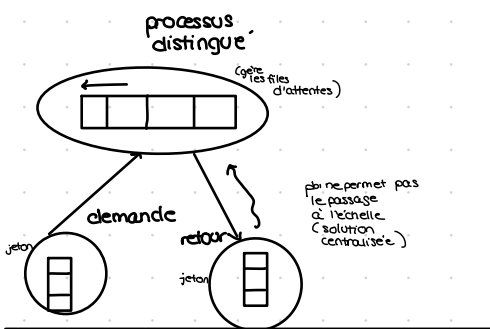
Liveness (vivacité) = satisfaite dans le futur

↳ un processus peut rentrer dans sa section critique au bout d'un temps fini

↳ terminaison par exemple

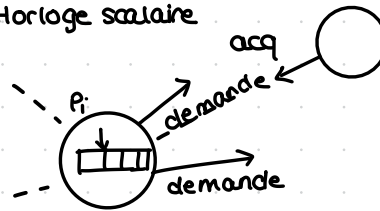
(ex (pour election de leader)
s: à aucun moment on a 2 ou plus leaders
t: terminaison

Solution centralisée:



Horloge de Lamport avec l'ordre total:

Horloge scalaire



(n, i)
↓
num du processus
valeur de l'horloge de Lamport

canaux FIFO (on conserve l'ordre des messages)

Avant d'entrer en section P_i , un processus attend que l'estampille soit la plus petite de toutes les estampilles des messages qu'il a reçu

impliquant que la demande est causalement la plus ancienne et unique

type_message: demande, sortie, acquittement

message: 3 champs: type

{ horloge
numéro de processus

file_d_attente: tableau $[1, \dots, n]$ de message; (initialisé à $file[i] = (sortie, 0, i)$;

un processus P_i pour demander à rentrer en s.c. →

gérer l'horloge;

diffuser (demande, horloge, i);

attendre $file[i]$ contienne le plus petit compte (estampille, numéro) de toute la file d'attente;

dedans := ura_i ;

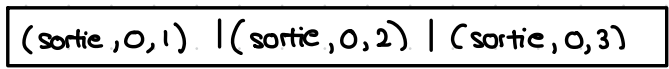
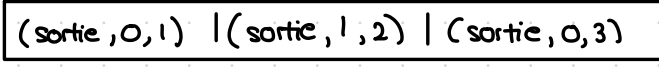
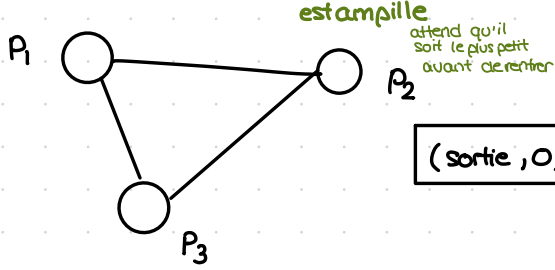
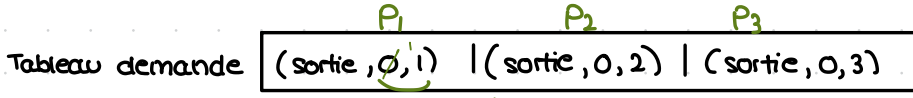
< section critique >

un processus qui reçoit un message de demande de p_i →
tant que de dans attendre
gérer l'horloge;
envoyer(acquittement, horloge, i);

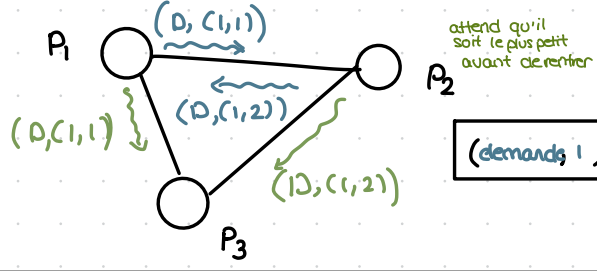
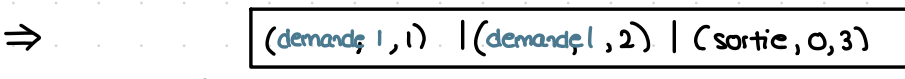
un processus p_i qui sort de section critique →
gérer l'horloge, ^(incrémenter)
diffuser(sortie, horloge, i);

