

Cours 3

Terminaison distribuée (fin)

Vagues (début)

Quelques idées

Un chef: répartit le travail et/ou donne les autorisations de travail

Critique?

Quand on a fini on « consulte » tout le monde?

Consultation?

- pourquoi une simple consultation générale ne fonctionne pas?

Détection de la terminaison...

Consultation:

consulter tout le monde (une « **vague** »)

- un (ou plusieurs) initiateur i
- consultation de tout le monde après l'initiateur (= **causalement après** l'ordre de Lamport \leq))
- une « décision » d est un événement qui « récupère » le résultat de la consultation: il est **causalement après** l'arrivée de la consultation sur chaque processus. Le **décideur** peut être ou non l'initiateur.

$$\forall p : \exists a \in E_p : i \leq a \leq d$$

Notations

- On ne s'intéresse qu'au schéma de communication nécessaire pas au contenu des messages (qu'on notera $\langle \rangle$) - ils peuvent contenir toutes les infos
(Dans une consultation on pose une question et on attend une réponse, cette question et les réponses correspondront à la valeur de $\langle \rangle$)
- On considère les événements (émission et réception) de la vague
- On notera (en général) e_p le premier événement de la vague sur p
 - Si cet événement est une émission d'un message, p est un initiateur
 - Sinon cet événement est une réception d'un message de la vague
- Il peut y avoir un ou plusieurs initiateurs
- Il peut y avoir un ou plusieurs décideurs (mais tous les décideurs doivent vérifier $\forall p : \exists a \in E_p : i \leq a \leq d$)
- (On peut facilement diffuser à tout le monde la décision et assurer ainsi que tout le monde soit décideur)

La consultation?

Consultation: on verra comment faire des consultations (vagues-traversées) ensuite.

La consultation est essentiellement un parcours distribué de graphe.

Exemple: un jeton qui passe par tous les processus sur un anneau

- Le jeton passe par tous les sites et retourne à son point de départ:
 - Chaque site maintient un booléen b_p : le jeton est un booléen B ,
 - Quand le jeton arrive sur p : le jeton passe au site suivant avec la valeur $B := B \wedge b_p$.
- Après un tour $B = \bigwedge_{p \in \Pi} b_p$, (en fait $B = \bigwedge_{p \in \Pi} b_p^{t_p}$ où $b_p^{t_p}$ est la valeur de b_p au moment où le jeton passe en p).

Simple consultation?

- Pourquoi une simple consultation: « avez-vous terminé localement? » n'est-elle pas suffisante?

Détection de la terminaison

Il ne suffit pas que la consultation ait eu des réponses TL_p pour tous les p (pourquoi?):

- il faudrait avoir $\exists t : \forall p \in \Pi : TL_p^t$ (pour le même t !)

Dans une consultation simple positive on a obtenu pour chaque $p \in \Pi$, que- au moment t_p de la réponse de p - $TL_p^{t_p}$ était vrai...

Proposez des solutions?

Geler?

Pour résoudre ce problème de la consultation on peut « geler » le système:

Un initiateur qui a localement terminé lance une détection de terminaison

- Une première vague « gèle » le travail (arrêt provisoire sur le processus)
- Quand un décideur sait que le travail est gelé chez tout le monde:
- Une deuxième vague interroge sur l'état au moment du « gel » (oui si terminaison locale)
 - Si le décideur de cette deuxième vague Si tout le monde répond « oui »: un décideur de la vague sait qu'il y a terminaison
 - (Sinon dans tous les cas il y a une diffusion « dégeler » les calculs)

Critique?

Détection de la terminaison

Comment faire?

