

THESE

Présentée au

CENTRE DE DEVELOPPEMENT DES TECHNOLOGIES AVANCEES

Pour l'obtention du grade de

MAGISTER EN CYBERNETIQUE

(Option : Architecture des systèmes)

**ETUDE ET EVALUATION D'ALGORITHMES
DE CONTROLE DISTRIBUE :
ELECTION**

Par:

Melle Taboudjemat Nadia

Soutenue le 25 Mai 1992 Devant le jury Composé de

| | |
|------------------------------|-------------------|
| <i>Mr H. Khelalfa</i> _____ | <i>Président</i> |
| <i>Mr M. Benhamadi</i> _____ | <i>Examineur</i> |
| <i>Mr M. Benmihoub</i> _____ | <i>Examineur</i> |
| <i>Mr S. Lezzar</i> _____ | <i>Examineur</i> |
| <i>Mr N. Badache</i> _____ | <i>Rapporteur</i> |

A toute ma famille,
et tous mes amis.

Remerciements

Je remercie **très** sincèrement:

Mr. H. KHELALFA, chef du laboratoire de recherche et de développement en **informatique(LRDI-CERIST)**, qui **me** fait l'honneur de presider le **jury** de cette these.

Mr. M. BENHAMADI, Directeur du CERIST, pour avoir mis a ma disposition les moyens necessaires a la bonne **marche** de mon travail, pour la confiance qu'il m'a **témoignée**, et sa participation **au jury**.

Mr. M. BENMIHOUB, charge de recherche au laboratoire intelligence artificielle **au** CDTA, **Mr. S. LEZZAR**, chef du laboratoire de genie logiciel **au** CDTA, qui ont **accepté** de participer **au jury** de cette these.

Je tiens à exprimer ma reconnaissance à **Mr. N. BADACHE**, charge de cours à l'**USTHB** qui m'a été d'une aide precieuse en dirigeant mon travail.

Je remercie l'ensemble du personnel du CERIST pour leur aide et sympathie.

Que tous ceux qui m'ont **soutenuet** encouragée **trouvent** ici **l'assurance** de mes remerciements **sincères** et amicaux.

RESUME

La nature **complexe** des algorithmes distribués engendre, dans certains cas, une discordance entre les résultats théoriques et les résultats pratiques. Parfois des mises en **oeuvres(même simulées)** permettent de déceler des comportements inattendus, voire incorrects.

Cette thèse est consacrée à l'étude des algorithmes d'élection. Le **problème** est, en effet, fondamental dans le calcul distribué. Les applications qui utilisent ce mécanisme **sont nombreuses** et variées.

Le **but** de cette thèse a **été** de faire un inventaire de ces algorithmes, en les répartissant par classes en fonction de la topologie du **réseau**. Leur spécification en **Estelle**, leur évaluation en utilisant le logiciel Echidna et l'analyse de **leurs performances** ont conduit **aux résultats** suivants:

- confirmer des comportements prévus par la théorie,
- **infirmer** certaines prévisions théoriques,
- déceler des erreurs dans certains algorithmes et proposer **une solution**,
- proposer une classification des algorithmes **évalués** en fonction de leurs **performances**,
- **enfin**, cette étude a permis de mettre en évidence **les** mécanismes **utilisés** dans la conception d'algorithmes d'élection. **Elle** pourrait aider à asseoir une méthodologie de construction d'algorithmes distribués d'élection.

Mots **clés**: Algorithmes distribués - Election - Simulation - Evaluation - Analyse - Complexité - Performance - Conception d'algorithmes distribués.

ABSTRACT

Because of their complexity, distributed algorithms generate in **some** cases, conflicts or differences **between theoretical** and practical results. Implementation of this **type** of algorithms may produce unexpected behaviour.

This thesis is devoted to study election algorithms. The **problem** is fundamental and very elemental in distributed computing. **Our** aim is to make an inventory of the algorithms, divide **them by classes** according to the network topology. The algorithms are specified in Estelle and evaluated under Echidna. The analysis **results are**:

- confirmation of expected **behaviours**,
- invalidation of **theoretical results**,
- detection of errors,
- classification of the algorithms,
- and at last, this **work** brings out the mechanisms and the **techniques** used to design election algorithms. It may **be** helpful for establishing a building methodology of leader finding solutions.

Key words: Distributed algorithms - Election - Simulation - Evaluation - Analysis - Complexity - Performance - Distributed algorithms design.

ملخص:

ان الطبيعة المعقدة للخوارزميات الموزعة تؤدي في بعض الحالات الى عدم التنسيق بين النتائج النظرية والنتائج التطبيقية . ففي بعض الاحيان يعطي تحقيق هذه الخوارزميات نتائج غير منتظرة أو غير صحيحة. و نظرا لأهمية الخوارزميات الإنتخابية في الحساب الموزع واستعمالها في تطبيقات عديدة و متنوعة خصصت هذه الأطروحة لدراستها.

وتهدف الدراسة إلى جمع هذه الخوارزميات و تصنيفها الى أقسام حسب شكل الشبكة. وبعد ترجمتها الى لغة برمجة الخوارزميات المتوازية ESTELLE و تقييمها عن طريق البرمجة ECHIDNA، ثم تحليل مدى فعاليتها، تم التوصل إلى النتائج التالية:

- ايجاد نتائج غير معارضة للتنبؤات النظرية،
- تنفيذ بعض هذه النتائج،
- الكشف عن بعض الاخطا في هذه الخوارزميات و اقتراح حلول لها،
- اقتراح تصنيف الخوارزميات التي تم تقييمها حسب فعاليتها،
- وأخيرا، سحت هذه الدراسة بإبراز الميكانيزمات المستعملة لإنشاء الخوارزميات الإنتخابية وقد يساعد ذلك على إقامة منهجية لتصميم هذه الخوارزميات.

الكلمات الدالة :

الخوارزميات الموزعة - انتخاب - تقييم - تحليل - تعقيد - فعالية - تصميم الخوارزميات الموزعة.

S O M M A I R E

| | |
|--|----|
| INTRODUCTION | 1 |
| CHAPTER 1: NOTIONS FONDAMENTALES DE L'ALGORITHMIQUE DISTRIBUEE | |
| I CONTRAINTES DE LA DISTRIBUTION | |
| 1 Absence De Mémoire Commune | 3 |
| 2 Absence De Référentiel Temporel Global | 3 |
| 3 Non Fiabilité De L'environnement D'exécution | 4 |
| II QUELQUES CONCEPTS DE BASE DE L'ALGORITHMIQUE DISTRIBUEE | |
| 1 Le Jeton Circulant | 4 |
| 2 L'estampillage | 4 |
| 3 Le Calcul Diffusant | 4 |
| 4 Le Transfert De Connaissances | 5 |
| III LES QUALITES D'UN ALGORITHME DISTRIBUE | |
| 1 Hypothèses Sur Le Réseau De Communication | 5 |
| 2 Degré De Répartition Isymétrie De L'algorithme | 6 |
| 3 Résistance Aux Pannes | 6 |
| 4 Etat Global Ou Local | 6 |
| 5 Complexité En Trafic Et En Temps | 7 |
| IV LES ALGORITHMES D'ELECTION | 7 |
| CHAPTER 2: EVALUATION SOUS ECHIDNA | |
| I MOTIVATIONS ET OBJECTIFS | 10 |
| II LES OUTILS D'EVALUATION | |
| 1 Le Logiciel | |
| 1.1 Langage de spécification Estelle | 11 |
| 1.2 Présentation d' Echidna | 17 |
| 2 Le Matériel | 18 |
| III MISE EN OEUVRE DE LA SIMULATION SOUS ECHIDNA | |
| 1 Le Modèle D'élection | 19 |
| 2 Les Critères D'évaluation | 19 |
| 3 Le Modèle Estelle Implémenté Sous Echidna | 19 |
| 4 La Simulation Sous Echidna | 20 |
| 5 Collecte Des Informations De Mesure | 21 |
| a) Évaluation du temps | 21 |
| b) Le principe d'évaluation: | 23 |

CHAPTER 3: ETUDE DES ALGORITHMES D'ELECTION

| | | |
|-------|---|----|
| I | DEFINITIONS ET TERMINOLOGIE | 25 |
| II | LANGAGE DE DESCRIPTION ET QUELQUES NOTATIONS | 26 |
| III | PRESENTATION D'UN ECHANTILLON D'ALGORITHMES | |
| 1 | ALGORITHMES D'ELECTION POUR TOPOLOGIE EN ANNEAU | |
| 1.1 | Introduction | 27 |
| 1.2 | Hypotheses | 28 |
| 1.3 | Algorithme De Chang Et Roberts | 28 |
| 1.4 | Algorithme De Dolev, Klawe Et Rodeh | 30 |
| 1.5 | Algorithme De Hirshberg Et Sinclair | 32 |
| 1.6 | Algorithme_P De Korach, Rotem Et Santoro | 34 |
| 1.7 | Algorithme De Fredreickson Et Lynch Pour Un Reseau Synchrone | 35 |
| 1.8 | Algorithme De Spirakis, Tampakas Et Athanasios | |
| 1.9 | Algorithme De Frankin | 37 |
| 1.10 | Algorithme_D De Bodlaender Et Van Leeuwen | 37 |
| 1.11 | Algorithme De J.Van Leeuwen Et R.B.Tan | 38 |
| 2 | ALGORITHMES POUR TOPOLOGIE COMPLETE | |
| 2.1 | Introduction | 38 |
| 2.2 | Hypotheses | 39 |
| 2.3 | Algorithmes De Garcia-Molina | |
| 2.3.1 | Principe De L'algorithme Brutal | 40 |
| 2.3.2 | Principe De L'algorithme(invitation) | 40 |
| 2.4 | Algorithme De Korach, Moran Et Zaks | 40 |
| 2.5 | Algorithmes De Afek Et Gafnı | |
| 2.5.1 | Principe | 41 |
| 2.5.2 | Algorithme A | 42 |
| 2.5.3 | Algorithme B | 44 |
| 2.5.4 | Algorithme C | 46 |
| 2.6 | Algorithme De Loui , Matsushita Et West | 47 |
| 2.7 | Algorithme De Attiya, Leeuwen Et Shmuel | 50 |
| 3 | ALGORITHMES D'ELECTION POUR TOPOLOGIE QUELCONQUE | |
| 3.1 | Introduction | 50 |
| 3.2 | Hypotheses. | 51 |
| 3.3 | Les Sites Sont Dotes D'identites Distinctes | |
| 3.3.1 | Algorithme De Helary, Maddi et Raynal (exploration en profondeur) | 51 |
| 3.3.2 | Algorithme De Helary, Maddi et Raynal (exploration en largeur) | 55 |
| 3.3.3 | Algorithme De Sven Skyum | 57 |
| 3.3.4 | Algorithme De Lavallee Et Lavault | 57 |
| 3.4 | Algorithmes Pour Reseaux Anonymes | |
| 3.4.1 | Algorithme De Christian Lavault | 59 |
| 3.4.2 | Algorithme De Yossi Et Afek | 60 |

| | | | |
|----------------------------------|-----------------------------|--|------------|
| IV | TABLEAU RECAPITULATIF | | |
| | 1 | Hypotheses Sur Le Reseau De Communication | 62 |
| | 2 | Hypotheses Sur Le Reseau Des Processus | 62 |
| | 3 | Degré De Répartition(ou Symétrie) | 63 |
| | 4 | Connaissance De L'état Global Ou Local | 63 |
| | 5 | Topologie En Anneau | 63 |
| | 6 | Topologie Complete | 64 |
| | 7 | Topologie Quelconque | 64 |
| V | LES TECHNIQUES UTILISEES | | |
| | 1 | Election quand Les processus sont dotes d'identites distinctes | 65 |
| | 2 | Technique d'élection dans les réseaux anonymes | 68 |
| CHAPTER 4: EVALUATION ET ANALYSE | | | |
| I | EVALUATION DES PERFORMANCES | | |
| | 1 | Topologie En Anneau | |
| | | 1.1 Algorithmes De Chang Et Roberts | 72 |
| | | 1.2 Algorithme De Dolev, Klawe Et Rodeh | 74 |
| | | 1.3 Algorithme De Hirshberg Et Sinclair | 75 |
| | | 1.4 Algorithme_P De Korach, Rotem Et Santoro | 77 |
| | 2 | Topologie Complete | |
| | | 2.1 Algorithmes De Afek Et Gafni [AFE 851 | |
| | | 2.1.1 L'algorithme A, | 80 |
| | | 2.1.2 L'algorithme B | 81 |
| | | 2.1.3 L'algorithme C | 82 |
| | | 2.2 Algorithme De Loui , Matsushita Et West | 83 |
| | 3 | Topologie Quelconque | |
| | | 3.1 Algorithme De Helary, Maddi et Raynal (exploration en profondeur) | 84 |
| | | 3.2 Algorithme De Helary, Maddi et Raynal (exploration en largeur) | 88 |
| | | 3.3 Comparaison Des Deux Explorations (en profondeur et en largeur) | |
| | | 3.3.1 Election Collective | 91 |
| | | 3.3.2 Election Individuelle | 93 |
| II | ETUDE COMPARATIVE | | |
| | 1 | Topologie En Anneau | |
| | | 1.1 Election Collective | 95 |
| | | 1.2 Election Individuelle | 96 |
| | 2 | Topologie Complete | |
| | | 2.1 Election Collective | 96 |
| | | 2.2 Election Individuelle | 97 |
| CONCLUSION | | | 115 |
| BIBLIOGRAPHIE | | | |
| ANNEXE | | | |

INTRODUCTION

Ces deux dernières **décennies** ont vu croître la demande en puissance de calcul de telle manière **que** la réponse la plus immédiate et la moins **chère** pour les **constructeurs** a été de multiplier les **ressources** en calcul. On a vu alors la croissance rapide des **réseaux** locaux d'ordinateurs à taille réduite et à grands débits de transmission qui ont favorisé l'introduction du concept de système distribué. Un système distribué est un système **unique** dont les composants (**données**, contrôle,...) sont distribués **sur un** ensemble d'équipements informatiques **géographiquement** repartis et **interconnectés** par un **réseau** de communication. Il ne s'agit plus, en effet, d'offrir un moyen de communication et une interface entre les différents systèmes **éventuellement hétérogènes** mais plutôt de réaliser un **système** unique dont les composants sont repartis et coopèrent à la réalisation d'une **tâche commune**. En fait, c'est l'association entre l'évolution des **réseaux** locaux et le développement de la programmation parallèle qui a induit l'apparition d'une **nouvelle** algorithmique **basée** sur l'échange de messages entre **activités** coopérantes. Dans les systèmes repartis, ces algorithmes parallèles sont appelés algorithmes distribués.

Dans l'algorithmique distribuée il y a des algorithmes qui réalisent généralement les **fonctions** de **contrôle de base** inhérentes à **tout système** informatique (exclusion **mutuelle** par exemple), et des algorithmes qui **réalisent** des fonctions de contrôle inhérentes à la distribution (détection de la terminaison, l'élection, l'interblocage, la gestion de **données dupliquées**).

Indépendamment du support **sur lequel** ils **s'exécutent** et leur **finalité** propre, les algorithmes distribués peuvent être définis par l'équivalence suivante [RAY 85]:

algorithme distribué = processus + messages.

Concevoir un algorithme distribué est un exercice difficile, se convaincre de son **bon** fonctionnement est dans l'état **actuel** de l'art plus ardu. En effet, la nature **complexe** de ce **type** d'algorithmes (**parallélisme**, non **déterminisme**, communication par messages,...) ne **permet** pas de **s'assurer** facilement de leur bon fonctionnement. **C'est ce constat** de la difficulté de maîtrise qui justifie les nombreux travaux d'équipes **sur** la vérification et l'évaluation des algorithmes distribués.

Des outils spécifiques ont **été conçus** pour l'évaluation. Celle-ci **permet** de **mesurer** les performances, elle offre **également** un regard original **sur certains** aspects des algorithmes distribués, et **permet** d'en mieux saisir les **particularités**.

Une fonction de contrôle (le mot contrôle s'oppose ici à **calcul**) qui résout le **problème** de distinction d'un processus (ou site) parmi un ensemble de processus coopérant à la réalisation

d'un but **commun**, peut prendre deux formes distinctes selon le but recherche [RAY 85]:

- la realisation d'une exclusion mutuelle lorsque le privilege **doit être** equitalement attribue au processus (et ne peut donc **être possédé** indefiniment par un **même** processus).

- l'election d'un processus par les autres processus lorsque **le** privilege peut être attribue une fois pour toutes.

Cette these est **consacrée** a l'etude et a l'evaluation des algorithmes distribues **traitant** le probleme d'election. **Elle** est organisee en quatre chapitres.