$\forall j$,
à la réception de { acquittement, sortie, j } →
si la cellule j ne contient pas une demande, mettre à jour $file(j)$;

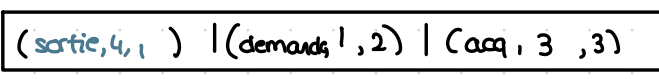
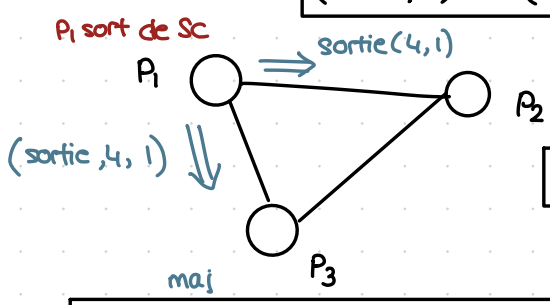
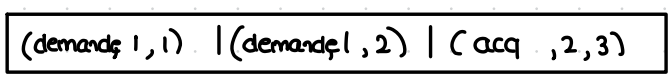
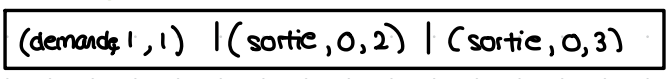
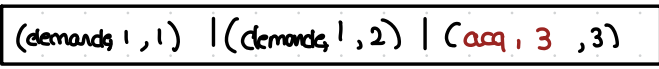
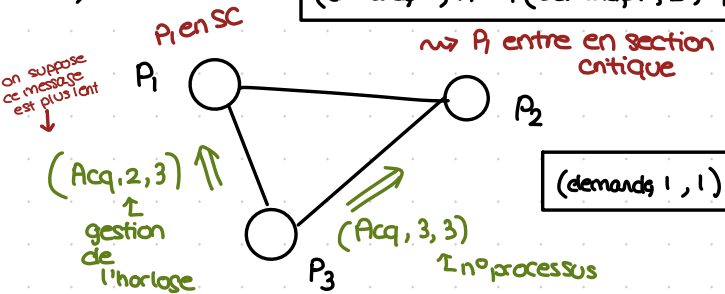
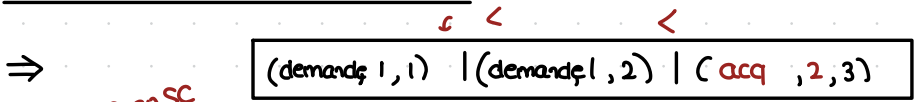
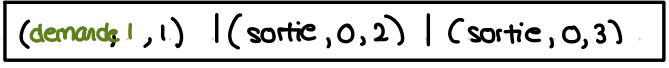
Algorithme d'exclusion mutuelle de Lamport (exemple de déroulement de l'algo)



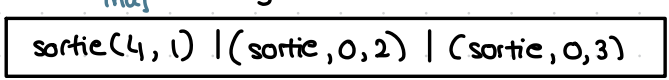
P_1 demande à entrer en S.C.
 P_2 demande à entrer en S.C.
 Durant toute l'exécution, P_3 ne demande pas à entrer en S.C.



! canaux FIFO
 (l'ordre des demandes est conserve')



⇒ P_2 entre en SC



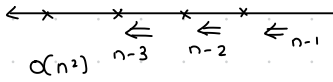
Complexité:

en messages

nb de messages émis pour 1 entrée en S.C.:

$$\left. \begin{array}{l} \rightarrow (n-1) \text{ (P}_1 \text{ diffuse son message)} \\ + (n-1) \text{ (diffuse son message de sortie)} \\ + (n-1) \text{ (acquittements)} \end{array} \right\} \Rightarrow 3(n-1)$$

(nb d'arêtes)
(n-1 si graphe complet)



Algorithme de Ricard et Agrawala:

Privilège: avoir la plus petite estampille parmi les demandes

↳ exclusion mutuelle est réalisée \Rightarrow seul un processus avec une plus petite peut entrer en S.C. (algo + FIFO)
la plus petite demande est unique (Lamport'élargi)

Fairness: Le processus avec la plus petite estampille finit par entrer en S.C.