- On fait des consultations **successives** (le début de la suivante ayant lieu (causalement) après la précédente)
- TG est stable il suffit donc d'assurer qu'on a eu un instant t_0 tel que TG^{t_0} était vrai:
 - On change la question:
 - est-ce que depuis la dernière consultation TL_p est restée vraie? Si tous répondent oui, la dernière consultation ayant eu lieu avant le début de celle qui a obtenu ce « oui »

soit t_p^0 l'instant de la consultation précédente sur p

Soit t_p^1 l'instant de la consultation sur p

On a la garantie que $\forall t (t_p^0 \leq t \leq t_p^1) : TL_p^t$

L'instant t_i du début de la consultation est tel que $\forall p : t_p^0 \leq t_i \leq t_p^1$ et

donc au temps t_i $TL_p^{t_i}$ soit $\bigwedge_{p \in \Pi} TL_p^{t_i}$ et donc TG^{t_i} .

Détection de la terminaison

- Avec un jeton passant par tous les sites
 - Si le site qui reçoit le jeton a été actif depuis le dernier passage: le booléen est mis à « faux », Sinon on transmet le jeton sans changement
 - L'initiateur (inactif depuis le dernier passage) met le jeton à « vrai »)
 - Si le jeton arrive à l'initiateur avec « vrai » la terminaison est détectée

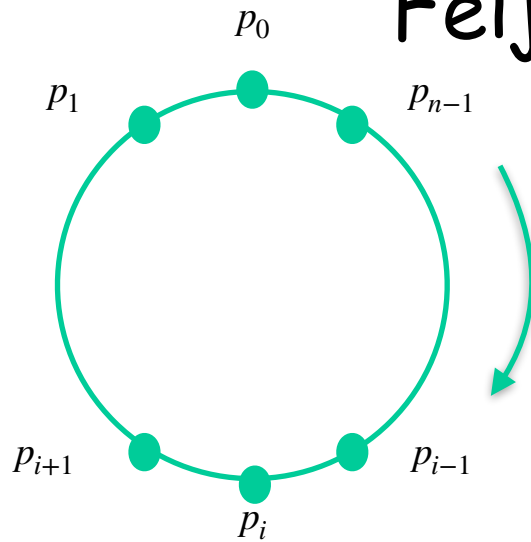
Remarque: on peut faire de nombreuses variantes de l'algorithme

Détection de la terminaison

Remarques:

- Ecrire un algorithme plus précis pour la détection de terminaison.
- Prouver les propriétés de vivacité et de sûreté.
- La forme de la consultation peut être quelconque (pourvue qu'elle vérifie les propriétés précédentes). Un jeton circulant par tous les sites peut faire l'affaire.
- Si au lieu d'avoir un téléphone (ou une communication instantanée) on avait une communication asynchrone (comme le mail) comment pourrait-on faire?

Une variante: algorithme de Dijkstra, Feijen, Van Gasteren



anneau construit sur

$\Pi = \{p_0, \dots, p_{n-1}\}$, dans l'anneau

le jeton passe de p_i à $p_{(i-1) \bmod n}$

p_0 est l'initiateur

Une variante: algorithme de Dijkstra, Feijen, Van Gasteren

$$Term \equiv \forall p : state_p = passive$$

(Ici la terminaison locale est la valeur de la variable

state: $TL_p \equiv (state_p = passive)$)

Construction par invariant:

- On cherche une propriété P qui est initialement vraie et qui reste vraie à chaque « step » telle que quand le parcours a réussi (retour du jeton): $P \Rightarrow Term$
- On va procéder par étapes par des « affaiblissements » successifs pour obtenir l'invariant

Une variante: algorithme de Dijkstra, Feijen, Van Gasteren

Soit t l'identité du processus (p_t) qui a le jeton

- $P_0 \equiv \forall i (N > i > t) : state_{p_i} = passive$
 - P_0 est vrai initialement ($t = N - 1$)
 - Quand le jeton revient en p_0 et $state_{p_0} = passive$, $Term$ est vrai:
 - $(t = 0 \wedge state_{p_0} = passive \wedge P_0) \Rightarrow Term$
 - On en déduit:
 - Règle 1: un processus ne transmet le jeton que s'il est passif
 - Mais P_0 n'est pas un invariant: il peut être invalidé par un processus p_i avec $i \leq t$ qui « réveille » un processus p_j ($j > t$) (le jeton n'est pas encore passé par p_i , mais p_i a réveillé p_j qui a déjà eu le jeton)