Le chapitre I est un chapitre preliminaire, dans lequel sont **présentés** principalement, quelques notions de base **inhérentes** a la distribution, ainsi que des exemples montrant l'utilisation des algorithmes d'election pour resoudre des problemes pouvant se produire dans un **système distribué**.

Le chapitre II est **consacré** a l'etude du langage de specification de **systemes** distribues Estelle [BUD 871, et de l'outil de simulation Echidna [JEZ 89a, JEZ 89b]. **Les** conditions de simulation y sont **décrites**.

Le chapitre III **comporte** l'etude detaillee des algorithmes d'election pour differentes topologies de reseaux de communication : anneau, complete et quelconque.

Le chapitre IV quant a lui, est reserve à l'analyse des **résultats**.

Le document est **achevé** par une conclusion et des annexes.

CHAPITRE 1

NOTIONS FONDAMENTALES DE **L'ALGORITHMIQUE** DISTRIBUEE

Dans ce chapitre, est donné un **aperçu** sur l'environnement de conception des algorithmes distribués, puis, sont présentés quelques outils, **conçus** dans ce **domaine**, et nécessaires à la construction de ce type d'algorithmes. La **nécessité** d'avoir des algorithmes d'élection est **montrée** sur des exemples.

I CONTRAINTES DE LA DISTRIBUTION

Le développement d'un algorithme distribué doit tenir **compte** des contraintes d'environnement qu'impose la structure matérielle sur laquelle il s'exécute [COR 81, RAY 85]: l'absence de **mémoire** commune, l'absence d'un référentiel temporel global, la panne physique des composants matériels du système.

1 Absence De **Mémoire** Commune

La structure matérielle d'un système distribué est **constituée** d'un ensemble de machines indépendantes et communicant à travers des lignes de communication. Il n'y a donc pas de **mémoire** commune qui puisse servir à la communication et à la **synchronisation** des composants de l'algorithme ou du système distribué. La seule possibilité de communication des composants reste l'**échange** de messages entre les **différents** sites du système.

Une conséquence **directe** de la communication par messages (problèmes de **décalage** de transmission) est qu'un site **donné** ne peut connaître que d'une manière approchée l'état d'un autre site. **Donc** l'**état** global d'un système distribué ne peut être **perçu** avec exactitude.

2 Absence De **Référentiel** Temporel Global

La plupart des problèmes de **contrôle** (de synchronisation en général) dans un système se résolvent par un **ordonnement** entre les événements de celui-ci. En effet, dans une structure distribuée les horloges physiques ne peuvent servir de

referentiel temporel global a **cause** de leurs derives les **unes par rapport aux autres**. La notion de temps logique [LAM 781 a été introduite pour la datation et l'ordonnancement entre les Cvenements du **systeme**.

3 Non **Fiabilité** De L'environnement **D'exécution**

Un des principaux apports des solutions distribuees est de pouvoir maintenir le **systeme en fonctionnement** en cas de pannes physiques (**portant sur les sites physiques ou sur les moyens de communication**) de quelques composantes du **systeme**. Un algorithme distribue **doit tenir compte** de ces problemes.

II QUELQUES CONCEPTS DE BASE DE **L'ALGORITHMIQUE** DISTRIBUEE

Avec l'algorithmique distribuee sont apparus des **techniques nouvelles**: le **jeton** circulant, l'estampillage, le **calcul** diffusant et le transfert de connaissances.

1 Le **Jeton** Circulant

Consiste en une marque particuliere appelee privilege **et** qui **circule sur un anneau** de processus. Le processus disposant du privilege entreprend alors **une action** correspondante au **protocole exécuté**. Le privilege **étant matérialisé** par une **certaine** configuration des variables des processus ou bien **par un message** special [LEL 771.

2 L'estampillage

En absence d'un **référentiel** temporel global, l'estampillage reste à la **base** de la plupart des algorithmes dans lesquels l'établissement d'un ordre partiel ou total entre les **événements** du **systeme** est **nécessaire**. Lamport [LAM 781 a propose le **principe** de l'estampillage (encore **appelé principe** des horloges logiques) qui **permet** de dater les **événements** relatifs aux interactions **entre** les processus (émission/reception de messages] de **manière cohérente**.

3 Le **Calcul** Diffusant

Le **principe** du **calcul** diffusant a **été** introduit par Dijkstra et Scholter [DIJ 80]. Les processus sont **connectés de manière** quelconque par les voies de communication, à l'exception d'un processus particulier, qui initialement ne peut **qu'émettre** des messages. Le **principe** du **calcul** diffusant est alors le suivant: Initialement seul ce processus peut **mettre** des messages, et **ensuite** tout autre processus ne peut en **émettre** que s'il en a **lui-même reçu**. Ce **principe** est souvent utilisé **sur un** arbre de recouvrement du graphe des processus. Chacun des processus **réexpédie** le message **reçu** vers chacun de ses **successeurs** et attend la reception des reponses avant d'envoyer a son tour une reponse a son **père**. Le **calcul** se termine lorsque le processus **racine** a **reçues réponses** de tous ses fils.

4 Le Transfert De Connaissances

Les processus constituant un algorithme distribue, bien qu'ils **coopèrent** à la realisation d'une **tâche** commune ne disposent initialement que d'informations locales (**généralement** des informations sur les voisins). Pour une prise de decision, un processus a besoin d'acquérir des informations qu'ont les autres processus **afin** de constituer une vue partielle ou **globale** sur l'état du systeme. C'est le **principe** du transfert des connaissances introduit dans [HEL 86a].

La plupart des **protocoles** de transfert des connaissances se basent sur la technique de la vague. Un processus desirant **acquérir** des informations sur le systeme diffuse a un processus voisin un message approprié. Ce message sera **enrichi**, a sa reception, des connaissances locales du site recepteur et diffuse aux voisins de ce dernier, jusqu'a atteindre l'ensemble des processus presents dans le systeme. Le risque de multiplication excessive des messages de transfert des connaissances peut **être évité** par une technique de **contrôle** de ce type de transferts [HEL 86a].

III LES QUALITES D'UN ALGORITHME DISTRIBUE

Les **qualités** d'un algorithme distribué se mesurent par les criteres suivants [RAY 85]:

1 Hypotheses Sur Le **Réseau** De Communication

Un algorithme distribue pose toujours des hypotheses sur le maillage du **réseau** de communication d'une part, et sur le **comportement** des voies de communication d'autre part.

a) Maillage du réseau de communication

La plupart des algorithmes distribues supposent une topologie particuliere, **peu** d'entre **eux** s'exécutent sur un reseau quelconque. Les plus tout-antes sont:

- Maillage en anneau: un processus ne connaît que ses deux voisins immediats (gauche et droit). L'anneau forme par les processus peut être uni- ou bi-directionnel selon que les messages circulent dans un seul sens ou dans deux sens.
- Maillage en arbre: structure hierarchique ou a tout noeud est associe **un** processus qui ne peut communiquer qu'avec son pere, et le processus racine ne communique qu'avec ses fils.
- Maillage complet: tout processus peut communiquer avec tout autre processus.

Une topologie quelconque est une structure **MIXTE** des topologies précédentes.

b) Propriétés comportementales des liaisons

La plupart des algorithmes distribués font des hypothèses sur le **comportement** des voies de communication. Les plus fréquemment **rencontrées** s'ennocent comme suit:

- Il n'y a pas de duplication de messages au cours de la transmission.
- Les messages ne sont pas alter-es.
- Il n'y a pas de dessequencement.
- Le **délai** d'acheminement des messages bien qu'aleatoire, il est **fini**; tout message est **reçu au bout** d'un temps **fini**. Il n'y a pas de perte de messages.

2 Degré De Repartition (**Symétrie** De L'algorithme)

Le **degré** de repartition **permet** de définir les niveaux de symétrie des **rôles joués** par les différents processus du système. Dans [RAY 85], on distingue plusieurs niveaux de symétrie. Ils sont **présentés** ici par ordre croissant:

- Non-symétrie: les processus **n'exécutent pas tous le même** texte (exemple du système client-serveur).
- Symétrie de texte: les processus exécutent **tous le même texte**, sauf que l'identité du processus apparaît dans le texte. Ceci peut conduire à des comportements différents.
- Symétrie forte: les processus exécutent **tous le même** texte sans aucune référence à leurs identités. **Selon** les messages **reçus**, ils peuvent avoir des comportements identiques ou différents.
- Symétrie **totale**: les processus exécutent **tous le même** texte et engendrent des comportements identiques.

3 Résistance Aux Pannes

L'algorithme distribué est d'autant plus intéressant qu'il résiste **aux pannes**. Ce **critère** est lié au précédent vu que la panne d'un site dans un algorithme non symétrique risque **d'être** fatale à l'algorithme.

4 Etat Global Ou Local

Un algorithme est plus performant si pour un processus de l'algorithme, la prise de décisions n'utilise **pas** un état global. Outre la réduction du nombre de messages échangés **cela**

peut permettre d'assurer une meilleure résistance aux pannes.

5 Complexité En Trafic Et En Temps

Pour une même fonction réalisée, un algorithme distribué est d'autant plus intéressant qu'il minimise le nombre de messages (trafic) engendrés sur le réseau. Le temps nécessaire à la réalisation d'une opération de contrôle est un autre critère de mesure de performance d'un algorithme distribué (temps d'élection).

IV LES ALGORITHMES D'ELECTION

L'élection est utilisée dans un système réparti quand un site est appelé à jouer un rôle de supervision, ou quand un site, sans être le coordinateur de l'ensemble des sites, joue un rôle particulier et que la panne de ce site est un événement auquel l'application doit réagir. Les applications qui utilisent ce mécanisme sont très nombreuses. Dans ce qui suit, sont donnés en bref, quelques exemples d'utilisation de ce type d'algorithmes.

Exemple 1:

C'est par exemple le cas dans les systèmes téléinformatiques, où la gestion du réseau de communication nécessite un ensemble de protocoles assurant un service fiable et de bonne qualité. Parmi ses fonctions, la couche réseau (3^e niveau de la structure de réseau définie par la norme ISO) consiste à sélectionner le chemin que doit suivre un message entre deux adresses. Le routage est un des éléments primordiaux qui permet la reconfiguration du système en cas de panne, d'anomalie ou tout simplement de congestion. Parmi les techniques de routage, il y a celles qui sont caractérisées par un centre de contrôle [PUJ 871 dont le rôle consiste à "élaborer" des tables de routage de façon statique ou dynamique. La défaillance du nœud de contrôle nécessite la désignation d'un autre nœud pour assurer cette fonction.

Exemple 2:

De même, dans certains réseaux locaux en anneau [PUJ 871 une station a un rôle important pour l'administration de l'anneau (contrôle d'accès aux ressources de communication, récupération des messages perdus, etc.). C'est le cas par exemple, du réseau local "token ring" d'IBM. Cette politique est décrite par la norme IEEE 802.5. Lorsqu'une station est arrêtée ou est en panne, elle est automatiquement retirée de l'anneau. Si c'est l'administrateur, il faut en élire un nouveau.

Exemple 3:

Certaines applications necessitent l'utilisation de protocoles de diffusion fiable. Les protocoles de diffusion sont des algorithmes distribues qui permettent l'envoi d'un **même** message à un ensemble de processus destinataires en une seule operation. L'ordre de reception des messages au niveau des sites peut **être** quelconque, FIFO, global, causal ou atomique.

Deux protocoles de diffusion caracterises par un ordonnancement global sont **donnés** dans [CHA 841 et [NAV 881. L'ordonnancement global assure que les messages sont **reçus** dans le **même** ordre **par tous les** sites recepteurs. **Ceci peut être utilisé, par** exemple, pour simplifier la conception de mecanismes du **contrôle** de concurrence et des procedures de recouvrement dans les **systèmes** de bases de donnees distribues.

Le **protocole** de diffusion ordonnee propose **par** Chang et Maxemchuk [CHA 841 est base sur un site particulier appele "**site jeton**". Ce dernier joue le role d'un entonnoire **par** lequel transitent **tous** les messages diffuses. Le site **jeton** est charge d'assurer que les messages sont **tous reçus**, et, dans le **même** ordre **par** les destinataires. La fiabilite de l'algorithme **est** assuree en attribuant, de **façon** circulaire, ce privilege a un ensemble de sites qui constituent la "**liste jeton**". Lorsqu'un site **détecte** une defaillance, il invite les autres sites a reconstituer une nouvelle "**liste jeton**" et a **élire** un nouveau "**site jeton**".

Exemple 4:

Dans les bases de donnees **réparties**, le **problème** du **contrôle** de la concurrence, **dû** a l'execution parallele de transactions accedant a des donnees dupliques, peut Ptre **résolu** en associant à chaque ensemble de donnees (**dupliquées** sur les divers sites du **réseau**) un site particulier. Le **rôle** de ce site, appele site primaire [BER 81] est de **contrôler** les **accès** qu'effectuent **les** transactions de **façon** a assurer la coherence mutuelle de ces donnees (respect des contraintes d'**intégrité**) et la convergence de leurs copies vers une **même** valeur.

Exemple 5:

Un autre exemple est celui de l'algorithme **distribué** de **synchronisation** d'horloges TEMPO **conçu** pour **synchroniser** les horloges du **système** UNIX **distribué** de Berkeley. Tempo est base sur une structure maitre-esclave **composée** de processus (time deamons) s'executant sur un **réseau** local . de machines individuelles; un processus **par** site. Un processus maitre (master time daemon) **mesure** la difference entre l'heure **indiquée** **par** l' horloge locale de son site et **celles** indiquees **par tous les** autres sites. Le maitre **calcule** l'heure **globale** du **réseau** comme **étant** la moyenne des valeurs fournies par toutes **les** autres **horloges**(non defaillantes) du reseau. Il transmet **ensuite** a chaque processus **esclave** la valeur **avec** laquelle ce

dernier corrigera l'horloge locale a son site de residence.

Ce processus de **synchronisation** d'horloge est **répété** périodiquement. Une description détaillée de Tempo est donnée dans [GUS 84, GUS 85]. Pour **s'assurer que** Tempo **fournit** un service **continu** et fiable, il est nécessaire d'implémenter un algorithme d'élection [GUS 86] qui élira un nouveau maître si le site sur lequel s'exécute le processus maître courant **tombe** en panne, ou si le réseau des processus est partitionné en **sous-réseau**.

Dans ces systèmes, la panne du site jouant un rôle particulier **nécessite** l'élection d'un nouveau site pour le remplacement du site défaillant et le maintien du système en fonctionnement.

Exemple 6:

Un autre cas nécessitant l'élection d'un site parmi les autres, est celui de plusieurs sites d'un réseau coopérant à la réalisation d'une **tâche** commune, où un des composants **tombe** en panne [GAR 82]. Il peut **être** alors nécessaire de **reconfigurer le système afin, soit** de poursuivre la **tâche** en question, **soit** de sélectionner de nouvelles **tâches** qu'il faudrait assigner aux différents sites. La reconfiguration du système **doit être** à la charge d'un site unique, **sinon** il est possible d'aboutir à une confusion. **Donc**, pour assurer la reconfiguration **correcte** du système, il faut **élire** un site coordinateur qui s'en chargera.

CHAPITRE 2

EVALUATION **SOUS** Echidna

I MOTIVATIONS ET OBJECTIFS

Le processus de développement **d'un** algorithme distribue passe toujours par **l'étape** de validation-évaluation de celui-ci. Cette validation-évaluation se fait **généralement** par une preuve formelle des **propriétés** de l'algorithme et une analyse théorique de **complexité** (**trafic**, temps) et **ensuite** par une simulation de l'algorithme sur une machine centralisée.

La simulation [GRO 87] est une technique largement utilisée dans **tous** les domaines lorsque la complexité d'un problème rend les tentatives de solution analytique inopérantes. La vérification d'algorithmes **distribués même** simples rentre dans ce cadre du fait de l'explosion combinatoire des comportements (ensemble infini de comportements) de ces algorithmes. Toutes les **méthodes** de validation automatique ont leurs limites; la validation est toujours réduite à une **partie** du **comportement**, et la **partie** prise en **compte** par les simulations est souvent plus importante que les parties prises en **compte avec** d'autres méthodes. La simulation d'algorithmes **distribués** peut aussi aider à la conception de ce type d'algorithmes. En effet la simulation peut traiter des **modèles** d'algorithmes **assez complets** et **permet** de **détecter** efficacement des erreurs **éventuelles** sur un sous-ensemble des comportements possibles du **protocole**. Elle **permet** de compléter la preuve formelle, et de tirer des conclusions plus **générales** quant au **comportement** de l'algorithme **distribué** réel. En effet, l'implantation d'un prototype dans un environnement **d'exécution réel** n'est **pas** toujours possible et de toute **façon** les conclusions qui peuvent en être **tirées** dépendent étroitement de la mise en oeuvre [ADA 88].