(pas de perte de message) (pas de blocage)

- processus "précédé" par les autres processus

i) $K < n-1$ (algorithme)

ii) au bout d'un certain temps, $K = K-1$

Privilège implicite (condition)

R.A. privilège explicite \Rightarrow jeton

Le 'jeton' est un tableau qui contient le nb d'entrées en S.C. de chaque processus

jeton (horloge interne de chaque processus)

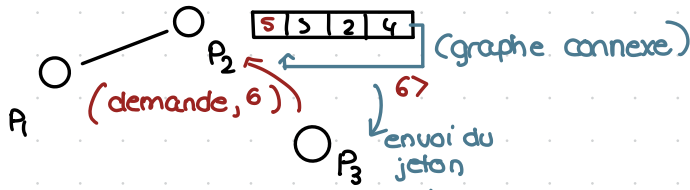
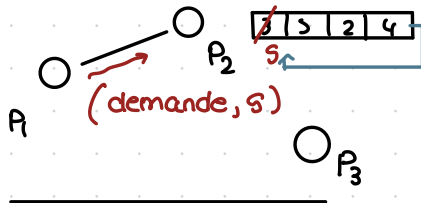
2	0	5	8
---	---	---	---

- Lorsqu'un processus qui "veut" entrer en S.C. libère le jeton à la sortie de S.C., il incrémente son compteur (dans le jeton)

- Au départ, le jeton est placé dans un processus particulier, et toutes ses cellules contiennent 0.

Si la possession du jeton garantit l'entrée en S.C. et si une entrée en S.C. est impossible sans possession du jeton, l'unicité du jeton garantit l'exclusion mutuelle des S.C.

- Un processus qui veut entrer en S.C. diffuse un message de demande, avec le compteur de ses S.C. passés.



algo:

message = (entier, émetteur);

jeton = tableau [site] de entier;

procédure attendre_jeton (v: jeton)

{ attend l'arrivée d'un message de type jeton et range sa valeur dans la variable v }

tampon

choisir le prochain dans l'anneau

constantes, variables

{ pour p_i }

const: numéro initialisé à i ;
var: tampon; demande : jeton;
estampille: horloge initialisé à 0;
jeton - présent: booléen initialisé à (mon - numéro = 1);
dedans: booléen initialisé à faux;
requête : message;

lors de demande d'entrée en SC →

```
si jeton - présent = faux alors
  debut
  estampille := estampille + 1;
  requête . date := estampille;
  requête . émetteur := mon - numéro;
  diffuser (requête);
  attendre - jeton (tampon);
  fin
  jeton présent := vrai;
  dedans := vrai;
```

lors de sortie en S.C →

```
var j : entre site;
tampon (mon - numéro) := estampille;
dedans := faux;
pour j := mon - numéro + 1 jusqu'à n, 1 jusqu'à mon - numéro - 1
  faire:
    si demande (j) > tampon (j) et jeton - présent
      alors | jeton - présent := faux;
             | envoyer (tampon, j);
```

Lors de l'arrivée d'une requête →

```
var h : numéro;
K := requête - émetteur;
demande [K] := max (demande [h], requête . date);
* au cas où il reçoit une demande "primée"
si jeton - présent et dedans = faux alors
  {}
```

→ **complexité**: en message: $(n-1) + 1 = n$

exclusion mutuelle: exclusion du jeton dont la possession est une condition nécessaire et suffisante d'entrer en section critique

Fairness: exploration circulaire

Election sur un anneau

but: désigner un processus particulier, parmi des processus numérotés

spécialisation: **Fairness**: l'algorithme se termine explicitement

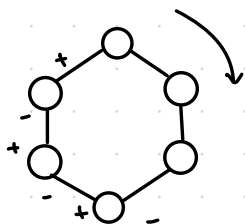
Safety: jamais 2 leaders et 1 seul à la terminaison.

→ graphe en forme d'anneau : - **unidirectionnel** (successeur)
- **bidirectionnel** (successeur + prédécesseur)
- **orienté** ssi étant donné 2 voisins p_i et p_j ,

Algorithme de Chang et Roberts mon - successeur (p_i) = $p_j \iff$ mon - prédécesseur (p_i) = p_j

Élection de leader

Réseau : anneau



Protocole:

- Terminaison (implicite ou explicite)
- Jamais plusieurs leaders et un seul à la terminaison.