Une variante: algorithme de Dijkstra, Feijen, Van Gasteren

Affaiblir P_0 :

- un processus a une couleur black (white sinon) s'il a pu réveiller un autre processus (envoi d'un message): soit $P_1 \equiv \exists j(t \geq j \geq 0) : colors_{p_j} = black$
- $P_0 \vee P_1$ est vraie initialement et:

$$(P_0 \vee P_1)$$

$$\wedge color_{p_0} = white \wedge t = 0 \wedge state_{p_0} = passive$$

$$\Rightarrow Term$$

(Mais $(P_0 \vee P_1)$ n'est pas -encore- un invariant)

Règle 2: un processus qui envoie un message devient black

Une variante: algorithme de Dijkstra, Feijen, Van Gasteren

Il faut aussi assurer que la couleur est « enregistrée »: si un processus est black:

Règle 3: Le jeton devient black s'il rencontre un site black (sauf p_0)

$P_2 \equiv$ le jeton est black

$P_0 \vee P_1 \vee P_2$

- Si pour $t = 0$ le jeton est white P_2 est faux et donc en $t = 0$ si le jeton est white et $color_{t_0}$ est white et $state_{p_0}$ est passive, alors

Term est vrai.

- $P_0 \vee P_1 \vee P_2$ est bien un invariant !

Une variante: algorithme de Dijkstra, Feijen, Van Gasteren

- $P_0 \vee P_1 \vee P_2$ est bien un invariant, mais... un jeton black empêche toute détection de terminaison.
- Règle 4: quand p_0 reçoit un jeton black, il initie un nouveau parcours de jeton avec un jeton white
- Règle 5: un processus devient white dès qu'il a transmis le jeton

Résumé

R1: ne faire quelque chose que si on est passive

R2: si on envoie un message on devient black

R3: un processus black ne transmet que des jetons black (sauf p_0)

R4: quand la vague se termine p_0 commence une autre vague

R5: après envoi du jeton le processus devient white

Une variante: algorithme de Dijkstra, Feijen, Van Gasteren

```
var statep : (active, passive) ;  
    colorp : (white, black) ;  
  
Cpq: { statep = active }  
begin (* p sends a basic message, which is received by q *)  
    colorp := black ;    (* Rule 2 *)  
    stateq := active  
end  
  
Ip: { statep = active }  
begin statep := passive end  
  
Start the detection, executed once by p0:  
begin send ⟨tok, white⟩ to pN-1 end  
  
Tp: (* Process p handles the token ⟨tok, c⟩ *)  
    { statep = passive } (* Rule 1 *)  
begin if p = p0  
    then if (c = white ∧ colorp = white)  
        then Announce  
        else send ⟨tok, white⟩ to pN-1    (* Rule 4 *)  
    else if (colorp = white) (* Rule 3 *)  
        then send ⟨tok, c⟩ to Nextp  
        else send ⟨tok, black⟩ to Nextp ;  
    colorp := white    (* Rule 5 *)  
end
```

Extrait de:
G. Tel *Introduction to distributed algorithms*

Exemple en TLA+

Pour ceux qui aiment la vérification, TLA+ (temporal logic of actions) :

([site de TLA+](#))

- langage de spécification formel développé par L. Lamport
- langage de programmation « pluscal » et traducteur
- modèle (spécification)
- TLC model checker
- TLAPS proof system

EWD840: l'algorithme de Dijkstra, Fijne, van Gasteren)

TLA+ specification of an algorithm for distributed termination detection on a ring, due to *Dijkstra*, published as *EWD 840: Derivation of a termination detection algorithm for distributed computations* (with *W.H.J. Feijen* and *A.J.M. van Gasteren*).