Un grand nombre d'algorithmes **distribués d'élection** ont été **publiés** pour diverses topologies de **réseaux** de communication et **avec** divers concepts de base. Il serait intéressant de faire un inventaire de ces algorithmes en les répartissant par **classe** en **fonction** de la topologie du **réseau**, et de les implémenter (**simuler**) et d'analyser leurs performances **afin** de les comparer.

Les resultats escomptes sont:

- **trouver** les meilleurs algorithmes,
- confirmer ou **infirmer** les previsions theoriques quant a leur bon fonctionnement et leur performance,
- **mettre en evidence** des comportements eventuels, non prévus par la **théorie**,

L'évaluation des algorithmes est **réalisée** a l'aide du logiciel Echidna [JEZ 89a, JEZ 89b]. C'est **un système conçu** dans le but de fournir les outils nécessaires à l'experimentation **et/ou** la simulation d'algorithmes distribues. Echidna est **construit autour** d'un compilateur d'algorithmes distribues **décrits** dans le langage de specification Estelle CISO 891. Echidna **permet** principalement la simulation du **parallélisme**.

II LES OUTILS D'EVALUATION

1 Le Logiciel

L'évaluation **consiste** a specifier les algorithmes **étudiés** en langage Estelle de specification de systemes distribues. Ces algorithmes sont **ensuite implémentés(simulés)** sous le systeme Echidna qui a **été développé** sur un hypercube IPSC/2 d'INTEL, sur des reseaux de stations SUN, et **sur les** micro-ordinateurs PC.

1.1 Langage de spécification Estelle

La **difficulté** a **prévoir toutes** les configurations possibles lors de la conception d'un algorithme distribue, **l'ambiguité**, la **longueur**, la lourdeur et les differences d'interpretation d'une description en langage naturel montrent **l'interêt** a specifier les algorithmes distribues dans un langage de description formelle [JAR 87]. Pour **cela**, de **nombreux** modeles de description des algorithmes (ou des protocoles) ont **été développés** de façon a pallier aux problemes ci-dessus. Ils peuvent **être** classes en trois grands types de techniques:

1. Les systemes a transitions, c'est-a-dire d'une part les **automates** d'etats finis et d'autre **part**, les **modèles** de graphes à marquage comme les **réseaux** de Petri [RAZ 84, JAR 85];
2. Les **langages** de haut niveau sequentiels qui ont **été** plus ou moins adaptes au **domaine** par la possibilite d'exprimer **l'arrivée d'événements** externes [BOU 88].
3. Les techniques algebriques et logiques [AUD 88].

Estelle est un **langage basé** sur le modele de transitions d'etats etendu par le langage Pascal **avec** certaines restrictions plus des elements d'expression des **comportements** paralleles. Une specification Estelle est la description du **comportement**

d'un système **fermé**; c'est-à-dire sans interaction **avec l'extérieur**.

Tous les objets Estelle sont fortement types; la notion de typage Pascal est étendue aux objets parallèles. On parle alors d'un type de canal, de comportement et d'un type de module.

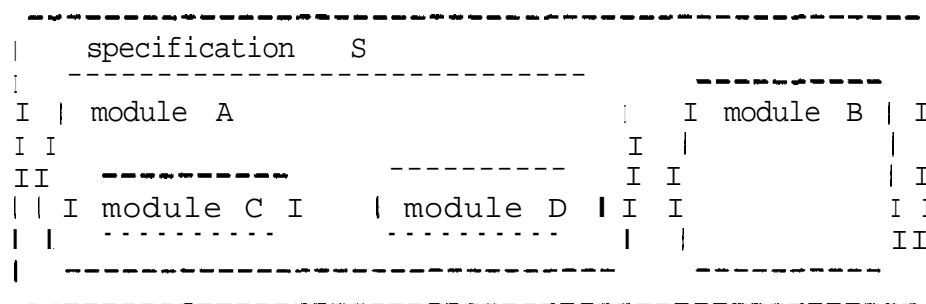
Dans le paragraphe suivant sont **présentés** les concepts fondamentaux d'Estelle. Une description plus **détaillée** d'Estelle est donnée dans [COU 86, BUD 87, ISO 891].

a) Les concepts de base

Ce paragraphe présente les concepts de base spécifiques du langage Estelle (expression des comportements **parallèles**). Il s'agit du sous-ensemble Estelle **implémenté** sous Echidna.

b) Notion de module

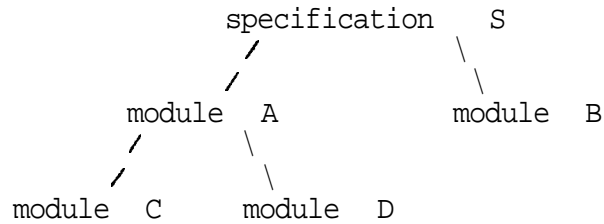
Une spécification Estelle est une collection de composants communicants **appelés** instances de modules.



Un module est défini dans la spécification par:

- une interface de communication **avec** son environnement,
- et, par son comportement à la définition d'un **entête(header)**, et d'un **corps(body)**.

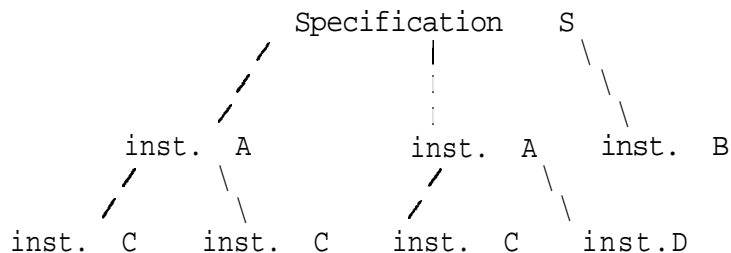
La définition d'un module peut comprendre des définitions d'autres modules (modules fils). Des définitions de modules en structure d'arbre **hiérarchique** sont ainsi obtenues, la **racine(module principale)** représentant toute la spécification. Il y a toujours une instance unique du module principal. On peut avoir **par** exemple, l'arbre des définitions de modules suivant:



L'entete d'un module definit l'interface du module avec son environnement **exterieur**. En d'autres termes, il **decrit** sa visibilite **exterieure**. Une instance de module est l'execution d'un module. Il y a **toujours une** instance unique de la specification qui est un module particulier. Une instance de module ne peut **être créée** que par **une** instance du module englobant. **Ceci** fait que l'arbre des definitions est un modele pour **l'arbre hierarchique** des instances de modules.

Exemple:

Par exemple, pour la specification **S**, un **arbre** d'execution possible est l'arbre suivant **avec** deux instances du module A et une instance du module B. La premiere instance de A initialise deux instances C, alors que la deuxieme lance une instance C et une instance D.



L'**entete** d'un module definit sa visibilite exterieure **sous** la **forme** d'un ensemble de points d'interaction, et de parametres.

```

Module type-entete1 (<parametres>);
  ip <ensemble de points d'interaction>
End;
  
```

Pour un **même** entete de module il peut y avoir differents comportements internes et **sous-structures**. **A un corps(body)**, **est toujours** associe un entete.

Exemple: Body type B for type-entete1;

Le **comportement** interne du module est donne par un modele de transition d'etats defini par:

- . un ensemble d'etats de **contrôle**,
- . **et**, un ensemble de transitions.

Syntaxiquement, le corps d'un module est composé de :

- une **partie** declarations
- une **partie** initialisations;
- et, une **partie** transitions.

```

Body a For b;
  <declarations>
  <initialisations>
  <transitions>
end;
```

1. **Partie** declarations:

La **partie** declarations comprend les declarations Pascal **usuelles**(**constantes**, **types**, variables, procédures et fonctions) et des declarations d'objets Estelle: canaux, **entêtes** et corps de modules, **états** de **contrôle** et variables modules.

La description d'un canal **contient** une enumeration des interactions qui peuvent **circuler** sur celui-ci dans un sens ou dans l'autre. Syntaxiquement, la definition d'un canal est donnée comme suit:

```

Channel nom-type-canal (id-role-ext1, id-role-ext2);
  By id-role-ext1:tdeclaration d'interactions>;
  By id-role-ext2:tdeclaration d'interactions>;
```

Exemple:

```

Channel canal(rôle1,rôle2);
  By rôle1 : M1;...Ml;
  By rôle2  : N1,...,Np;
  By rôle1, rôle2 : K1,...,Kq;
```

La declaration d'une interaction **comporte** le nom de l'interaction et éventuellement une liste de paramètres types.

Exemple: Requête(X:integer; Y:boolean);

Deux points d'interaction P1 et P2 peuvent être déclarés comme suit:

```

Ip      P1 : canal(rôle2);
        P2 : canal(rôle1);
```

P1 utilise le canal de communication de type canal et joue le rôle rôle1. P1 peut émettre tous les messages Mi et Kj, et recevoir les messages Ns. P2 peut émettre Ns, et recevoir Kj et Mi.

Exemple :

```
Channel canal1(rôle1, rôle2);
  By rôle1 : M1;... Mn;
  By rôle  : ;
```

Cet exemple montre la déclaration d'un canal unidirectionnel où seule l'extrémité qui jouera le rôle "rôle1" pourra émettre des interactions.

A chaque point d'interaction est associée une file d'attente non bornée de discipline FIFO.

Exemple de déclaration de corps de module:

```
Body A1 for A;

  Const tconstantes Pascal>;
  Type <declaration de types Pascal>;
  Var < " de variables Pascal>;

  Channel canal(rôle1, rôle2);
    By rôle1 : request(X:integer; Y:boolean);
    By rôle2 : response;

  State active, wait;
  stateset aw-[active, wait];

  Module B(moi: integer);
    I ip P1: canal(rôle1);
    | P2: canal(rôle2);
  End;

  Body B1 For B;
    I .
    | .
  End;

  Modvar int1, int2 : b;
end ;
```

Les variables de types modules introduites par le mot-clé Modvar servent de référence aux instances de module d'un certain type d'entête. Elles permettent de désigner les instances.

2. Partie initialisation:

La partie initialisation spécifie les valeurs de certaines variables avec lesquelles sera créée chaque instance de module. L'initialisation d'une variable module signifie la création d'un module fils. Elle établit aussi les liens de communication entre points d'interaction.

Exemple:

```

Initialize
  TO état-initial (état initial du module)
Begin
  <initialisations>
  <établissement des liens de communication>
End;
```

3. Partie transitions:

La **partie** transition est un ensemble de déclarations de transitions. La définition d'une transition est **constituée** d'une liste de gardes qui exprime la condition de tir de la transition, d'un bloc comportant des déclarations Pascal et d'une séquence d'instructions Pascal et Estelle **constituant** la **partie** action de la transition.

Syntaxe d'une transition

```

Trans
  <condition>
  {déclarations locales}
  <action>
```

La condition d'une transition est exprimée à l'aide des clauses suivantes:

```

.FROM      : définit l'état de départ;
.WHEN      : exprime une attente d'un message;
.PRIORITY  : assigne une priority à la transition, elle
             servira lors du choix des transitions
             tirables;
.PROVIDED  : exprime la satisfaction d'une expression
             booléenne;
```

c) Les instructions Estelle

1. Creation d'une instance:

L'initialisation d'une instance se fait grâce à l'instruction **INIT**:

```

INIT idf-module-variable WITH idf-body(<liste par
effectifs>)
```

La variable "module-var" dans l'instruction "**INIT**" fait référence à la nouvelle création **d'instance.de** module. Il peut y avoir plusieurs corps de modules associés à la **même entête** de module. L'instruction "**INIT**" **différencie** entre eux en spécifiant "idf-body".

2. Connexion de points d'interaction

La connexion des points d'interaction se fait grâce à l'instruction CONNECT.

Exemple: Connect **P1** TO P2

Cette demande de connexion ne peut être satisfaite que si:

- . les deux points d'interaction font référence au **même** canal et jouent des rôles opposés ,
- . aucun des deux points d'interaction référencés n'a déjà **été connecté**.

Une interaction ne peut être mise qu'à travers deux points d'interaction connectés.

3. Emission d'une interaction

La file d'attente **associée** à un point d'interaction est mise à jour lors des émissions et réceptions de messages. En Estelle il est possible d'émettre un message grâce à l'instruction OUTPUT qui **spécifie** le message attendu et le point d'interaction sur lequel il est attendu.

Exemple: OUTPUT **P1.m**

4. Réception d'une interaction

La clause WHEN d'une garde **permet** la réception d'une interaction en deux **étapes**:

- vérification de l'existence en **tête** de file du message attendu pour l'évaluation de la garde,
- **consommation** du message, si la transition est tirée.

5. des restrictions au langage Pascal sont admises, par exemple, il n'y a **pas** d'instructions de manipulation de fichiers.

Dans l'annexe est **présentée** une spécification Estelle de l'algorithme de Chang et Roberts [CHA 791].

1.2 Présentation D'Echidna

Echidna [JAR 89, JEZ 89] est **construit** autour d'un compilateur d'algorithmes **distribués décrits** en Estelle et fournissant du code pour machines multi-processeurs et d'un ensemble d'outils pour observer le **comportement** de l'algorithme distribué **sous** expérimentation.

L'exécution de la spécification Estelle va **donc**, consister à exécuter en parallèle les **tâches** sur chaque site. Pour mettre en œuvre ce **modèle** d'exécution, le compilateur traduit la description formelle Estelle en un ensemble de structures de **données**, exprimées en **langage C**, qui est **ensuite** compilé par le

compilateur C local, puis lie **avec** le **noyau** d'exécution distribuée (**NED**) avant **d'être** chargé et **exécuté** sur le **système cible**.

Le NED a la charge de connaître et de **gérer** la communication non locale **entre tâches** Estelle. Il est totalement écrit en C, et **découpé** en trois modules:

1. L'exécuteur:

Un **modèle exécutable** est obtenu après édition de lien **entre le** NED et le code **généré** par le compilateur. Il est chargé **sur tous** les sites. Parmi les **fonctions** de l'exécuteur on distingue:

- l'analyse des **commandes**,
- la création de la **tâche** englobante (**racine** de l'**arbre**) et son initialisation en **exécutant** sa "partie initialisation".
- le lancement du scheduler. Le scheduler **assure l'enchaînement** des **tâches initialisées** sur chaque noeud du **réseau**. Le code du scheduler est écrit en C.

2. Le runtime Estelle:

Il a pour **rôle** de reconstruire **une** machine virtuelle permettant l'exécution d'un programme Estelle à partir du **runtime** C qui **existe autour** de la machine. Il **doit** fournir les primitives de manipulation d'objets Estelle et Pascal qui n'existent pas en C, **un générateur** aléatoire (**pour le choix** non déterministe parmi les transitions **tirables**) et reconnaître la communication distante pour l'acheminer **par le réseau**.

3. L'interface **avec** le **système**:

Le **but** est de fournir de manière homogène l'abstraction "station d'expérimentation" quelque **soit** le **système** distribué **sous-jacent**. Il **fournit** des primitives permettant d'émettre des messages vers les autres sites et de tamponner les messages qui arrivent des **autres** sites. Il **fournit également** **une** horloge locale, **gère** les **entrées/sorties** et l'identificateur de site. **Seule cette** partie du NED est à adapter à chaque **nouveau système**.

2 Le Matériel

Le travail **présenté** dans cette thèse a été réalisé sur la **version V0.2(1989)** d'Echidna implémentée sur un micro-ordinateur du type PC **sous** le système d'exploitation MS/DOS.

III HISE EN **OEUVRE** DE LA **SIMULATION SOUS** Echidna

Ce paragraphe presente les conditions d'evaluation ainsi que les parametres pris en compte.

1 Le **Modèle** D'election

Deux **modèles** sont prevus:

- Election collective: l'ensemble des sites presents dans le systeme lancent, en parallele, une operation d'election.
- Election individuelle: un site particulier du systeme initialise l'operation d'election.

2 Les **Critères D'évaluation**

Deux criteres de performances ont **été mesurés**:

1. **Complexité** en nombre de **messages(trafic)**:

c'est le nombre de messages **échangés**, necessaires a la determination de l'identite de l'elu.

2. **Complexité** en temps **d'élection**:

C'est la **durée** necessaire a la determination de l'identite du site Clu.

Le **calcul** d'une election se termine dans un seul site qui se **chargera** alors d'en informer les autres sites. Les **complexités évaluées** dans ce document ne tiennent pas compte de cette phase de diffusion de l'identite de l'**élu**, **car**, elle ne fait pas toujours **partie** du **protocole** d'election.

3 Le **Modèle** Estelle **Implanté** **Sous** Echidna

Une description Estelle est un arbre dont les noeuds sont des **tâches** et la **racine** la specification englobante. Une **tâche** est **composée** d'une specification (module] et d'un corps (body). Un module est **créé** lors de sa specification dans le module immediatement englobant, et son corps lui est **attribué** lors de l'initialisation de la **tâche**. Les fils d'une **tâche** sont les instances des **modules(variables** de type module).