Algorithme de Chang et Roberts (extinction sélective)

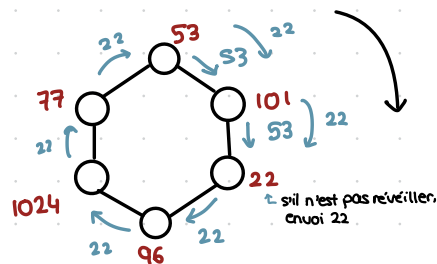
L'élu est le processus avec le plus petit identificateur (entiers)

cas asynchrone. anneau orienté

Lors d'un réveil un processus va "tenter" de faire parcourir

à son numéro le tour de l'anneau

on fait en sorte que le plus petit numéro peut faire le tour



Pseudo-code :

constante : mon_numéro

variable : participant initialisé à faux

leader : booléen initialisé à faux

Lors de réveil (par la couche supérieure) →

si non participant alors

envoyer (+, mon_numéro);

participant := vrai;

Lors de réception de (+, val) sur le lien →

si non participant alors

participant := vrai

si val < mon_numéro alors envoyer(+, val);

sinon envoyer (+, mon_numéro)

sinon si val < mon_numéro alors

envoyer(+, val);

si val = mon_numéro alors

leader := vrai;

< stop >

// sinon on ne fait rien

Preuves:

1) Terminaison:

- le processus avec le plus petit numéro se réveille (à un certain moment)

On considère la plus petite distance avec un processus qui se réveille et le processus avec le plus petit numéro

↳ induction sur la distance

- le plus petit numéro fait (au bout d'un certain temps) le tour de l'anneau

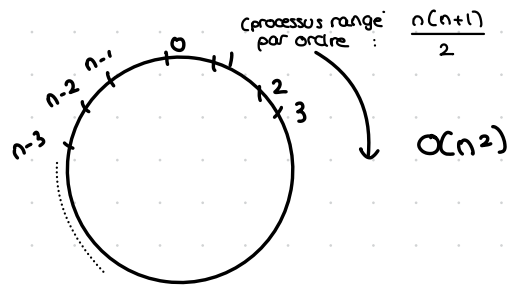
(induction sur la distance parcourue par le plus petit numéro)

2) Complexité

en messages en fonction du nombre de processus : n^2

borne supérieure : $n \times n$

une exécution particulière qui atteint cette borne:



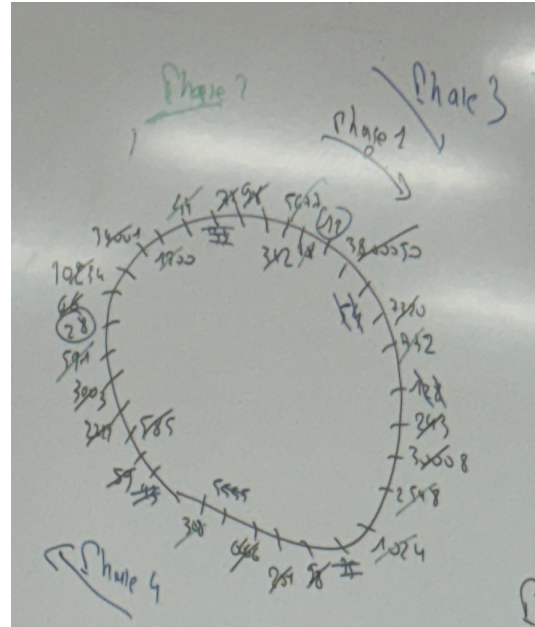
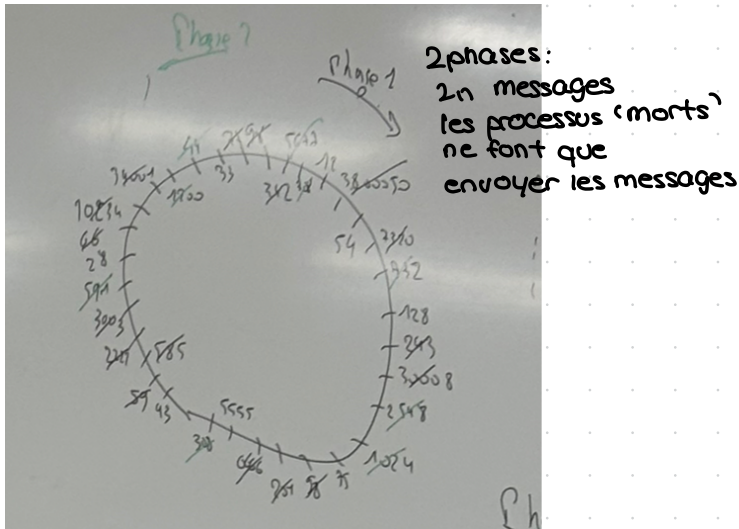
Algorithme de Peterson

anneau bidirectionnel (mais orienté)

phase: exécution où le processus envoie + reçoit (dans un sens ou l'autre)