EXTENDS *Naturals*

CONSTANT *N*

ASSUME *NAssumption* $\triangleq N \in \text{Nat} \setminus \{0\}$

VARIABLES *active, color, tpos, tcolor*

Nodes $\triangleq 0 \dots N - 1$

Color $\triangleq \{\text{"white"}, \text{"black"}\}$

TypeOK \triangleq

$\wedge \text{active} \in [\text{Nodes} \rightarrow \text{BOOLEAN}]$

status of nodes (active or passive)

$\wedge \text{color} \in [\text{Nodes} \rightarrow \text{Color}]$

color of nodes

$\wedge \text{tpos} \in \text{Nodes}$

token position

$\wedge \text{tcolor} \in \text{Color}$

token color

Initially the token is black. The other variables may take any "type-correct" values.

Init \triangleq

$\wedge \text{active} \in [\text{Nodes} \rightarrow \text{BOOLEAN}]$

$\wedge \text{color} \in [\text{Nodes} \rightarrow \text{Color}]$

$\wedge \text{tpos} \in \text{Nodes}$

$\wedge \text{tcolor} = \text{"black"}$

Node 0 may initiate a probe when it has the token and when either it is black or the token is black. It passes a white token to node $N - 1$ and paints itself white.

InitiateProbe \triangleq

$\wedge \text{tpos} = 0$

$\wedge \text{tcolor} = \text{"black"} \vee \text{color}[0] = \text{"black"}$

$\wedge \text{tpos}' = N - 1$

$\wedge \text{tcolor}' = \text{"white"}$

$\wedge \text{active}' = \text{active}$

$\wedge \text{color}' = [\text{color} \text{ EXCEPT } ![0] = \text{"white"}]$

H Fauconnier

Algo. Répartie

A node i different from 0 that possesses the token may pass it to node $i - 1$ under the following circumstances :

- node i is inactive or
- node i is colored black or
- the token is black.

Note that the last two conditions will result in an inconclusive round, since the token will be black. The token will be stained if node i is black, otherwise its color is unchanged. Node i will be made white.

$PassToken(i) \triangleq$

$\wedge tpos = i$

$\wedge \neg active[i] \vee color[i] = \text{"black"} \vee tcolor = \text{"black"}$

$\wedge tpos' = i - 1$

$\wedge tcolor' = \text{IF } color[i] = \text{"black"} \text{ THEN "black" ELSE } tcolor$

$\wedge active' = active$

$\wedge color' = [color \text{ EXCEPT } ![i] = \text{"white"}]$

token passing actions controlled by the termination detection algorithm

$System \triangleq InitiateProbe \vee \exists i \in Nodes \setminus \{0\} : PassToken(i)$

An active node i may activate another node j by sending it a message. If $j > i$ (hence activation goes against the direction of the token being passed), then node i becomes black.

$SendMsg(i) \triangleq$

$\wedge active[i]$

$\wedge \exists j \in Nodes \setminus \{i\} :$

$\wedge active' = [active \text{ EXCEPT } ![j] = \text{TRUE}]$

$\wedge color' = [color \text{ EXCEPT } ![i] = \text{IF } j > i \text{ THEN "black" ELSE @}]$

$\wedge \text{UNCHANGED } \langle tpos, tcolor \rangle$

Any active node may become inactive at any moment.

$$\begin{aligned} Deactivate(i) &\triangleq \\ &\wedge active[i] \\ &\wedge active' = [active \text{ EXCEPT } ![i] = FALSE] \\ &\wedge UNCHANGED \langle color, tpos, tcolor \rangle \end{aligned}$$

actions performed by the underlying algorithm

$$Environment \triangleq \exists i \in Nodes : SendMsg(i) \vee Deactivate(i)$$

next-state relation: disjunction of above actions

$$Next \triangleq System \vee Environment$$

$$vars \triangleq \langle active, color, tpos, tcolor \rangle$$

$$Spec \triangleq Init \wedge \Box[Next]_{vars} \wedge WF_{vars}(System)$$

Non-properties, useful for validating the specification with *TLC*.

$$TokenAlwaysBlack \triangleq tcolor = \text{"black"}$$

$$NeverChangeColor \triangleq \Box[UNCHANGED \ color]_{vars}$$

Main safety property: if there is a white token at node 0 then every node is inactive.

$$\begin{aligned} terminationDetected &\triangleq \\ &\wedge tpos = 0 \wedge tcolor = \text{"white"} \\ &\wedge color[0] = \text{"white"} \wedge \neg active[0] \end{aligned}$$

$$\begin{aligned} TerminationDetection &\triangleq \\ &terminationDetected \Rightarrow \forall i \in Nodes : \neg active[i] \end{aligned}$$

Liveness property: termination is eventually detected.