Sous Echidna, seul un sous-ensemble "statique" d'**Estelle** [JEZ 891, qui **permet** toujours la creation dynamique de l'arbre, **mais** pas sa modification a **été implanté**. En effet, les **tâches** non terminales sont passives: **donc** il ne **reste** sur chaque site **que** l'ensemble des **tâches** terminales **créées** et **placées** dynamiquement **sur ce** site. L'execution du programme Estelle va **donc** consister, a partir de la, a executer en parallele les **tâches** sur chaque site. Cette implementation d'**Estelle** est suffisante pour la specification et la mise en **oeuvre** des

5 Collecte Des Informations De Mesure

Dans le **système** Echidna l'instruction "Trace" [JEZ 89a, JEZ 89b] **est** la seule extension au langage Estelle (qui ne possède pas d'instructions d'Entrées/Sorties) qui a été faite. Sa syntaxe est celle d'un writeln Pascal, mais sa sémantique est propre au contexte de l'expérimentation. En effet, pour ne pas perturber l'expérimentation, une option prévue dans la commande de lancement d'exécution permet de forcer l'instruction "Trace" à effectuer les Entrées/Sorties en mémoire afin d'être transférées physiquement seulement après que l'expérience soit terminée.

L'instruction "Trace" a été utilisée pour recueillir les informations désirées (nombre de messages et temps d'élection) et qui sont enregistrées au cours de l'exécution de l'algorithme.

Il faut noter également que l'évaluation des algorithmes a été faite indépendamment de toute application.

Le schéma de la figure II.2 résume les différentes étapes par lesquelles transite un algorithme exécuté sous Echidna.

a) Evaluation du temps

Une des caractéristiques des systèmes distribués est l'absence d'un référentiel temporel global. Afin de pallier à cet inconvénient, de nombreux travaux ont été menés sur ce problème qui consiste à réaliser une horloge (logique) globale fournissant un temps abstrait global, servant de référentiel commun à tous les sites [LAM 78, JUD 87, JEZ 98a]. Cette perception du temps, est nécessaire essentiellement pour résoudre les problèmes de synchronisation (mécanismes d'allocation de ressources telle que l'exclusion mutuelle, la coordination de processus, les relations de précedence causale entre événements, etc).

Dans cette thèse, il est question de mesurer le temps d'exécution des algorithmes. Pour être pratique, il n'est pas nécessaire d'essayer de mesurer le temps d'exécution "reel" d'un algorithme. Ceci pour deux raisons:

- D'une part, deux exécutions différentes aboutissent à des résultats différents, car elles sont liées à l'environnement d'exécution et à la charge du système qui varie d'un instant à l'autre.

- D'autre part, les résultats obtenus vont dépendre de la mise en œuvre. C'est-à-dire, deux implementations du même algorithme donneraient des performances différentes, si elles sont réalisées sur des systèmes différents.

Dans ces conditions, la comparaison des performances d'algorithmes distincts (c'est là le but de notre travail) pourrait aboutir à des résultats spécifiques, voire non

significatifs. En effet, les conclusions qui peuvent en être tirées ne seront valables que dans un contexte équivalent à celui de l'évaluation. Avec des mises en œuvre différentes il est possible d'obtenir des conclusions complètement différentes et même contradictoires.

Donc, pour réaliser l'évaluation du temps, un modèle basé sur le fonctionnement du système ECHIDNA(scheduler) a été défini.

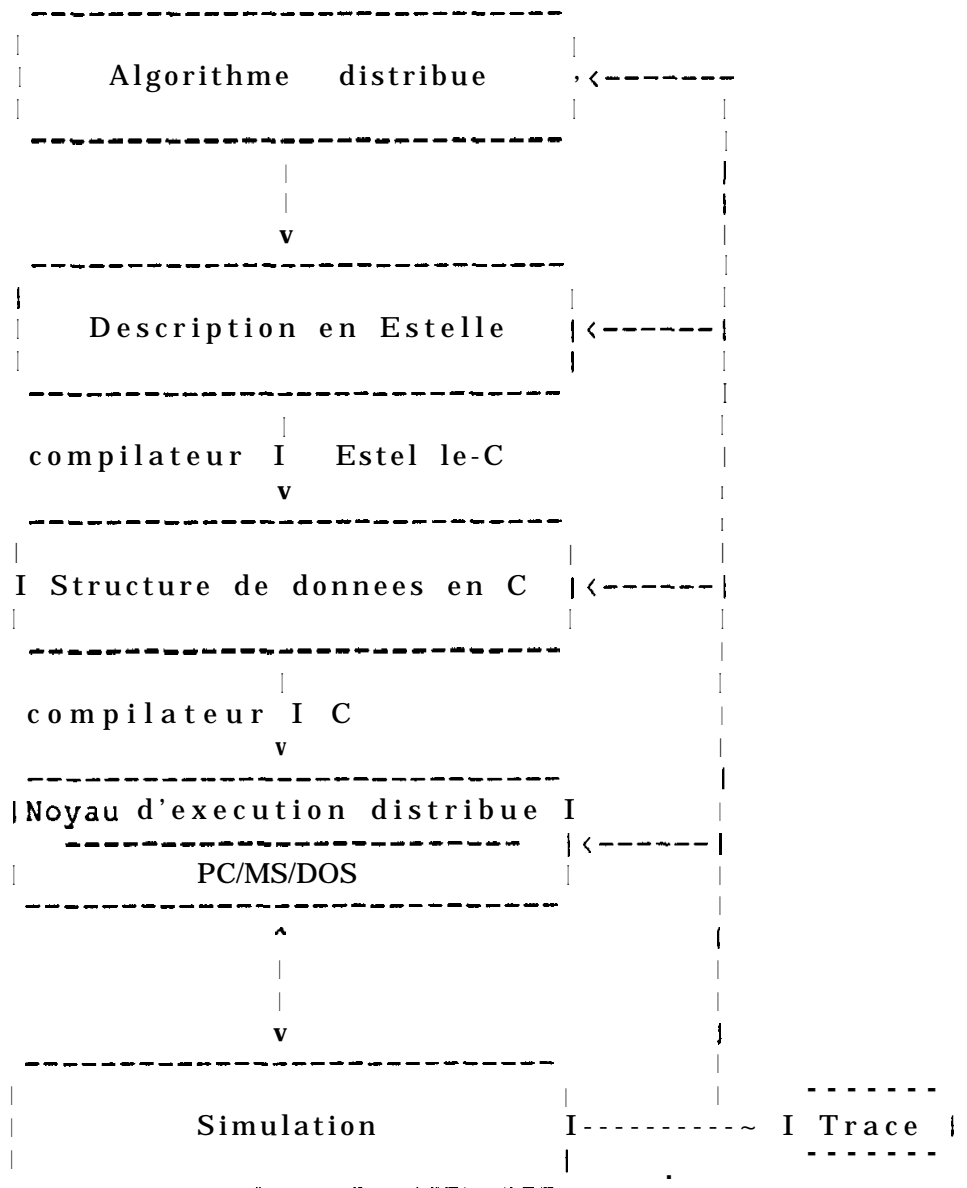


Figure 11.2: Simulation d'algorithme **distribué sous** Echidna

Examinons de près l'exécution de l'algorithme E qui est supposée s'effectuer en trois étapes :

```

étape1: { trans21, trans12, transmn }
étape2: { trans31, trans42, trans2n }
étape3: { trans72, transk6n }

```

1. Durant chacune des k étapes (k=1 à 31, tout message émis sera reçu durant la même étape. Il pourra (selon que la transition correspondante soit exécutée ou non) être traité dans l'étape suivante, car il s'agit d'un système centralisé.
2. On suppose que la durée d'exécution d'une étape représente une unité de temps.

L'évaluation du temps d'exécution d'un algorithme peut donc être représentée par le nombre d'étapes de calcul qui ont été effectuées. Pour l'algorithme E ci-dessus, le temps mesuré sera égal à 3 unités.

CHAPITRE 3

ETUDE DES **ALGORITHMES** D'ELECTIONI DEFINITIONS ET **TERMINOLOGIE**

1 Concept De Processus

Un processus est **considéré** comme **une** unite d'execution **élémentaire** d'un algorithme distribue ou parallele. Plusieurs processus peuvent operer simultanément. Dans un algorithme distribue, un processus est **considéré** comme un algorithme **séquentiel** (au point de vue de son execution) deterministe ou non, muni de primitives de communication **lui** permettant d'envoyer et de recevoir des messages, **afin de communiquer avec** les autres processus.

2 **Réseau** Anonyme

C'est un reseau dans lequel les processus ne sont pas dotes de numeros d'identifications **tous** differents, ou encore qui ont **tous** la **même** identite [LAV 87, MAT 891.

3 **Modélisation** Du **Réseau**

Le reseau de communication est modelise par un graphe $G=(V,E)$ connexe. Les noeuds (**V**) et les aretes (**E**) representent respectivement les processus et les canaux de communication.

4 Communication **Synchrone** Et Asynchrone:

Dans le mode de communication **synchrone**, une horloge **globale** est **connectée** a **tous** les sites du reseau. L'intervalle de temps entre deux pulsations **consécutives** de l'horloge est le "round". Au debut de chaque round, chaque site decide, en se basant sur son **état**, quels messages envoyer et vers **quelle** destination. Chaque site **recevant** des messages **durant** ce round, les utilise ainsi **que** son **contexte** local pour decider de l'action a entreprendre. Les sites reveilles spontanément commencent l'algorithme en entrant dans un **état** initial, et en attendant **le** debut du **prochain** round. Dans le mode asynchrone il n'y a pas d'horloge **globale**, et les **délais** d'acheminement des **messages** sont variables [AFE 85, FRE 85, ADA 881.

II LANGAGE DE DESCRIPTION ET NOTATIONS

Le texte de chaque algorithme est exprimé par la suite à l'aide d'un langage de type "événementiel". A chaque événement **ei** est associée une action **ai** décrivant le comportement d'un processus. Un tel événement est soit l'émission ou la réception d'un message **donné**, soit la validation d'une expression booléenne.

L'émission d'un message se fait par la primitive: **envoyer**<type de message>(message) à destination.

La notation suivante est utilisée:

$\forall i \in V : \text{envoyer}\langle \text{type de message} \rangle(\text{msg}) \text{ à } P_i;$
signifie envoyer à tous les P_i , tel que i soit dans l'ensemble V , le message **msg**.

La réception d'un message se fait par la primitive: recevoir <type message>(message) depuis origine.

La réception d'un message par P_i depuis P_j provoque une interruption de la façon suivante:
Lors de la réception <type message>(message) depuis j faire
-- le traitement du message --
fait.

Les notations suivantes sont utilisées tout le long de cette thèse:

- . **env-vg**: envoyer au voisin gauche,
- . **env-vd**: envoyer au voisin droit,
- . **faire-suivre**: env-vg le message reçu par la droite, et env-vd le message reçu par la gauche.
- . **Idi**: l'identité du processus P_i ,
- . **Pmax**: le processus d'identité maximale,
- . **n**: le nombre de processus,
- . **log x**: logarithme décimal de x ,
- . **lg x**: logarithme népérien de x ,
- . **Log x**: logarithme en base 2 de x ,
- . **Hn=0.693..Log n**.

III PRESENTATION D'UN ECHANTILLON D'ALGORITHMES

Le problème de l'élection a été étudié de façon intensive ces dernières années sous l'hypothèse d'identités distinctes des sites du réseau. Il s'agit alors d'un problème de calcul d'extremums, et divers algorithmes ont été publiés dans différents contextes et sous des hypothèses multiples:

- selon la topologie du réseau (anneau ou maillage complet le plus souvent);
- selon l'information sur la topologie globale du réseau qui est

mise a la disposition des sites (connaissance de la topologie particuliere du reseau, ou seulement de sa taille, connaissance des identites des voisins, ou connaissance seulement de l'orientation dans un reseau complet, **etc**);

- selon que le reseau **lui-même** est synchrone ou asynchrone, etc.

Dans ce chapitre, est **présenté** un echantillon d'algorithmes **publiés** pour les differentes topologies de reseau: anneau, a maillage complet, quelconque. Parmi ces algorithmes **certain**s ont été évalués sous Echidna.

1 ALGORITHMES D'ELECTION POUR TOPOLOGIE ANNEAU

1.1 Introduction

La plupart des algorithmes **étudiés** ne s'interessent qu'au cas où les processus du reseau sont dotes de numeros d'identification uniques. De plus, ils supposent **tous** que le reseau des processus ainsi que le systeme de communication sont fiables. Parmi ces algorithmes, il y en a **ceux** qui sont **basés sur une** structure en anneau unidirectionnel asynchrone [CHAN 79, DOL 82, PET 82], en anneau unidirectionnel synchrone [FRE 87b, GAF 89] **avec** une **complexité** en $O(n)$ messages, ou bidirectionnel asynchrone [HIR 80, LAV 89, BOD 86, FRA 82, LEE 87].

Deux algorithmes sont **donnés** dans [GOL 87] **avec** l'hypothese que le reseau bidirectionnel de connexion des processus pourrait être defaillant. Dans le premier algorithme le nombre n de **procesus** n'est pas connu, cependant, chaque processus P_i connait les identites ID_{i-1} et ID_{i+1} de ses deux voisins de gauche et de droite. Le second algorithme suppose la connaissance de n . La **complexité** attendue pour ces deux algorithmes est $O(n \log n)$ messages. En effet, cette borne est la meilleure borne pour **tous** les algorithmes proposes dans le cas d'un reseau **asynchrone** (c'est la **complexité minimale** possible dans le cas d'un tel reseau [PAC 84]).

Parmi les algorithmes **présentés** dans ce chapitre, quatre algorithmes sont **évalués sous** Echidna. L'algorithme de Chang et Roberts [CHA 79], c'est l'un des premiers algorithmes proposes pour resoudre le **problème d'élection**. L'algorithme d'Hirschberg et Sinclair [HIR 80], base **sur un** anneau bidirectionnel. Le troisieme algorithme est celui propose, simultanement, par Dolev [DOL 82] et Peterson [PET 82] sur **un** anneau unidirectionnel necessitant $O(n \log n)$ messages. **Ceux-ci** **présentent** une refutation de la conjecture de Hirschberg et Sinclair qui avaient **énoncé que seul un** anneau bidirectionnel pouvait permettre la realisation d'une **complexité** inferieure a $O(n^2)$ messages. L'algorithme de KORACH [KOR 81], par **contre**, est une variante probabiliste de l'algorithme de [CHA 79] conçue pour un anneau bidirectionnel.

1.2 Hypothèses

Les hypothèses faites sur le réseau de processus et le système de communication:

- Les processus ont des fonctionnements asynchrones et ne communiquent que par messages.
- Ils ont des identités différentes prises dans un ensemble totalement ordonné.
- Aucun processus ne tombe en panne pendant l'élection.
- Aucun message n'est perdu, ni altéré.
- Les délais de transmission sont quelconques, mais finis.
- Enfin, les processus n'ont pas d'informations sur le nombre n de processus, et ne connaissent que leurs identités.

Dans le cas où une ou plusieurs de ces hypothèses ne sont pas respectées, cela sera précisé.

Tous les processus exécutent le même texte, sauf que chaque texte fait référence à l'identité du processus qui l'exécute. Ces identités sont toutes distinctes. Les processus ont des comportements pouvant être distincts, l'algorithme a donc une symétrie de texte.

Remarque:

Les algorithmes sont présentés en commençant d'abord par les algorithmes basés sur une structure unidirectionnelle puis bidirectionnelle et en tenant compte de l'ordre chronologique de leur publication.

1.3 Algorithme De Chanq Et Roberts [CHA 791

1.3.1 Principe

L'algorithme est basé sur l'hypothèse que les processus sont structurés en anneau unidirectionnel de façon quelconque. Chaque processus P_i , qui connaît son propre numéro i , le transmet à son voisin de gauche P_j . Celui-ci, à la réception d'un tel message compare son propre numéro j au numéro reçu i . Si ce dernier est plus grand, il le transmet à son voisin de gauche, dans le cas contraire c'est son propre numéro qui est transmis: c'est le principe de l'extinction sélective de messages.

Initialement un ou plusieurs processus commencent une élection en envoyant un message: un tel processus se marque comme participant à l'élection. L'arrivée d'un message à un processus non marqué, provoque sa participation à l'élection. Le processus marqué, qui reçoit un message portant son propre

numero conclut alors que son numero a fait un tour complet sur l'anneau et est donc le plus grand; il est alors élu et diffuse son identite aux autres si cela est necessaire au fonctionnement de l'algorithme qui a déclenché l'election.