Chaque phase fait des victimes et laisse des survivants

(TODO)

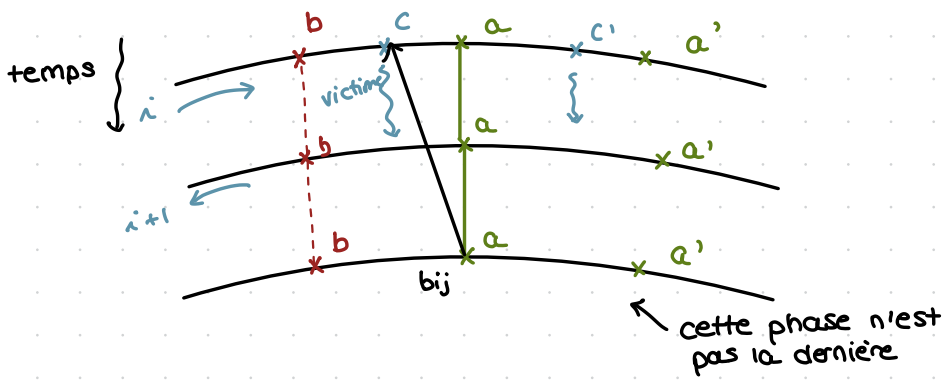


Idee: 1. Une exécution procède par une succession de phases, alternant les directions + et -

2. Une phase est un segment d'exécution dans lequel un processus émet un message et reçoit un message ou reçoit un message et envoie un message. (rq: chaque processus a son propre rythme de phase)
3. Chaque phase fait des victimes et laisse des survivants. Une victime est un processus qui reçoit une plus petite valeur que son numéro par la suite une victime ne fait que relayer les messages
4. Lorsqu'il termine une phase, un processus survivant entame la phase suivante en envoyant son numéro dans la direction opposée à celle de la phase précédente
5. Un processus qui reçoit son numéro devient leader

Éléments de preuve:

1. Le nombre de phases est fini
2. Au bout d'un certain temps, le plus petit numéro fait le tour de l'anneau.
Toutes les deux phases, le nombre de processus vivants au début des deux phases a été divisé par au moins 2 à la fin des 2 phases. (nb de phases $O(\ln(n))$)



a est le numéro d'un processus

Durant la phase $i+1$, a a reçu une valeur plus grande

Durant la phase i , a a reçu une valeur b et puisque a est vivant, $b > a$

b est le plus petit numéro vivant à la fin de la phase $i+1$ dans le sens - par rapport à a .

→ il existe une injection ($c = c'$) entre les processus survivants à la fin d'une phase et les victimes faites dans les 2 phases précédentes

⇒ le nb de processus est divisé par 2 après

Detection de terminaison

Le système utilise des messages de contrôle.

Il doit annoncer que l'application est terminée, si c'est le cas

↳ le processus de l'application

sont passifs et il n'y a pas de travail en transit.

Problème: On a une application distribuée qui échange des messages (messages de travail)

Quand un processus reçoit un message de travail, il devient "actif"

S'il ne reçoit pas d'autres messages de travail, un processus actif devient passif dans le futur.

Sûreté: Le système n'annonce pas de fausse terminaison

Vivacité/Équité: le système doit l'annoncer au bout d'un nombre fini d'opérations.

Hyp: anneau orienté unidirectionnel, communications instantanées

Preuve par invariant

↳ prédicat : - vrai dans la configuration initiale

- quand il est vrai dans une certaine configuration, cela implique que le problème est résolu

Règles:

R₁: Un processus a une couleur, blanc ou noir, initialement blanc

R₂: Le jeton a une couleur, blanc ou noir, initialement blanc.