Liveness \triangleq

$$(\forall i \in Nodes : \neg active[i]) \leadsto terminationDetected$$

The following property asserts that when every process always eventually terminates then eventually termination will be detected. It does not hold since processes can wake up each other.

FalseLiveness \triangleq

$$(\forall i \in Nodes : \Box \Diamond \neg active[i]) \leadsto terminationDetected$$

The following property says that eventually all nodes will terminate assuming that from some point onwards no messages are sent. It is not supposed to hold: any node may indefinitely perform local computations. However, this property is verified if the fairness condition $WF_{vars}(Next)$ is used instead of only $WF_{vars}(System)$ that requires fairness of the actions controlled by termination detection.

SpecWFNext $\triangleq Init \wedge \Box[Next]_{vars} \wedge WF_{vars}(Next)$

AllNodesTerminateIfNoMessages \triangleq

$$\Diamond \Box [\forall i \in Nodes : \neg SendMsg(i)]_{vars} \Rightarrow \Diamond (\forall i \in Nodes : \neg active[i])$$

Dijkstra's inductive invariant

Inv \triangleq

$$\vee P0:: \forall i \in Nodes : tpos < i \Rightarrow \neg active[i]$$

$$\vee P1:: \exists j \in 0 \dots tpos : color[j] = \text{"black"}$$

$$\vee P2:: tcolor = \text{"black"}$$

Use the following specification to let *TLC* check that the predicate $TypeOK \wedge Inv$ is inductive for *EWD 840*: verify that it is an (ordinary) invariant of a specification obtained by replacing the initial condition by that conjunction.

CheckInductiveSpec $\triangleq TypeOK \wedge Inv \wedge \Box[Next]_{vars}$

\ * Modification History

\ * Last modified Tue Jun 28 18:17:45 CEST 2016 by merz

\ * Created Mon Sep 09 11:33:10 CEST 2013 by merz

Algorithme de Safra

Et si la communication est asynchrone?

(Non-instantanée: des messages peuvent être en « transit »).

- Pourquoi cela ne marche pas?
- Que faire?

La propriété stable devient (B est le nombre de messages en « transit » (=envoyés et non reçus))

$$\forall p : state_p = passive \wedge B = 0$$

- Rendre la communication « synchrone »: à chaque message on attend un accusé de réception:
 - Il n'y a plus de messages (de l'algorithme de détection) dans les canaux
- Compter (estimer) le nombre de messages en transit
 - B est le nombre de message en transit : algorithme de Safra

Algorithme de Safra

Chaque processus maintient un compteur de messages mc_p

Le jeton a une couleur (black ou white) et le nombre estimé des messages en transit

Règles:

Règle M: quand p envoie un message il incrémente mc_p , quand il reçoit un message il décrémente mc_p

Règle 1: le jeton contient le nombre estimé de messages en transit: un processus p transmet le jeton que s'il est passive et ajoute mc_p à la valeur du jeton

Règle 2: un processus qui reçoit un message devient black

Règle 3: le jeton a aussi une couleur, quand un processus black transmet le jeton, le jeton devient black

Règle 4: La vague échoue si elle revient en p_0 avec un jeton black ou si $mc_{p_0} + q \neq 0$ (q est la valeur du jeton). Si la vague échoue p_0 initie une nouvelle vague