1.3.2 Description formelle

a) Contexte local d'un processus P_i

| | |
|-----------|---|
| Constante | mon-numero: valant i ; |
| Variable | participant: bouleen initialise a faux; coordinateur: entier; (utile si l'algorithme qui a declenche l'election nécessite la connaissance du vainqueur final) |

b) Types de messages

- Election: indique que la valeur qu'il transporte est candidat a l'election.

- Elu: necessaire pour diffuser l'identite du gagnant.

c) Algorithme exécuté D a r un processus P_i

```

-----
| Lors de
|   decision de provoquer une election
|   faire participant:= vrai;
|   |   env-vg(élection.mon-numéro)
|   fait;
|   reception de (élection.i)
|   faire selon le cas
|   |   i>mon-numéro: env-vg(election.i);
|   |   |   participant:= vrai;
|   |   itmon-numero et non-participant:
|   |   |   env-vg(élection.mon-numéro);
|   |   |   participant:= vrai;
|   |   i=mon-numero: env-vg(élu.i);
|   fait;
|   reception de (élu,i)
|   faire coordinateur:- i;
|   |   participant:= faux;
|   fait; si i<>mon-numéro alors env-vg(C1u.i) fsi;
-----

```

. complexité thdorique

En messages: l'algorithme de Chang et Roberts utilise, pour l'étape d'élection, un seul type de message. Dans le cas favorable (les processus sont places sur l'anneau dans l'ordre croissant de leurs numéros, et le processus P_{max} lance l'élection) le nombre de messages est $O(n)$. Dans le cas le plus defavorable (ordre decroissant et tous les processus lancent l'élection) $n(n-1)/2$ soit $O(n^2)$ messages sont necessaires. Dans les autres cas le nombre moyen est de $O(n \lg n)$.

En temps: le temps requis pour l'élection est de $O(n)$.

1.4 Algorithme De Dolev, Klawe Et Rodeh [DOL 821

Ce même algorithme fût publié simultanément par Dolev [DOL 821 et Peterson [PET 821.

1.4.1 Principe

- Les communications se font sur un anneau unidirectionnel de droite à gauche,
- les messages sont reçus dans leur ordre d'émission.

L'algorithme procède par phases, à chaque phase l'anneau unidirectionnel des processus est réduit au moins de moitié (un processus sur deux devient passif) jusqu'à ce qu'il ne reste qu'un processus actif qui aurait déterminé alors l'identité du processus élu.

Le principe d'une phase est le suivant:

Un processus P_j attend d'avoir les numéros associés à ses deux voisins de droite, soient P_k et P_l . Puis il se place du point de vue de P_k (qui est entre P_j et P_l) il teste si celui-ci est maximum local auquel cas il transmettra le numéro, sinon P_j deviendra passif et son rôle se limitera alors à transmettre les messages reçus (lorsque P_j reste actif son voisin de droite devient nécessairement passif).

1.4.2 Description formelle

a) Contexte local d'un processus

Constante mon-numero valant j ;
 Variable Cstat: (actif, passif) initialisé à passif;
 max : entier initialisé à mon-numero; {numéro du processus du point de vue duquel le processus P_j se place}
 voisindroit: entier; {identité du premier voisin actif à droite de P_j }

b) Types de messages

Deux types de messages sont utilisés pour véhiculer les identités des deux processus voisins de droite actifs.

(b , numéro) où b vaut 1 ou 2 selon qu'il s'agisse du premier ou du deuxième voisin de droite.

c) L'algorithme exécuté par P_i

```

-----
| Lors de
|   decision de provoquer une election
|   faire   env-vg(1,mon-numéro);
|           état:= actif;
|   fait;
|
|   reception de (1, numero)
|   faire si état=actif alors
|       |           si numéro<>max alors
|       |           |   env-vg(2,numéro);
|       |           |   voisindroit:=  numero;
|       |           |   sinon (max est le plus grand numero
|       |           |   |   il est diffuse aux processus)
|       |           |
|       |           fsi
|       |   sinon env-vg(1, numero);
|   fsi;
|   fait;
|
|   reception de (2, numero)
|   faire si état=actif alors
|       |           si voisindroit>numéro et voisindroit>max
|       |           |   alors max:= voisindroit;
|       |           |   |   env-vg (1, voisindroit)
|       |           |   |   sinon état:= passif;
|       |           |
|       |           fsi
|       |   sinon env-vg(2, numero);
|   fsi;
|   fait;
-----

```

. **Complexité théorique**

A chaque phase de l'algorithme, l'anneau est réduit de moitié et il se termine quand il ne restera qu'un seul processus actif dans l'anneau. Ce qui nécessite $(\text{Log}n+1)$ phases. $2n$ messages sont @changes (2 par processus (1, numéro) et (2, numéro)) durant les $\text{Log} n$ phases, et n messages dans la dernière phase (un seul processus est actif et il n'y a donc pas un deuxième voisin). La complexité de l'algorithme est, par conséquent, de $2n\text{Log}n + n$ messages.

1.5 Algorithme De Hirshberg Et Sinclair [HIR 801

1.5.1 Principe

Les processus sont places de **façon** quelconque sur un anneau bidirectionnel.

Le **principe** de l'algorithme est celui d'élection primaire sur des ensembles de processus de plus en plus grands **jusqu'à** contenir **tous** les processus places sur l'anneau. Ainsi, lorsque P_j lance une election, il se declare **candidat** et cherche à savoir si son numero est **supérieur** a ceux de ses deux voisins **immédiats** P_k et P_l (dont il ne connaît pas les **numéros**) en leur diffusant son propre numero. Si l'un des deux autres processus P_k ou P_l est vainqueur (k ou l supdrieur a j), il devient **candidat** s'il ne l'est pas déjà. Si les deux autres processus se **déclarent** battus, P_i **réitère** sa candidature sur un ensemble plus grand. Les messages de candidature vont progresser a chaque iteration sur des chemins dont la longueur **croit** selon une puissance de 2.

Le role d'un processus battu a une iteration **donnée** se limite a faire suivre les messages d'un **côté** vers l'autre sur l'anneau.

1.5.2 Description formelle

a) Contexte local d'un processus P_i

```

Constante  mon-numero  valant  i;
Type       position=(non-concerné, candidat, battu, élu);
Variable   état: position initialisé à mon-numero;
           lgmax:  entier;
           vainqueur:  entier;
           nbrep:  0..2 initialise a 0;
           repok:  booléen initialisd à vrai;

```

b) Types de messages

Deux types de messages sont utilises: candidature et réponse.

- 1) (candidature, numero, lg, lgmax);
 - numero indique le numero associe au processus qui a lance le message de candidature,
 - lgmax indique la longueur du chemin que le message doit à priori parcourir sur l'anneau,
 - lg indique la longueur déjà parcourue.
- 2) (réponse, bool, numero);
 - bool a la valeur vraie si la **réponse** est favorable, faux sinon,
 - numero **désigne** l'identite du processus **déstinataire** (c'est une **façon** de l'adresser sur l'anneau puisque ce numero est unique).

Or un processus lançant sa candidature sur des chemins de longueur 2^L provoque le transfert de quatre messages (2 dans chaque direction). Par conséquent, le nombre maximum de messages échangés est borne par $8n+8n \log n$ ce qui est $O(n \log n)$.

1.6 Algorithme-P De Korach, Rotem Et Santoro [KOR 81]

1.6.1 Principe

L'anneau est bidirectionnel, l'ordre des messages est conserve et les delais de transmission des messages sont tous égaux.

L'algorithme execute par un processus P_i procède en étapes. Une variable locale Max garde l'identité maximale reçue par P_i . Les étapes d'élection se déroulent comme suit: A la réception d'un message contenant l'identité ID_j du processus P_j venant de la droite ou de la gauche de P_i : si ID_j est supérieur a Max, le message est retransmis et continue son parcours. Si deux messages arrivent simultanément des deux directions, celui qui véhicule la valeur la plus petite est ignoré. Le processus P_i qui reçoit une valeur ID_i égale a Max est élu. Il diffuse un message de fin d'élection.

L'initialisation consiste a choisir, avec une probabilité égale a $1/2$, la direction droite ou gauche vers laquelle le message contenant ID_i est émis par le processus P_i .

1.6.2 Description formelle

a) Contexte d'un processus P_i ..

Constante ID_i : entier valant i ;
 Variable max: entier initialisé a i ; {identité de l'élu}
 participant, élu: booléens initialisés a faux.

b) Types de messages

Deux types de messages sont utilisés: (numéro) pour la phases election, et (stop) pour signaler fin election.

c) L'algorithme exécuté var un processus P_i

```

-----
| Lors de
|   decision de provoquer une election
|   faire Max:= IDi; élu:= faux; participant:= vrai;
|   | choisir une direction d appartenant a {droite, gauche};
|   I emettre(IDi) sur d;
|   fait;
|
|   reception de (numéro)
|   faire Cas participant: si numéro > Max alors
|   |   max:= numero;
|   |   fairesuivre(numero)
|   |   sinon si numéro = max alors
|   |   |   env-vg(stop)
|   |   fsi;
|   fsi;
|
|   non-participant: provoquer election avant de I
|   |   traiter le message;
|   Fcas;
|   fait;
|
|   reception de (stop)
|   faire
|   | I si élu=vrai alors fairesuivre(stop) fsi;
|   fait;
-----

```

. Complexité théorique

Le nombre moyen de messages prévu par les auteurs [KOR 81] est $3n \ln n/4$ ($\ln n = 0.693 \dots \log n$) et $n/2$ dans le pire des cas, mais dans [Lav 89C] la valeur asymptotique calculée est $\sqrt{2} n \ln n/2 + O(n)$. Cette évaluation a été faite en supposant un fonctionnement synchrone des processus.

1.7 Algorithme De **Fredreickson** Et Lynch Pour Un **Réseau Synchrone** [FRE 87b]

. Principe

L'algorithme se déroule sur un anneau unidirectionnel avec communication synchronisée par des tops horloge.

Il est initialisé par des processus qui décident indépendamment les uns des autres de se réveiller. Le réveil ne se fait pas nécessairement en même temps pour tous les processus.

Une fois actif, un processus (appelé maintenant participant) génère un message, qui se déplacera sur l'anneau, transportant son identité IDi.

L'algorithme est basé sur deux idées. La **première-idée** consiste à transmettre les messages à des vitesses différentes. Le message comportant l'identité ID_i circule à la vitesse d'un transfert tous les 2^i tops c'est-à-dire que le message est retenu pour une durée de 2^i tops au niveau de chaque processus, qu'il visite, avant d'être retransmis.

Un message qui est rattrapé par un autre message plus rapide est éliminé de la course. De même, lorsqu'un message transportant un ID_i , atteint un processus j , il est éliminé si $j < i$ et si, en plus j aurait déjà généré son propre message. L'idée de faire transférer les messages à des vitesses variables est également exploitée dans [GAF 85a] qui propose un algorithme semblable à celui-ci.

La **deuxième-idée** est d'avoir une phase préliminaire avant la phase des vitesses variables. Dans cette phase, tous les messages circulent à la même vitesse: un transfert par top. Quand un processus décide de participer à l'élection, il envoie un message vers son voisin. Ce dernier est transmis sur l'anneau jusqu'à la rencontre du participant suivant. Dès cet instant, le message entre en deuxième phase et l'algorithme se poursuit comme expliqué précédemment.

. Complexité théorique

En messages elle est de $O(4n)$. En temps elle est de $O(n2^i)$ rounds, où i est l'identité du processus Clu et c'est la plus petite identité du réseau.

1.8 Algorithme De Spirakis, Tampakas Et Athanasios [SPI 891

Cet algorithme suppose connue au moins une borne supérieure sur la valeur de n . Dans [SPI 891 est considéré le cas d'un anneau de processus dont les identités ne sont pas toutes distinctes: le réseau est anonyme.

. Principe

L'algorithme est basé sur un anneau unidirectionnel et consiste en deux étapes majeures:

(a) Etape de création de candidats: dans cette, étape les processus, choisissent par tirage aléatoire une identité dans l'intervalle $[1, \dots, n]$. À la fin de cette phase, seul un ensemble de processus deviennent candidats. Les autres deviennent des relais dont le rôle se limitera, par la suite, à retransmettre les messages reçus. Dans cette étape, il existe une probabilité de blocage de l'algorithme lorsqu'aucun processus ne devient candidat.

(b) Etape d'élection: Les candidats échangent des messages comportant un numéro de phase, et leur identité ID . Au début de chaque phase, après le choix de ID , et l'augmentation du numéro de la phase, un message contenant ces informations est émis. Si

le message réussit à parcourir **tout** le réseau sans rencontrer un autre **candidat** dont la phase **et/ou** l'identité sont supérieures aux siennes, alors son initiateur devient **élu**. **Sinon** une autre ID est choisie et le message est **réémis après** incrémentation de la phase.

Un processus **candidat** i recevant un message (phase j , ID j) supérieur à (phase i , ID i) n'est plus candidat. IL **réemet** le message (phase j , ID j).

. Complexité théorique

Le **trafic engendré** est en moyenne $O(n)$ messages.

1.9 Algorithme De Franklin [FEW 821

. Principe

La communication est **synchrone**. Cet algorithme est basé sur un **principe** analogue à celui de Dolev [DOL 821, mais ici l'**anneau** est bidirectionnel. Une étape de l'algorithme se déroule comme suit:

Un processus **actif** teste **lui-même** si son numéro est plus grand **que ceux** de ses deux voisins de gauche et de droite. Si ce n'est pas le cas, il devient **passif**.

À chaque étape, l'anneau des processus **actif** voit sa taille **divisée au moins par deux**: les processus **passifs relayant** simplement les messages **reçus** d'un **côté** vers l'autre.

L'algorithme se termine lorsqu'un processus **reçoit un message** de lui-même.

. Complexité théorique

Comme à chaque étape, il y a $2n$ messages transmis et il y a $(\text{Log}n+1)$ étapes, au total $2n(\text{Log}n+1)$ messages sont **échangés** plus n messages pour informer l'ensemble des processus de l'identité du vainqueur.

Si l'**unité** de temps **d'exécution** est supposée être le délai de transfert d'un message d'un processus vers ses voisins, alors dans le pire des cas, le temps **d'exécution** est $O(n\text{Log}n)$.

1.10 Algorithme-D De Bodlaender Et Van Leeuwen [BOD 861

. Principe

C'est une version **déterministe** de l'algorithme-P décrit ci-dessus. L'**idée** est de remplacer la phase d'initialisation qui choisie une direction de **façon** aléatoire, 'en une phase complètement **déterministe**. Chaque processus P_i envoie le message contenant ID i dans la direction où se trouve le voisin de plus petite **identité** (au début de l'algorithme, un processus P_i **émet un message** contenant ID i vers ces deux voisins et **reçoit** à son tour ID $i-1$ et ID $i+1$). Ceci va **éliminer** les plus petits voisins. **Dès** la fin de la phase d'initialisation, seuls les processus qui sont maximums locaux survivront.

. Complexité théorique

$\sqrt{2}$ $nHn/2$ messages en moyenne et $n^2/4$ dans le pire des cas.

1.12 Algorithme De Leeuwen Et Tan [LEE 871

. Principe

A chaque direction sur l'anneau est associée une file d'attente de messages et l'ordre des messages est préservé.

L'algorithme **opère** de la manière suivante:
 A tout moment chaque processus tient à jour un registre ID contenant la plus grande identité de processus, et un registre DIR **booléen** qui pointe vers la "direction" de l'anneau où se trouve un processus de plus petite identité. Le "grand" **candidat** se trouvera sur un **côté** du processus tandis que DIR **pointera** dans la direction **opposée**. Les messages contenant le "grand" **candidat** sont **envoyés** dans la direction où se situe le "petit" **candidat** encore **actif**. L'idée est **donc** de se débarrasser du **candidat** dont l'identité est la plus petite. C'est la **même idée** qui est utilisée dans l'algorithme-D dans la phase d'initialisation.

. Complexité théorique

La complexité en messages est de $1,44...n \text{Log} n + O(n)$.

2 ALGORITHMES POUR TOPOLOGIE COMPLETE

2.1 Introduction

Dans un réseau complet chaque nœud est **connecté** à **tous** les autres nœuds ($n-1$ liaisons sont **reliées** à chacun des **nœuds**). Ce type de réseau est **modélisé** par un graphe $G=(V, E)$ complet.

L'algorithme d'élection le plus évident dans ce cas est le suivant: Chaque nœud envoie des messages contenant son identité à **tous** ses voisins. Ces derniers choisissent comme **élu** celui dont l'identité est la plus **élevée**. La complexité de cet **algorithme** est $O(n^2)$ messages dans le pire des cas.

Pour les algorithmes **présentés** dans cette section, la complexité en messages est **améliorée**, Cependant, la complexité temporelle augmente vu que la **fréquence** de transmission des messages est **réduite**.