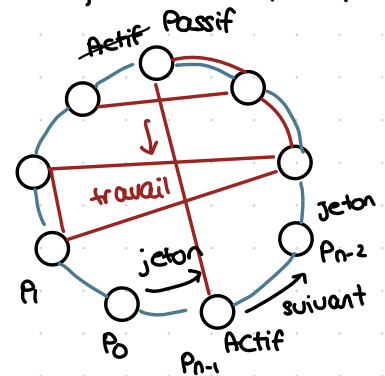
R₃: Un processus qui reçoit un message de travail devient noir

R₄: Seuls les processus passifs transmettent le jeton

R₅: Un processus blanc transmet le jeton

Un processus noir transmet le jeton noir, et redevient blanc

R₆: Si le jeton revient blanc à P₀, lui-même blanc et passif, P₀ <annonce> la terminaison. (...)



Pseudo-code de l'algorithme Dijkstra:

Application

lors d'un envoi d'un message par p de travail reçu par q → état := actif (transmission instantanée)
spontanément → état := passif;

lors de l'envoi d'un message de travail par p reçu par q →

couleur_p := noir;

état_q := actif;

quand état = actif → état := passif

lors de l'initiation d'une détection de terminaison par p₀ →

envoyer (jeton, blanc) à P_{n-1};

quand p reçoit le jeton (jeton, c) et est passif →

si p = p₀ alors

si (c = blanc) et (couleur p₀ = blanc) alors

<annoncer>

sinon envoyer jeton blanc à P_{n-1}

sinon si (couleur_p = blanc) // et est passif alors

envoyer (jeton, c) à suivant sur l'anneau

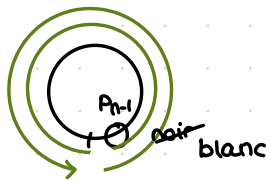
sinon // la couleur de p est noir

envoyer (jeton, noir) à suivant

couleur_p := blanc

Complexité:

nombre maximum de messages de contrôle après que l'application est terminée (latence),



2 tours de l'anneau (-1)

Algorithme de Safia (transmission des messages asynchrones)

Utilisation de Compteurs de message

(règles de Dijkstra+ :)

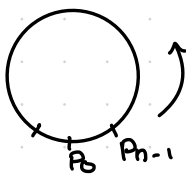
R₁: Chaque processus a un compteur de la différence entre nombre de messages de travail émis et nombre de messages de travail reçu

R₂: le jeton relève les compteurs.

Detection de terminaison (Safia)

{ pour le processus p }

var state = (actif, passif);
 color : (blanc, noir);
 me : entier init à 0;



quand state = actif \rightarrow
 envoyer (message);
 me := me + 1;

quand un message mes arrive \rightarrow

receve (mes);
 me := me - 1;
 color := noir;

quand état = actif \rightarrow
 état := passif;

pour lancer une détection de terminaison \rightarrow

envoyer (tok, blanc, 0) à P_{n-1} ;

pour un processus pour traiter le token $\langle \text{tok}, c, q \rangle$

quand état = passif \rightarrow

si $p = p_0$ alors si $\langle c = \text{blanc} \rangle$ et $\text{color}_p = \text{blanc}$ et $me_p + q = 0$ alors $\langle \text{Annoncer} \rangle$
 sinon envoyer (tok, blanc, 0) à P_{n-1}

sinon si color = blanc alors envoyer (tok, c, $q + me_p$) à suivant
 sinon envoyer (tok, noir, $q + me_q$) à suivant;

color := blanc;

(le jeton fait au plus 2 tours après la terminaison)

Gerard Tel - Distributed algorithm - pu par invariant

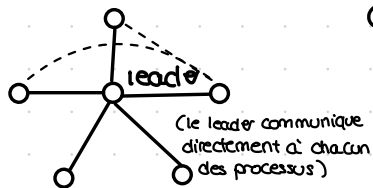
Seconde version des compteurs (Mattern)

(au plus 1-tour après la terminaison)



Algorithme du crédit (Mattern)

Calcul diffusant:



Quand le leader reçoit le signal de l'application, il positionne une variable crédit = 1.

Tous les processus sont initialement passifs

Ils deviennent actif par réception d'un message de travail du leader

Quand le leader envoie un message de travail au processus p , il joint à ce message le réel $c/2$ et conserve un crédit $c/2$

Quand un processus p avec un crédit c_p envoie un message à q , il attache $c_p/2$ au message et conserve $c_p/2$

(invariant: la somme des crédits des processus et les messages ^{+transportés par} = 1)

Quand un processus avec un crédit c_p reçoit un message transportant un crédit c , il fait $c_p := c_p + c$

Quand un processus devient passif, il renvoie son crédit au leader

Le leader a une variable retour, initialisée à 0 et qui est modifiée par les crédits rendus au leader

Quand retour = 1, le leader annonce la terminaison.