Règle 5: un processus devient white après l'envoi du jeton

Algorithme de Safra

```

var  $state_p$  : (active, passive) ;
     $color_p$  : (white, black) ;
     $mc_p$  : integer init 0 ;

 $S_p$ : {  $state_p = active$  }
begin send  $\langle mes \rangle$  ;
       $mc_p := mc_p + 1$  (* Rule M *)
end

 $R_p$ : { A message  $\langle mes \rangle$  has arrived at  $p$  }
begin receive  $\langle mes \rangle$  ;  $state_p := active$  ;
       $mc_p := mc_p - 1$  ; (* Rule M *)
       $color_p := black$  (* Rule 2 *)
end

 $I_p$ : {  $state_p = active$  }
begin  $state_p := passive$  end

Start the detection, executed once by  $p_0$ :
begin send  $\langle tok, white, 0 \rangle$  to  $p_{N-1}$  end

 $T_p$ : (* Process  $p$  handles the token  $\langle tok, c, q \rangle$  *)
{  $state_p = passive$  } (* Rule 1 *)
begin if  $p = p_0$ 
then if  $(c = white) \wedge (color_p = white) \wedge (mc_p + q = 0)$ 
then Announce
else send  $\langle tok, white, 0 \rangle$  to  $p_{N-1}$  (* Rule 4 *)
else if  $(color_p = white)$  (* Rules 1 and 3 *)
then send  $\langle tok, c, q + mc_p \rangle$  to  $Next_p$ 
else send  $\langle tok, black, q + mc_p \rangle$  to  $Next_p$  ;
       $color_p := white$  (* Rule 5 *)
end
end

```

Extrait de:
G. Tel *Introduction to distributed algorithms*

Quelques remarques

De la communication instantanée à une communication asynchrone (du téléphone au courrier)?

- attendre de recevoir un accusé de réception
- (Méthodes similaires à l'algorithme de Safra)
- (Avoir du broadcast (non-instantané) ne change pas grand chose)

Centralisé- décentralisé:

- On peut aussi envisager des solutions avec un chef qui répartit le travail

Consultation

Vagues...

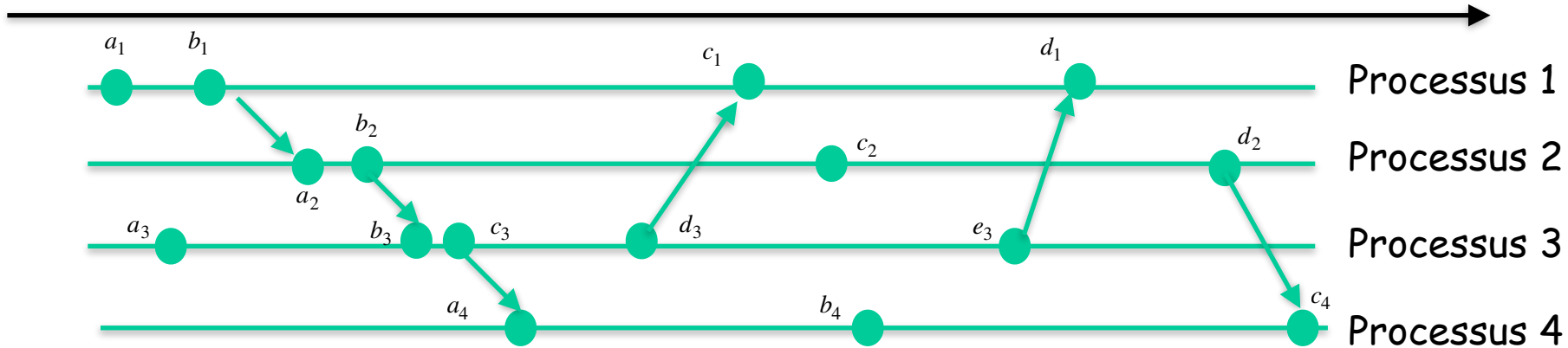
Algorithmes de Vagues

« Consultation » de tous les processus:

- Un **initiateur** (ou plusieurs) lance la consultation
- Un **décideur** (ou plusieurs) récupère le résultat de la consultation. Tous les processus ont été consultés.
 - Une consultation est après le début de la consultation:
 - si c_p est un événement de la consultation sur p , et si i est un événement du démarrage de la consultation alors $i \leq a_p$ (causalité)
 - Tous le monde a été consulté:
 - Si d_p est l'événement de décision (=fin de la consultation) sur le décideur p alors pour tout q il existe un événement de consultation c_q sur q tel que $c_q \leq d_p$
(pour tout q il existe un chemin de communication de l'initiateur i à la décision sur p passant par q)

Vagues

Temps



a_1 initiateur c_1 décision : n'est pas une vague

a_1 initiateur c_4 décision : est une vague

Vague: il existe i événement initiateur et d événement de décision tel que

$$\forall p \in \Pi \exists c_p \in E_p : i \leq c_p \leq d$$

des vagues...

Propagation d'information with Feedback (PIF):

- Les initiateurs ont un message M
- Diffuser ce message et attendre que le message soit diffusé (notification)

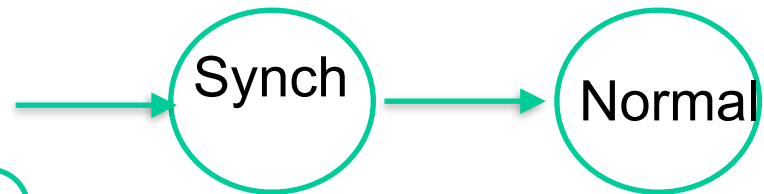
On peut vérifier:

- Tout algorithme de vague permet PIF
- Tout algorithme de PIF est un algorithme de vague
 - Soit une exécution de vague:
 - initiateurs de la vague= initiateurs (ceux qui ont initialement le message M)
 - M est ajouté aux messages de communication de la vague, quand M est reçu par p il est considéré comme diffusé en p
 - Une décision = notification
 - On a bien un algorithme de PIF

des vagues...

Resynch;

- Tous les processus avant d'aller dans l'état normal doivent être simultanément dans l'état synch



- Un algorithme de vague permet Resynch
- Un algorithme de Resynch est un algorithme de vague
 - Décision = état normal

des vagues...

Calcul du minimum

- (X, \leq) un ensemble ordonné
- Chaque processus p a une valeur $r_p \in X$
- Calculer $J = \min\{r_p \mid p \in \Pi\}$

- Tout algorithme de calcul du minimum est un algorithme de vagues
- Tout algorithme de vague permet le calcul du minimum

des vagues...

Chaque processus maintient une variable i_p :

- initialement $i_p = r_p$
- Quand p envoie un message, p ajoute au message sa variable i_p
- Quand p reçoit un message avec la valeur v , $i_p := \min(i_p, v)$
 - On notant $i^{(x)}$ la valeur de i pour l'événement x on a:
 $a \leq b \Rightarrow i^{(a)} \geq i^{(b)}$
 - D'où:
 - $\forall p : i^{(d)} \leq i^{(e_p)} \leq r_p$ (d événement de décision et e_p premier événement de la vague sur p) et donc $i^{(d)} \leq \min\{r_p \mid p \in \Pi\}$
comme $i^{(d)}$ est un des r_p on déduit $i^{(d)} = \min\{r_p \mid p \in \Pi\}$

des vagues

Election:

- choix d'un leader= décider de l'identité d'un processus unique le leader.
- algorithme d'élection de leader:
 - en général algorithme symétrique: tous les processus exécutent le même code (mais ils ont de identités distinctes qu'ils connaissent)
 - en général décentralisé: tout processus peut prendre l'initiative
- Election : Calcul du minimum des identités

des vagues...

Topologie: déterminer la topologie du réseau

Tables de routages

Connectivité

Parcours d'un jeton (traversée : vague particulière sans concurrence des messages)

Exercices...

on suppose que tous les processus ont des identités.
On veut un algorithme distribué qui calcule sur un noeud (tous les noeuds) la liste de toutes les identités. (décision: cette liste est calculée)

(1) montrer qu'un tel algorithme est nécessairement une vague

(2) partant d'un algorithme de vague donner un algorithme permettant de calculer la liste des identités