Les algorithmes rencontrés, **conçus** pour la topologie complète, sont basés sur des réseaux à communication bidirectionnelle. Comme pour le cas de l'anneau, rares sont **ceux** qui résistent à **certains** types de pannes [GAR 821. Les algorithmes proposés dans [AFE 85] **présentent** une complexité en $O(n \text{Log} n)$ messages et la meilleure borne temporelle en $O(n)$ unités de temps. L'idée de base exploitée dans ces algorithmes

[AFE 85, KOR 84] est la notion de **quantité-de-travail** qui est utilisée pour choisir le **candidat** à **éliminer** parmi les processus concurrents à l'élection. L'algorithme de Loui [LOU 86] par **contre** utilise la notion du sens d'orientation pour réduire la complexité à $O(n)$ messages. Cette notion a **été** reprise dans [ATT 89, MAS 89] pour, respectivement, **des** réseaux en anneau auxquels des liaisons supplémentaires ont **été rajoutées**, et des réseaux où peuvent apparaître des pannes de processeurs.

2.2 Hypotheses

Les hypothèses suivantes sont communes à **tous** les algorithmes, d'autres hypothèses seront **précisées** en cas de besoin.

- A chaque noeud est associé un numéro identificateur unique (ID) **puisé** dans un ensemble totalement ordonné.

- Au niveau de chaque processus, chaque port d'entrée et chaque port de sortie a un identificateur unique connu par le processus.

- Initialement (avant le début de l'algorithme, en dehors des identificateurs de ses ports, un processus n'a aucune information sur le **réseau**. En particulier, il ignore les identificateurs des processus auxquels **il est connecté**. Les processus n'ont pas de connaissance quant à la topologie du **réseau** ou au nombre total de processus.

- Les liaisons de communication sont fiables. Les messages **émis** sur ces lignes atteignent leur destination en un **délai** imprévisible, **mais fini**.

- L'ordre des messages est préservé.

- Les liaisons sont bidirectionnelles.

- Le mode de communication est asynchrone.

- Les algorithmes **présentés** sont caractérisés par une symétrie de texte.

Les algorithmes seront **présentés**, comme pour la topologie en anneau, dans l'ordre chronologique de leur publication.

2.3 Algorithmes De **Garcia-Molina** [GAR 82]

L'étude faite dans [GAR 82] sur la **définition de** l'élection suivant le type de pannes peut **être considérée** comme une première classification de ce type d'algorithmes. Dans cette étude, l'auteur a **défini** l'élection suivant deux environnements d'exécution, un environnement **avec** des pannes de communication (formation de **sous réseaux**) et des **arrêts** temporaires de sites et un environnement exempt de ces deux types de pannes. Il propose un algorithme pour chaque **type** d'environnement,

respectivement, l'algorithme "brutal" et l'algorithme "invitation".

2.3.1 Principe de l'algorithme brutal

Dans le premier algorithme, dit "brutal" (bully algorithm). Les hypothèses suivantes sont **considérées**:

- Le réseau de communication est fiable.
- Les messages sont **traités** par les processus **dès** leur **arrivée** **avec** une **durée nulle**.
- Une **hypothèse implicite** de connaissance **globale** est également faite, les processus connaissent les **identités** des autres processus **ainsi** que la structure du réseau.

L'algorithme consiste, alors pour un processus P_i , à envoyer un message de candidature à **tous** les processus P_j tels que $i < j < n$. Il se déclare battu et il attend une candidature pendant un temps T au **delà** duquel il relancera sa candidature. **Si par contre** au bout de la première période T , il ne **reçoit** aucune **réponse**, il conclut que ces processus sont en panne. Il demande aux processus P_k tels que $k < i$ d'être à l'écoute car il ne connaît pas leur **état** et il les informe **ensuite** de son élection en leur **envoyant** un message Clu .

Le nombre maximum de messages **générés** par cet algorithme est en $O(n^2)$.

2.3.2 Principe de l'algorithme (invitation)

Dans le **deuxième** algorithme, il est supposé que le **réseau** peut se partitionner en plusieurs **sous-réseaux** et qu'un site peut **arrêter** son traitement et le reprendre **ensuite** comme si rien ne s'est passé. On peut **donc** avoir à un instant **donné** plusieurs groupes de processus **indépendants** **avec** chacun un coordinateur. Le **principe** de l'algorithme consiste alors pour un coordinateur d'un groupe **donné** à **inviter** **dès** que possible (**dès** la réapparition d'un site **isolé** et/ou le rétablissement d'une **ligne** en panne) les autres coordinateurs à rejoindre son groupe. À la réception d'une invitation, un processus peut **accepter** ou refuser une **telle** invitation. Au terme de l'algorithme, si le réseau est redevenu connexe, un seul coordinateur subsistera.

2.4 Algorithme De Korach, Moran Et Zaks [KOR 84]

. Principe

Initialement **tous** les processus se considèrent rois. Au terme de l'algorithme, seul le processus **élu** restera dans cet **état**, **tous** les autres processus deviendront "citoyens". Durant l'exécution de l'algorithme, chaque roi est **racine** d'un arbre représentant son royaume. Les autres nœuds de cet arbre sont citoyens du royaume. Un roi tente d'agrandir son royaume en **envoyant** des messages vers d'autres **rois**, leur demandant de joindre leurs royaumes au sien.

Six types de messages sont **utilisés**:

réveil: un noeud qui **reçoit** ce message se reveille et démarre l'algorithme. Au plus un message **réveil** atteint un processus quelconque.

demande: ce message est **envoyé** par un roi i , a travers un canal non utilise, au roi j pour l'inviter a rejoindre son royaume.

acceptation: ce message est **émis** par le roi i , comme reponse au roi j pour l'informer de l'acceptation de son invitation.

mis-a-jour: un roi i **recevant** un message d'acceptation emet le message mis-a-jour, pour aviser ses nouveaux citoyens.

citoyen: ce message est **émis** par un citoyen en reponse a une demande de son roi.

. Complexité théorique

L'algorithme **génère** au maximum $5n \log k + O(n)$ messages (k étant le nombre de processus ayant **reçu** le message **réveil**).

2.5 Algorithmes de Afek et Gafni [AFE 85]

Trois algorithmes asynchrones A, B et C ont **été** proposes dans [AFE 85]. Chacun d'eux est **conçu** dans le but de résoudre les problèmes poses **par** le précédent, de telle **façon** que l'algorithme C atteigne la **complexité** la plus intéressante par rapport à B et à A. Les trois algorithmes sont **tous** bases sur le **même** principe.

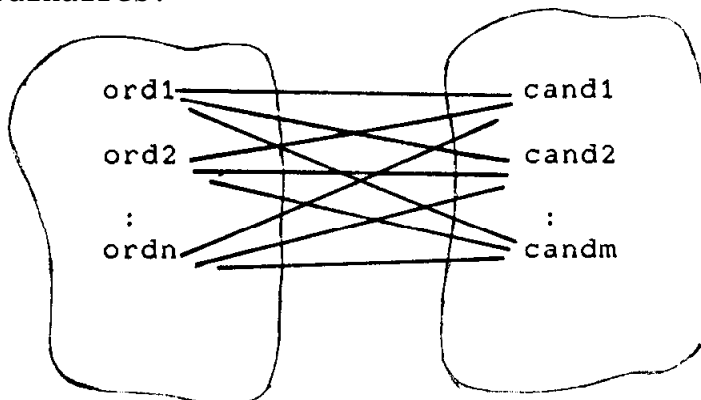
2.5.1 Principe

L'algorithme est initialisé par un ensemble quelconque de noeuds **candidats à l'élection**. Chaque **candidat** tente de construire une arborescence, dont il est la **racine**. Les autres noeuds seront **tous** ses **fil(s)** (directs car la topologie est complète). Pour **cela** il envoie des messages sur toutes les lignes qui lui sont **connectées**. Un seul **candidat** sera **élu**, tous les autres seront éliminés.

Un noeud se reveillant spontanément **génère** deux processus:

- un processus candidat,
- un processus ordinaire.

Tous les autres noeuds **génèrent** seulement un processus ordinaire. Chaque processus **candidat** est **relié** à tous les processus ordinaires.



Les candidats utilisent une variable appelée 'niveau' qui reflète le nombre de noeuds battus. Ces derniers tiennent à jour des variables 'niveau', qui gardent la trace du candidat ayant le niveau le plus élevé, qui ils ont été visités.

Un candidat qui atteint un noeud dont le niveau est plus élevé que le sien est éliminé. Au contraire, si le niveau du candidat est plus élevé ou égal à celui du noeud hôte, le niveau de ce dernier est remplacé par le niveau du candidat qui devient alors son père(racine), et qui tente d'éliminer l'ancien père.

Les trois algorithmes différent en deux points:
 (1) la façon dont la variable niveau est mise à jour,
 (2) la règle d'élimination.

Dans l'algorithme A, le niveau est égal à u nombre de processus ordinaires battus. Dans l'algorithme B, il est égal au nombre de candidats concurrents éliminés, et dans l'algorithme C, il représente la combinaison des deux précédents.

Deux variables de type pointeur, pere et pere-potentiel, sont tenues à jour par chaque processus ordinaire. Elles représentent, respectivement, le père le plus récent et le candidat qui, actuellement, tente de joindre le processus en question à son arborescence afin de devenir son nouveau père.

2.5.2 Algorithme A

Dans cet algorithme, la variable **niveau** d'un candidat est le nombre de processus ordinaires dont il est la racine. Elle est incrementée de un, lorsque le candidat a réussi à agrandir son arborescence d'un noeud. La règle utilisée est la suivante:

Pour que P puisse capturer un noeud V dont le propriétaire courant est le candidat Q:

```

-----
| (1) (niveau, ID) de P doit être plus grand (lexicographiqu-|
| ement) que (niveau, ID) du père courant de V. |
| (2) l'ancien père de V doit être éliminé. |
| |
| Autrement dit: |
| Si (niveau(P), ID(P)) < (niveau(V), ID(Q)), P est éliminé. |
| Si (niveau(P), ID(P)) > (niveau(V), ID(Q)) alors |
| | (1) V reçoit niveau de P. |
| | (2) P est envoyé vers Q. |
| Arrivé à Q: |
| Si (niveau(P), ID(P)) < (niveau(Q), ID(Q)), P est éliminé. |
| Si Q a été éliminé par un autre candidat alors capturer V. |
| Si (niveau(P), ID(P)) > (niveau(Q), ID(Q)) alors |
| | (1) Q est éliminé. |
| | (2) P capture V. |
-----

```

. Description formelle

a) Gontexte-local d'un-processus-Pi

Constante ID: entier valant i; (* identite du noeud i *)

Variable niveau: entier initialise a 0;
 ID-propretaire: entier initailisé à 0 ;
 non-traverse: ensemble des liaisons adjacentes;
 père, pere-potentiel: pointeurs initialises à nil.

b) Types de messages

Un seul type de messages est utilisé: [niveau, ID].

c1 L'algorithmme-exécuté-par-un-noeudi

```

-----
| Candidat:
| Lors de
| décision de provoquer une election
| faire
| l:= elt(non-traversé); (* un element dans non-traverse *)
| I envoyer(niveau,ID) à l; (* l: lien de communication *)
| fait;
| reception de (niveau', ID') de puio l'
| faire si (ID=ID') alors (* capture réussie *)
|     niveau:= niveau+l;
|     non-traverse:- non-traverse-l;
|     si non-traverse <> Ø alors
|         l:= elt(non-traversé);
|         envoyer(niveau,ID) sur l ;
|         sinon annoncer"élu, f i n  algorithmme";stop
|     fsi
| sinon si ((niveau',ID')<(niveau,ID)) alors
|     ignorer ce message;
|     sinon (1) envoyer(niveau',ID') sur l ' ;
|         (2) ignorer tous les messages
| fait; fsi; fsi ; antérieurs
-----

```

```

-----
| I Ordinaire:
| Lors de reception de (ID',niveau') sur l'
| Faire cas (niveau',ID')<(niveau,ID-propretaire): ignorer
| le message;
| I (niveau',ID')>(niveau, ID-propretaire):
| I père-potentiel<---ID';
| I niveau<--- niveau';
| I si père=id then père<--- pere-potentiel fsi;
| I envoyer (ID',niveau') vers père;
| I (niveau',ID')=(niveau, ID-propretaire):
| I père<--- pere-potentiel;
| Fait; fais; evoyer(ID',niveau') vers père;
-----

```

. Complexité théorique

Le nombre de messages est $(4n \lg n = 2,773n \lg n)$, la complexité temporelle est $O(n)$.

2.5.3 Algorithme B

La variable `niveau` représente ici le nombre total de processus **candidats éliminés**, c'est-à-dire, le nombre de racines battus.

Pour qu'un candidat `p` puisse capturer le noeud `v`, il faut que la condition `niveau(v) > niveau(p)` soit vérifiée. Quand le processus candidat `p` a atteint le processus ordinaire `v` qui est déjà capturé par un autre candidat `q`:

```

-----
| Si niveau(p) < niveau(v), p est éliminé; |
| Si niveau(p) > niveau(v), v est capture par p, |
|         niveau(v) := niveau(p); |
| Si niveau(p) = niveau(v), p est envoyé à q. |
| Arrivé à q : |
|   Si (niveau(p), ID(p)) < (niveau(q), ID(q)), p est éliminé. |
|   Si q a déjà été éliminé, p est éliminé aussi car cela |
|   signifie qu'il existe un candidat x dont la variable |
|   niveau est supérieure en valeur à celle de p. |
|   Si (niveau(p), ID(p)) > (niveau(q), ID(q)), alors |
|       (1) q est éliminé, |
|       (2) p incrémente sa variable niveau, |
|       (3) p capture v. |
-----

```

. Description formelle

a) Contexte-local-d'un-processus-Pi

Constante ID: entier valant `i`; (* identité d u noeud *)
 Variable niveau: entier initialisé à 0 ;
 ID-propriétaire, id-potentiel: entiers
 initialisés à 0 ;
 non-traverse: ensemble des liaisons adjacentes;
 père, père-potentiel: pointeurs initialisés
 à nil.

b) Types-de-messages

Un seul type de messages est utilisé: `(niveau, ID)`.

c) L'algorithme-exécuté-par-un-noeudi

```

-----
| Candidat:
| Lors de
| décision de provoquer une élection
| faire
|   e:= elt(non-traverse); { un element dans non-traverse }
|   envoyer(niveau,ID) sur e;
| fait;
| réception de (niveau', ID') sur e'
| faire si (ID=ID') alors { capture réussie }
|   niveau:= niveau+1;
|   non-traverse:= non-traverse-e;
|   si non-traverse 0  $\neq$  alors
|     e:= elt(non-traversé);
|     envoyer(niveau,ID) sur e;
|   sinon annoncer"élu, fin algorithme";stop
|   fsi
| sinon si ((niveau',ID') < (niveau,ID)) alors
|   ignorer ce message;
|   attendre un autre message;
|   sinon
|     (1) envoyer(niveau',ID') sur e';
|     (2) ignorer tous les messages futures;
| fait; fsi;   fsi;
-----

```

```

-----
| Ordinaire:
| niveau<--- -1;
| Lors de la reception de (niveau', ID') sur e'
| faire cas niveau' < niveau : écarter message;
|   niveau' > niveau : { remplacer le père }
|     père:= e';
|     niveau<--- niveau';
|     ID-propriétaire:= ID';
|     id-potentiel:= 0;
|     père-potentiel:= nil;
|     envoyer(niveau',ID') sur e';{ vers le père }
|   niveau' = niveau:
|     si (ID' < ID-propriétaire) alors
|       | père:= père-potentiel;
|       | niveau' := niveau'+1;
|       | ID-propriétaire:= ID';
|       | père-potentiel:= nil;
|       | envoyer(niveau',ID') vers le père;
|     I sinon
|       | si pere-potentiel <> nil
|       |   alors écarter le message
|       |   sinon { pas de père potentiel }
|       |     id-potentiel:= ID';
|       |     père-potentiel:= e';
|       |     envoyer (niveau', ID') vers le père;
|   Fait; fcas; fsi; fsi;
-----

```

. Complexité théorique

$2n \log n + 2n$ messages et $O(n \log n)$ unités de temps sont attendus.

La croissance de la complexité temporelle est due au fait que, à la différence de l'algorithme A, le niveau d'un candidat n'est pas une fonction du nombre de noeuds qu'il a déjà capturés. Un candidat qui réalise un "travail" important (et consomme plus de temps) en construisant une arborescence, pourrait être éliminé par un candidat ayant réalisé un travail moins important. L'arborescence qui aurait été construite par ce dernier serait moins volumineuse, mais le nombre de racines qu'il aurait battues serait plus élevé.

L'algorithme A quant à lui, est caractérisé par le fait qu'un noeud puisse être éliminé plusieurs fois (voir règle d'extinction de l'algorithme A).