Invariants:

- 1) La somme de tous les crédits (messages et processus) est égale à 1.
- 2) Un message de travail a un crédit strictement positif.
- 3) Un processus actif a un crédit strictement positif

Pseudo-code:

{pour p }

var state p : (actif, passif); int si $p = p_0$ donc actif sinon passif;
 cred p : fraction init si $p = p_0$ alors 1 sinon 0;
 ret: fraction int 0; { pour p_0 seulement };

si state = actif \rightarrow $\sqrt{\text{piggybacking}}$

envoyer(mes, cred/2); cred = cred/2;

quand un message mes arrive à $p \rightarrow$

recevoir(mes, p); state p = actif;

cred p = cred p + c ;

Quand l'état de p passe d'actif à passif \rightarrow

envoyer(ret, cred p) à p_0 ;

cred p := 0;

preuves: si l'application n'est pas terminée, p_0 n'annonce pas la terminaison (fairness)

\hookrightarrow 1 ps actif

\rightarrow 0 ps actif mais un mess de travail circule. Dans ce cas, ret \neq 1

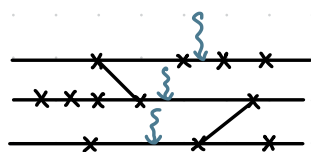
terminée \rightarrow 0 ps actif, pas de messages de travail. ret = 1.

Quand un message ret arrive à $p_0 \rightarrow$

recevoir(ret, c); ret = ret + c ;

si ret = 1 alors < Annoncer >;

Snapshot / Instantané / Prise d'état global.



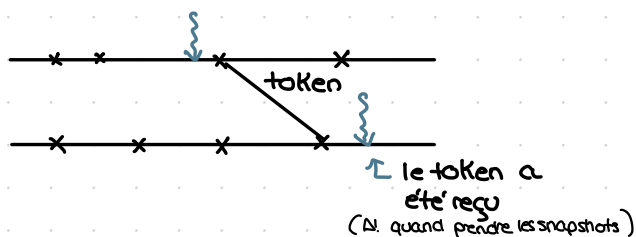
À un moment donné, l'état local de l'application est sauvegardé s^*

Snapshot: $\langle s_1^*, s_2^*, \dots \rangle$

\downarrow
configuration

Détection d'une propriété stable (e.g. terminaison, blockage ...)

Hypothèse: l'état d'un processus p contient la liste des messages émis vers le processus q , depuis le début de l'exécution $send_{p,q}^*$ et ainsi que la liste de messages reçus du processus q depuis le début de l'exécution $received_{p,q}^*$



(1) Def: Un snapshot est possible ssi pour tout processus p et q , $send_{p,q} \subseteq received_{q,p}$
 (l'ensemble des messages envoyés de p à q est inclus dans l'ensemble des messages reçus par q de p)
 (→ le schéma ne peut pas se produire)

(2) Def: Une coupe cohérente L est un ensemble d'événements clos par la relation de causalité (α_p)
 Relation d'ordre $L_1 \subseteq L_2$: Un snapshot correspond à une coupe maximale pour l'inclusion
 éléments pre-shot
 autres éléments post-shot

(3) Def: Un snapshot est significatif si la configuration C^* peut être atteinte par une exécution

On va prouver (1)+(2)+(3):

S^* est un snapshot possible

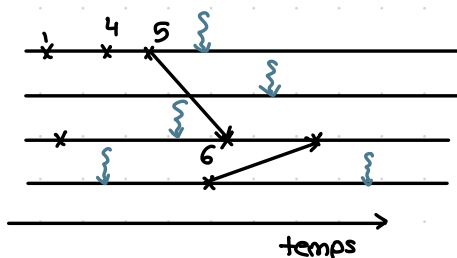
$e, e' \ e \rightarrow e'$

e' est la réception d'un message émis d'un événement e

On doit prouver que si e' est pre-shot, e est pre-shot.

S^* est une coupe cohérente

On suit une exécution qui peut atteindre S^*



Tout message pre-shot est reçu dans un événement pre-shot.