Dans l'algorithme C ces deux problèmes sont résolus par la combinaison des deux techniques.

Le but fixé par l'auteur consiste à améliorer la complexité temporelle de l'algorithme B, en conservant la complexité en messages.

2.5.4 Algorithme C

Une fois qu'il a battu un processus ordinaire un candidat **incrémente** sa variable niveau qui sera au moins égale à **Log(nombre total de noeuds battus)**. Plus précisément, la variable niveau est mise à jour comme suit:

niveau: * **maximum** {**Log**(nombre de noeuds non racine battus (processus ordinaires)), nombre de noeuds racines battus}.

Cette règle d'utilisation de la variable niveau permet de prendre en compte le travail, effectué par un candidat en nombre de processus ordinaires et/ou candidats concurrents battus.

. Description formelle

La description formelle de l'algorithme C est similaire à celle de B. Une variable **taille** est utilisée pour compter le nombre de noeuds capturés par candidat. Le seul changement qui est introduit c'est la mise à jour de la variable niveau d'un processus candidat (première clause then):

```
taille := taille+1;
niveau:= Max(niveau', Log (taille));
```

. Complexité théorique

Les complexités ont été estimées à $2n(\log n + \log \log n + 2)$ messages, et $O(n \log \log n)$ unités de temps.

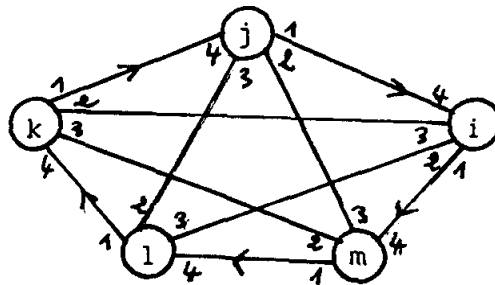
Remarque:

Le temps peut être amélioré à $O(n)$ unités de temps en apportant la modification suivante à l'algorithme [AFE 851: les candidats qui atteignent un niveau plus élevé que Logn tenteront de capturer n/Logn noeuds simultanément (en une unité de temps).

2.6 Algorithme de Loui, Matsushita et West [LOU 861**2.6.1 Hypothèse supplémentaire**

on admet que les processus ont un sens d'orientation globale.

Definition: On se fixe un cycle hamiltonien comprenant tous les processus. Chaque liaison est étiquetée par la distance le long de ce cycle, au processus se trouvant à l'autre extrémité d'un arc. En particulier, si le processus x est à la distance d du processus y , alors y est à la distance $(n-d)$ de x .

**2.6.2 Principe**

Initialement tous les processus sont actifs, ils deviendront tous passifs sauf l'élu (qui a l'identité la plus élevée). Un processus actif devient passif à la réception d'un message provenant d'un processus dont l'identité est supérieure à la sienne. Un processus actif qui reçoit un message émis par un processus dont l'identité est inférieure à sa propre identité ID , lui envoie en réponse un message contenant cette ID . Ceci permet de s'assurer que le processus en question passera dans l'état passif.

2.6.3 Description formelle**a) Contexte-local-d'un-processus- P_i**

Constante ID : entier valant i ; {identité du processus P_i }

Variable D : entier initialisé à $n-1$;
 E , $Newid$: entiers;
 $état$: (actif, passif, élu) initialise à actif;

b) Les primitives de communication

Envoyer(d:e,j): émet le message (e,j) le long de la liaison d au processus qui se trouve a la distance d.

Recevoir(e,j): attend l'arrivée du message (e,j).
 e: distance de l'émetteur au receptr, j: identité de l'émetteur.

c) Algorithme execute par un processus Pi

```

-----
| Lors de la
|   decision de provoquer une election
|   faire
|     I envoyer(n-D, n-D, ID); (* D a été initialisé à n-1 *)
|     fait;
|
| reception de (e, j)
|   faire
|     I Newid:=-j;
|     I si état=actif alors
|       I I           D:=e;
|       I I           cas Newid < ID : envoyer(n-D, n-D, ID);
|       I |           Newid = ID : état:=élu;
|       I |           annoncer "élu";
|       I I           Newid > ID : état:= passif;
|       I I           fcas;
|       I I           sinon { état=passif }
|       I |           E:=e;
|       I |           envoyer(n-D,(n-D+E), Newid)
|     I fsi;
|     fait;
|
-----

```

. Complexité théorique

Le nombre de messages est inférieur a $3,62n(c \rightarrow d C(n))$ et le temps est borne par $n + \text{Log}_\phi n + O(1)$ unités de temps; où $\phi = \frac{1}{2}(1 + \sqrt{5})$.

2.6.4 L'algorithme corrigé

L'algorithme donné par les auteurs (voir description détaillée ci-dessus) n'est pas correct(voir annexe). L'erreur se situe au niveau du calcul de distance. En effet, l'idée sur laquelle est basé l'algorithme consiste en ce qui suit:

Dans l'état passif, un processus y recevant un message (e,j) provenant du processus x (donc y est à la distance e de x), envoie un message à destination du processus z, qui l'avait, a un moment antérieur, mis dans l'état passif. Le processus z se trouve à la distance N-d de y, et a la distance e' de x. Il faut noter que le message envoyé par le processus passif y, évitera aux processus actifs de lui envoyer d'autres messages par la suite. Après l'émission de ce message, aucun processus

n'émettra un message à y. La distance e' incluse dans le message envoyé par y à z, permettra à z de communiquer directement avec x sans transiter par y (qui est passif). En tenant compte du sens d'orientation, deux schémas (figure III.1) sont possibles, selon que c'est x ou z qui précède immédiatement y le long du cycle.

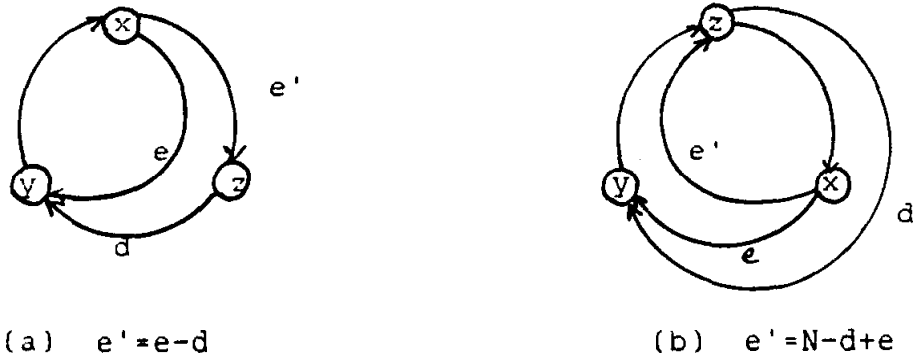


Figure III.1

Or dans l'algorithme, la valeur prise en compte est celle où la situation (b) se présente.

L'algorithme n'aboutit pas si le cas (a) se présente. Le même algorithme est reproduit avec la modification de la clause:

```

sinon {état=passif}
  E:=e;
  envoyer(n-D, (n-D+E), Newid)

```

qui devient:

```

sinon { état=passif }
  E:= e;
  si (e-d)>0 then envoyer(n-D,(e-D), Newid)
  else envoyer(n-D,(n-D+E), Newid)
fsi;

```

Remarque :

Masuzana et al. ont proposé dans [MAS 89] une version modifiée de l'algorithme de [LOU 86]. Cet algorithme a été transformé de façon à résister aux pannes de processus. L'idée sur laquelle est basée cette transformation consiste, pour un processus donné, à utiliser $(fp+1)$ jetons au lieu d'un jeton. Le nombre $fp < n/2$ représente les processus défectueux. En effet, l'utilisation de $fp+1$ jetons, permet de pallier au problème d'interblocage, car les jetons envoyés à des processus défectueux sont perdus.

2.7 Algorithme De Attiya, Leeuwen Et Shmuhel [ATT 891

Les auteurs de [ATT 891 exploitent ici la notion de sens d'orientation pour **généraliser** le travail **présenté dans** [LOU 861. Ils présentent un algorithme **général d'élection** pour des réseaux dont la topologie est l'anneau **mais** auquel, des liaisons ont **été** ajoutées. Les cas extrêmes **étant** l'anneau et le réseau a topologie complète ($n-1$ cordes ont **été** ajoutées a l'anneau [LOU 86]). Ils proposent un corollaire qui montre que l'ajout de $O(\log n)$ cordes a chaque processus du **réseau** suffit pour obtenir un algorithme utilisant au plus $O(n)$ messages.

. Principe

L'idée **utilisé** est la **même** que celle de [LOU 861, et le **principe** de base est celui de [FRA 821. A chaque phase, un ensemble de processus est **candidat à l'élection**. Durant l'**étape** i un processus qui est maximum **local**(son ID est **supérieur** a ceux de ses deux voisins **actifs** les plus **proches**), passe à l'étape $i+1$ et **échange** son ID **avec** ses deux voisins **actifs**. Ces derniers sont **repérés** sur l'anneau grace à un **échange** de messages permettant a chaque processus qui survit a une **étape donnée** de **mesurer** la distance, le long de l'anneau, entre lui et ses deux voisins **actifs** les plus proches.

3 ALGORITHMES D'ELECTION POUR TOPOLOGIE QUELCONQUE

3.1 Introduction

Dans cette **partie**, sont **présentés** des algorithmes s'exécutant sur des réseaux **modélisés** par un graphe quelconque. D'une part ces algorithmes calculent un extremum, et d'autre part, ils construisent une arborescence couvrante sur le graphe permettant a **tous** les processus de **communiquer** via des chemins uniques **avec le processus élu associé** a la **racine**.

Pour la topologie quelconque, les algorithmes **conçus** peuvent **être** divisés en deux grandes classes: des algorithmes proposés pour des réseaux anonymes, et des algorithmes **conçus** pour des réseaux dont les processus sont dotés d'identités toutes distinctes. Dans la deuxième **catégorie**, il y a les deux algorithmes proposés par [HEL 86b], basés respectivement sur les techniques d'exploration de réseaux en profondeur et en parallèle. Ils utilisent la notion de transfert de connaissances **avec** une **complexité théorique** en $O(n^2)$ dans le pire des cas.

S. Skyum [SKY 871 a également **publié** un **algorithme** basé sur une technique d'exploration de réseaux en profondeur. Ce dernier atteint une **complexité** en $O(n \log n) + O(E)$ messages et c'est aussi la **complexité suggérée** dans [LAV 891].

Dans le cas de réseaux anonymes, on peut **citer** les algorithmes de Matias et Afek [MAT 891 dont la **complexité** est en $O(Er \log n \log r)$ messages où $r = 1/E$ et $(1-E)$ est la probabilité

avec laquelle un leader unique est choisi. De son côté Lavault [LAV 87] avait proposé un algorithme qui construit un arbre couvrant le réseau, la racine étant le processus élu. Sa complexité est de $n^2 + O(n)$. Un autre algorithme est proposé par [LAV 90]. Il choisit un Clu avec une probabilité de $1-E$ et une complexité de $O(E+n \log n)$.

3.2 Hypotheses

- Le graphe $G(V,E)$ est quelconque et les liaisons sont bidirectionnelles.
- Il n'y a pas de perte de messages.
- Les messages sont **délivrés** au bout d'un temps non prévisible **mais fini**, de plus, ils ne sont pas **altérés**.
- La communication est asynchrone.
- Au cours du **calcul distribué**, un processus quelconque n'apprendra jamais la structure **globale** du réseau, ni **même** le nombre total de processus participants: sa connaissance topologique se **limitera** toujours à ses voisins directs.
- Il y a **symétrie** du texte de l'algorithme exécuté par l'ensemble des processus.

Ces hypothèses sont en général **respectées**, dans le cas contraire, d'autres hypothèses seront données.

3.3 LES SITES SONT DOTÉS D'IDENTITÉS DISTINCTES

Les algorithmes proposés dans ce paragraphe, supposent que chaque processus est **dôté** d'une identité unique.

3.3.1 Algorithme de Helary, Haddi Et Raynal [HEL 86b] (exploration en profondeur)

3.3.1.1 Principe

En plus des hypothèses **citées** précédemment, cet algorithme **admet** le déséquencelement de messages. Le **principe** sur lequel repose l'algorithme proposé **consiste** pour un processus P_i , à lancer une exploration du réseau dans le but de visiter **tous** les processus en établissant des chemins de ceux-ci vers **lui-même**. Comme l'algorithme peut **être démarré** par un ou plusieurs processus, une ou plusieurs explorations peuvent **être** simultanément en **cours**. À chacune de ces explorations est associé un poids correspondant au processus initiateur de celle-ci.

Une exploration E_k (issue de P_k) peut **être** vue comme une **activité** séquentielle; une séquence de messages **est** engendrée, chaque message est **traité** par le processus qui **le reçoit** et celui-ci, à son tour, peut éventuellement émettre un nouveau message de poids k . À chaque instant de l'exploration, **un** seul message de poids k est en transit sur une ligne ou en traitement par un processus. Les auteurs décomposent le processus d'élection en **étapes matérialisées** par la réception d'un message m par un processus P_i :

1. Etape d'extinction:

-ou bien P_i a déjà été atteint par une exploration de poids supérieur à k (ou a déjà lancé sa propre exploration E_i si $i > k$),

-ou bien P_i n'a pas encore été atteint, mais $i > k$.

Dans un cas comme dans l'autre, l'exploration E_k est éteinte par P_i ; de plus, dans le deuxième cas, P_i lance l'exploration de poids i .

2. Etape de conclusion:

Dans cette étape et les deux suivantes, on suppose que E_k est l'exploration de poids le plus fort ayant jamais atteint P_i ; k désignera la plus grande identité de l'ensemble des processus du réseau. L'information de contrôle amenée par le message indique à P_i que tous les autres processus ont été atteints par l'exploration. Dans ce cas, E_k s'arrête et on a alors nécessairement: $k = \max$; P_i peut donc conclure que l'exploration E_{\max} a établi une arborescence recouvrant le réseau et les autres processus sont informés de la fin de l'algorithme par un signal diffusé sur les branches de cette arborescence.

3. Etape de génération:

P_i doit faire progresser E_k , et possède des voisins non encore atteints par E_k . Il met alors à jour les liaisons relatives à E_k , et engendre un nouveau message d'exploration de poids k vers un de ses voisins non atteints (par exemple le plus grand).

4. Etape de renvoi:

P_i doit faire progresser E_k ; tous ses voisins ont été atteints par E_k , mais l'information de contrôle amenée par le message indique que des processus du réseau n'ont pas encore été atteints par E_k . P_i engendre alors un message de poids k permettant de remonter vers un processus non encore atteint.

3.3.1.2 Description formelle

al Contexte local d'un processus P_i

Constante voisins: initialisé à {identités des voisins de P_i };

Variable état: (initial, candidat, battu, élu) initialisé à initial;

pgvu: entier initialisé à i ;
{donne l'identité de la plus grande exploration, à la fin pgvu = identité maximum}

pred : entier initialise a nil;
succ: ensemble d'entiers initialise a A;

b) Type de messages

(i) explorer(k,z,s): k represente l'identite du processus ayant initialise l'exploration, z contient les identites des processus visités par ce message, et s définit l'ensemble des processus voisins immédiats des processus de z.

(ii) rebrousser(k,z,s): les cahmps de message sont les même que ceux de type explorer.

(iii) conclure servent a signaler la terminai son de l'algorithmme.

c) L'algorithmme exécuté oar un processus Pi

```

-----
| Procédure lancer-exploration |
|   debut soit x=maximum(voisins); |
|   |   état:=candidat; |
|   |   pred:=nil; |
|   |   succ:={x}; |
|   |   envoyer explorer(i,{i},voisins-{i}) a Px |
|   fin; |
| Lors de |
|   decision de provoquer une election |
|   faire |
|   | si état=initail alors lancer-exploration fsi |
|   | fait; |
|   reception de explorer(k,z,s) depuis Pj |
|   faire |
|   | cas pgvu>k: si état=initial alors |
|   | |   lancer-exploration fsi; |
|   | |   pgvu<k: étati:= battu; |
|   | |   pgvu:= k ; |
|   | |   pred:=j; |
|   | |   soit y=voisin-z; |
|   | |   cas y=∅ : succ:=∅ ; |
|   | |   | cas s=∅ : envoyer conclure à Pj; |
|   | |   | s<>∅: envoyer rebrousser(k,zU{i},s) |
|   | |   | fcas; à Pj; |
|   | |   | y<>∅: soit x=maximum(y); |
|   | |   | succ:= {x}; |
|   fait;fcas; fcas; envoyer explorer(k,zU{x}) a Px |

```

```

reception de rebrousser(k,z,s) depuis Pj
  faire si pgvu-K alors
  |   |   soit y-voisins s;
  |   |   cas y=∅: envoyer rebrousser(k,z,s) a Ppred
  |   |   y<>∅: soit x-maximum(y);
  |   |           succ:=succU(x);
  |   |           envoyer explorer(k,z,s-{x}) a Px
  |   |   fcas;
  |   fsi;
  fait;

reception de conclure depui Pj
  faire si pgvu=i alors état:= Clu fsi;
  |   ∀ x∈ (succ U pred)-{j}: envoyer conclure a Px
  fait;

```

. Complexité théorique

Dans le cas le plus favorable, oh Pmax est seul a lancer l'exploration:

- . n-1 messages du type explorer,
 - . 0 à n-1 messages du type rebrousser (0 est obtenu dans le cas d'un reseau en anneau),
 - . n-1 messages du type conclure.
- soit $3(n-1)$ messages ($2(n-1)$ pour un reseau en anneau).