Cours 11 01/04/2026

Algorithme de Chandy-Lanport

Suppose les canaux FIFO

À la demande de prise de snapshot, le snapshot est pris et un message de type particulier le 'token' est envoyé à tous les voisins
justification: les liens sont FIFO

À la réception d'un ^{marqueur} token, si il n'a pas déjà pris son snapshot, le processus le prend et envoie un token à tous ses voisins

var: token booléen initialisé à faux

À la réception d'un signal de la couche sup:

enregistrer l'état local;
token := vrai;
pour tout q dans voisins
envoyer <msg> à q

Algorithme de Lai et Yang (-)

À chaque message envoyé par l'application, le système a un champ indiquant si le message est pre-shot ou post-shot.

Si un processus reçoit un message (appli) qui est post-shot et si ce processus n'a pas pris son snapshot, il le prend immédiatement et ensuite, reçoit le message.

(-)

var token: booléen init à faux

à la réception de init →

enregistrer l'état local;
token := vrai;

lors de l'envoi d'un message m (pour l'application) →

envoyer (m, token);

lors de la réception d'un message <mes, c> →

si c et non token alors | enregistrer l'état local
| token := vrai;

traiter le message <mes>;

Défaillance:

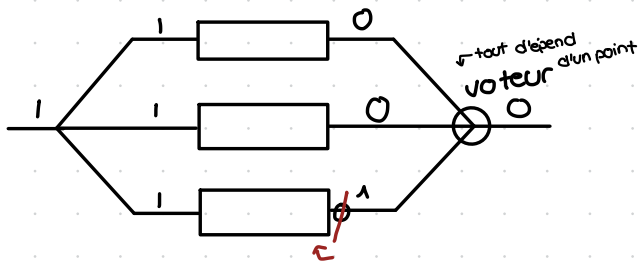
Des processus peuvent mal fonctionner (arrêt total, fonctionnement non conforme à la spécification, malveillance)

Arrêt Complet: A un point d'une exécution 1 ou plusieurs processus cessent de fonctionner

Omissions

Défaillances Byzantines:

Le problème des Généraux Byzantines



⇒ réaliser un "votant" distribué

Problème: Initialement, chaque processus a une valeur entière.

Il s'agit de concevoir un algorithme réparti à l'issue duquel

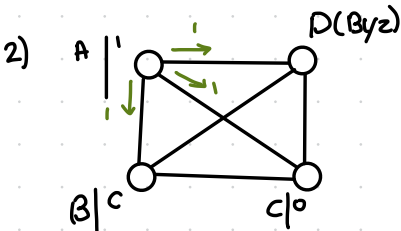
- i) Chaque processus non défaillant décide une valeur ^(faible)
- ii) Deux processus non défaillants ne décident pas des valeurs différentes.
- iii) Si tous les processus non défaillants ont la même valeur initiale, cette valeur est la seule valeur de décision possible.

preuves:

(tous les processus correctes choisissent la même valeur)

- 1) le problème est facile
- 2) le problème est difficile
- 3) le problème est impossible
- 4) Solution

consensus
byzantin



1 ou 0?

	A	B	C	D
A	1	1	1	
B	1	1	1	
C	1	1	0	0
D	1	0	0	0

4) Fisher/Lynch/Satens

Dans un système asynchrone, le problème du consensus est impossible, même avec un seul processus défaillant et un nombre arbitrairement grand de processus correctes.

cas synchrone et défaillances "arrêt total"

flood set: f une borne supérieure au nombre de défaillances ($f < n$)

chaque processus exécute une succession de $f+1$ phases

Initialement un processus initialise un ensemble d'entiers $\omega = \{\text{valeur initiale}\}$;

Chaque processus répète $f+1$ fois:

- envoyer w à tous
- recevoir les w_i et faire $w := w \cup \{w_i\}$;

À la fin de la phase $f+1$,

si $w = \{x\}$ alors décider x

sinon décider (0);

Si tous les processus ont la même valeur initiale, alors la seule valeur de décision possible est cette valeur initiale.

Le cas restant est lorsque les valeurs initiales sont différentes.

si à un point de l'exécution, les ensembles w_j des processus sont égaux, ($w_j = w_i \quad j=i$) alors ils restent égaux jusqu'à la fin du calcul.

Donc tous les processus qui décident décident la même valeur

Il existe au moins une phase au cours de laquelle aucun processus n'est défaillant.

Lors de cette phase, chaque processus reçoit le même ensemble de valeurs.

Lamport, Shostak, Pease