Dans le cas le plus defavorable:

- . au plus ($\sum_{k=1}^{n-4} k'$) + (n-1) messages du type explorer,
 - . autant messages rebrousser,
 - . et, n-1 messages conclure,
- soit $(n-1)(n-3)$ messages tous types confondus. La complexite est donc $O(n^2)$ dans le pire des cas.

La complexite temporelle est de $[2(n-1)+d]D$, où d est le diametre de l'arborescence construite ($1 \leq d \leq n-1$), et $2(n-1)D$ le temps de transfert des messages explorer et conclure. dD est le temps maximum pour réveiller Pmax. D est le délai d'acheminement des messages entre deux processus quelconques.

Pour un reseau complet, le nombre de messages est $3(n-1)$, et $(2n-3)D$ le temps d'election.

Pour un reseau en anneau, le nombre de messages est $n/2(n/2+1)-1$ c'est-a-dire $O(n)$, et le temps $(3n-4)D$.

3.3.2 Algorithme de Helary, Madi Et Raynal [HEL 86b] (exploration en largeur)

3.3.2.1 Principe

Dans cet algorithme l'exploration issues d'un processus P_k est lancée en **parallèle** vers **tous** ses voisins. Comme précédemment un processus ne laisse progresser vers ses voisins la composante de l'exploration qui l'atteint que si le poids k de **celle-ci est supérieur** au poids de la dernière exploration prise en compte.

Pour obtenir une arborescence couvrante, il est **nécessaire** que tout processus autre que P_{max} ait un seul prédecesseur et que deux processus **distincts** aient des ensembles de successeurs disjoints.

Tout processus P_i , lorsqu'il **reçoit** un message relatif a une exploration qu'il prend en compte pour la première fois, **mémorise** l'émetteur P_j du message comme étant son prédecesseur. Il **propage** l'exploration vers ses voisins, et attend d'avoir **reçu** des acquittements de **tous** ceux-ci, avant d'émettre l'acquittement vers p_j . Dans le cas où il **reçoit** un message relatif a l'exploration déjà prise en compte, il l'acquitte immédiatement en indiquant a P_j de ne pas l'inclure dans ses successeurs.

3.3.2.2 Description formelle

a) Contexte local d'un processus P_i

La constante voisins et les variables @tat, p_{gvu}, pred, succ ont la même signification que dans le cas de l'exploration en profondeur.

```
variable nbacq: entier non negatif initialise a 0;
           {permet de memoriser le nombre d'acquittements
            attendus et de contrôler l'envoi de
            l'acquittement vers son prédécesseur}
```

b) Type de messages

(i) explorer(k,z): k représente l'identité du processus ayant initialisé l'exploration, z désigne l'ensemble des processus dont on est sûr qu'ils ont été ou seront visités par des messages relatifs a cette exploration}

(ii) acquitter(k,x): x permet au recepateur de savoir s'il doit considérer P_i comme un de ses successeurs dans l'arborescence (x=terminé) ou non (x=déjàvu).

c) L'algorithme exécuté par un processus P_i

```

-----
procédure lancer-exploration
début
  |   &tat:-candidat;
  |   pred:=nil;
  |   succ:=voisins;
  |   nbacq:=cardinal(succ);
  |    $\forall x \in \text{succ}$ : envoyer explorer(i,voisins-{i}) à  $P_x$ 
fin;

lors de
  décision de provoquer une élection
  faire
    | si état=initail alors lancer-exploration fsi
  fait;
  réception de explorer(k,z) depuis  $P_j$ 
  faire
    | cas pgvu>k: si état=initial alors
    | |   lancer-exploration fsi;
    | |
    | |   pgvu-k : envoyer acquitter(z,déjàvu) à  $P_j$ ;
    | |   pgvu<k: état:= battu;
    | |   pgvu:= k;
    | |   pred:= j;
    | |   succ:= voisin-z;
    | |   cas succ=∅: envoyer acquitter(z,terminé)
    | |   |   à  $P_j$ ;
    | |   succ<>∅: nbacq:= cardinal(succ);
    | |    $\forall x \in \text{succ}$ : envoyer acquitter(k,zUsucc)
    | |   |   à  $P_j$ ;
    | |
    | fcas;          fcas;
  fait;
  réception de acquitter(k,x) depuis  $P_j$ 
  faire si pgvu=K alors
    | |   si x=déjàvu alors succ:=succ-{j} fsi;
    | |   nbacq:=nbacq-1;
    | |   si nbacq=0 alors
    | | |   cas k-i: état:= élu;
    | | |   |   k<>i: envoyer acquitter(x,terminé) à
    | | |   |   |   Ppred
    | | |   |   fcas;
    | | fsi;fsi;
  fait;
-----

```

. Complexité théorique

La complexité dans le pire des cas est en $O(n)$ messages, et de $3d_{\max}D$ en temps. d_{\max} est la distance maximale de P_{\max} à tout autre processus ($1 \leq d_{\max} \leq n-1$). $2d_{\max}D$ est le temps de transfert des messages explorer et acquitter, relatifs à l'exploration P_{\max} , et $d_{\max}D$ le temps nécessaire pour réveiller P_{\max} .

3.3.3 Algorithme de **Sven** Skyum [SKY 871]

. Principe

S. Skyum propose un algorithme basé sur une technique de **parcours** de **réseaux** en profondeur. Quand un processus v se réveille, il **crée** un **jeton** $\langle s, 0, v \rangle$ de niveau 0 et ce **jeton** commence un **parcours** d'exploration en profondeur. Le symbole s indique si le **jeton** **circule** en avant, en arrière ou chasse un autre **jeton**. Chaque fois qu'un **jeton** fait un pas en arrière sur un arc, **soit** il ferme l'arc, **soit** il le déclarera appartenant à l'arbre couvrant et il ne sera plus jamais traversé par un message.

Si un **jeton** $\langle s, l, u \rangle$ émis par un processus u atteint un processus qui, a déjà été traversé par un **jeton** d'un niveau supérieur, le **jeton** $\langle s, l, u \rangle$ est éteint. S'il arrive à un processus qui a été visité par le **jeton** $\langle s, l, w \rangle$ de même niveau, alors:

- ou bien u est inférieur à w , le **jeton** $\langle s, l, u \rangle$ s'arrête en v et attend qu'un autre **jeton** le rattrape. Ceci est réalisé en mettant le processus v à l'état (candidat, niveau 1).
- ou bien u est supérieur à w , le **jeton** $\langle s, l, u \rangle$ arrêtera son parcours et poursuivra $\langle s, l, w \rangle$ aussi longtemps qu'il n'est pas chassé, à son tour, par un troisième **jeton** de niveau 1, auquel cas $\langle s, l, v \rangle$ s'arrêtera et l'état du processus où a lieu l'arrêt devient (candidat, niveau 11).

Lorsque deux jetons de même niveau se rencontrent, un nouveau **jeton** de niveau $l+1$ est **créé**, et un nouveau parcours en profondeur du graphe est lancé, alors que les deux jetons de niveau 1 sont éteints. Si un **jeton** réussit à terminer son parcours, il déclarera le processus initiateur comme élu et on a un arbre de recouvrement qui relie tous les processus au processus élu.

. Complexité théorique

Elle est de $2E+2n \log n + O(n)$ messages, où E représente le nombre de liaisons connectant les processus du réseau.

3.2.4 Algorithme de Lavalée et Lavault [LAV 8911]

. Principe

Cet algorithme **construit** un arbre couvrant et désigne la **racine** de l'arbre comme l'élu. La construction de l'arbre couvrant se fait, progressivement, **par** combinaison de **sous**-arbres en fragments de plus en plus grands, en respectant la relation d'ordre total strict existant dans l'univers I auquel appartient les n identités.

Au **début**, chaque processus constitue un fragment dont il est la racine. Le processus de connexion de fragments consiste en ce qui suit:

a) un fragment F **donné** est **candidat** à la fusion avec un autre fragment G, et la racine de F **emet** une **requête** de connexion vers G.

b) une **telle requête** est **acceptée**(ou **rejetée**) par le premier processus du fragment G qui l'aurait **reçue**, en tenant **compte** de la relation d'ordre strict **existant** entre les identités.

. Complexité théorique

Cet algorithme **possède**, dans les réseaux à maillage **complet**, une complexité en messages de $1/2n \log n + O(n)$ dans le pire des cas et une complexité de $O(n)$ en moyenne. Pour la topologie en anneau, dans le pire des cas, la complexité en messages est $3n \log n + O(n)$, et en moyenne, elle est égale à $1/2n \log n + O(n)$.

Dans les réseaux **d'interconnexion** quelconque, un bon majorant de la complexité en messages est $2E + 7/2n \log n + O(n)$, mais dans le cas où n est connu, la borne supérieure est de $7/2n \log n + O(n)$. Pour le temps la **complexité** est de $O(\log n)$ unités.

3.4 ALGORITHME POUR RESEAUX ANONYMES

Sous l'**hypothèse d'identités** distinctes des sites du réseau, le problème d'élection a été ramené au **problème de calcul** d'un extremum et diverses solutions ont été proposées. Si l'on ne s'astreint pas à l'hypothèse **d'unicité** des identités des sites (**soit** à cause de l'absence d'identités, **soit** parce que les **identités** existent **mais** ne sont pas toutes distinctes), le problème d'élection d'un site ne peut être résolu par une **méthode** de calcul d'extremum.

Dans un **réseau** anonyme, initialement **tous** les processus n'ont aucune **identité** et, chaque processus ne connaît que les ports le reliant à ses voisins directs. Au **cours** de l'algorithme, il n'apprend d'eux que leurs identités. Il est **montré** dans [ANG 80] que si les sites sont **tous** anonymes (ou **possèdent** **tous** une **identité** commune), il n'existe alors aucune solution **déterministe** au problème de l'**élection distribuée**, que le réseau **soit** **synchrone** ou asynchrone. Par **conséquent**, toute solution au problème ne saurait être obtenue, dans ce cas, que par un algorithme non **déterministe**. Santoro et Lavallée [SAN 86, LAV 90] soulignent que le **problème** de l'élection peut être résolu par un algorithme probabiliste si les hypothèses de symétrie de sites et de leur connaissance de n (**nombre** total de sites) sont faites conjointement.

3.4.1 Algorithme de Christian Lavault [LAV 871]

Les processus connaissent un majorant du nombre de sites, Maxint, fixé par la machine.

. principe

La première étape de l'algorithme est l'étape d'activation-dénomination. En effet, initialement tous les processus sont dans un même état dit **état initial**. Après le passage spontané des processus vers l'état **actif**, comme ils sont anonymes, leur première **tâche consiste** à se donner une identité par génération d'un nombre aléatoire dans $[1..Maxint]$. Donc, plusieurs processus peuvent acquérir des identités égales. Après le choix de l'identité, les processus passent dans l'état candidat.

Dans la deuxième étape est construite une arborescence couvrante du graphe G associé au **réseau**. Un ou plusieurs fragments peuvent être simultanément **construits** sur G. Ces fragments seront fusionnés deux à deux, chacun cherchant à en absorber un autre, **jusqu'à** obtention d'une arborescence **constituée** de tous les processus du **réseau**.

Les identités des processus sont, **comparées** selon la relation d'ordre total non strict **existant** entre elles du fait de l'existence **d'identités** non-distinctes.

Dans la **mesure** où il est possible que plusieurs processus (et, éventuellement **tous** les processus, **avec** une probabilité infime de $1/(maxint)$) puissent tirer aléatoirement le même numéro d'identité entre 1 et Maxint, le processus de **négociation-connexion** peut **échouer** et la terminaison **correcte** de l'algorithme peut en **être affectée**. Pour éviter cela, une procédure permettant de **réaliser** autant de nouveaux tirages **aléatoires d'identités** qu'il est **nécessaire** a été prévue.

Par exemple, un sommet quelconque P_i (racine ou non), **reçoit** une demande de **négociation** de la part d'une racine quelconque dans le **réseau** possédant la même identité, l'algorithme suivant est exécuté:

1. le sommet P_i renvoie un message de refus à la racine qui demande la **négociation**.
2. Tant qu'il reste encore un fils **sortant** P_x de P_i possédant une **arête sortante** non-utilisée et **d'identité** différente de celle de P_x (puisque P_x est fils de $P_i \implies ID_x = ID_i$), la **négociation** est **tentée** à travers P_x .

Par **contre**, lorsque chaque fils **sortant** P_x de P_i ne possède que des portes sortantes dont les identités sont **égales** à l'identité de P_x , alors P_i **doit** changer **d'identité** et **effectuer** un nouveau **tirage** aléatoire uniforme dans

[1,..Maxint]-{IDi}.

3. L'ensemble fils-sortants finit par être vide, et il ne reste plus de portes à explorer. Il reste alors une seule racine qui diffuse un message fin à tous les sommets.

. Complexité théorique

Le nombre de messages attendu est minoré par $11n/2+O(n)$; il s'agit du cas où toutes les identités sont toutes distinctes. Il est majoré par $n+O(n)$; c'est le cas où toutes les identités sont égales.

Le temps d'exécution est minoré par $O(n\sqrt{11n\log n})$, et majoré par $O(\log n\sqrt{11n})$.

Remarques sur la terminaison:

Il y a une autre façon d'assurer la terminaison de l'algorithme lorsque chaque nœud connaît la taille n du système distribué. Comme chaque racine r connaît la taille de son propre fragment F , alors que F est fusionné à un autre fragment G ; il envoie sa taille à la racine de G . Ainsi, si un fragment atteint éventuellement la taille n , sa racine se déclare Clue, et l'algorithme est terminé. Une variable compteur de nœuds est initialisée à la valeur 1 et mise à jour après chaque fusion entre deux fragments.

3.4.2 Algorithme de Yossi Matias et Yehuda Afek [MAT 89]

. principe

Le processus P_i recevant une identité ID_i diffuse à tous ses voisins toutes les identités qu'il a reçues auparavant, sont inférieures à ID_j . Cela signifie que c'est le processus portant la plus grande identité qui capture le réseau. De plus, chaque processus mémorise la liaison par laquelle il a reçu le plus grand ID, ce sera son processus-père. Il y aura donc la construction d'une forêt couvrante.

Si le plus grand ID est unique, alors la forêt couvrante se réduit à un arbre dont la racine est le leader.

La difficulté de l'algorithme réside dans la procédure de sélection d'un ID. Elle est conçue de façon à ce que deux conditions soient satisfaites simultanément avec une grande probabilité:

- (a) un processus unique est distingué parmi les autres nœuds;
- (b) le nombre d'identités distinctes est faible.

La procédure choisir(IDi):

Chaque processus P_i sélectionne deux nombres t_i et s_i de la façon suivante:

- (1) t_i : est le nombre de lances avant l'obtention de face;

(2) si: est un nombre choisi au hasard dans l'intervalle $[1..d]$, oh $d = O(r \log r)$ avec $r = 1/\epsilon$ et $1-\epsilon$ étant la probabilité de succès de l'algorithme.

Donc l'identité du processus P_i est $ID_i = \langle t_i, s_i \rangle$.

L'algorithme ELECT exécuté par P_i

```

-----
| Procédure initialise(i)
| début
|   ! appel choisir(IDi);
|   ! max:=IDi; père:=0;
|   ! diffuser IDi vers tous les voisins.
| fin.
| Lors de la réception de ID a travers la liaison l
| faire
|   I ! si IDi n'est pas défini alors appel initialise(i) fsi;
|   I ! si ID <= max alors ignore
|   ! !           sinon max:=id; père:=1;
|   I ! !           diffuser ID a tous ses voisins
|   ! fsi;
| fait.
-----

```

. Complexité théorique

L'algorithme ELECT s'arrêtera éventuellement. Apt-es la terminaison! un leader unique est choisi, et un arbre couvrant est construit avec la probabilité $(1-\epsilon)$. La complexité en messages attendue est $O(E r \log n)$, E est égal au nombre d'arcs du graphe des processus. La complexité temporelle est $O(D)$ oh D est le diamètre du réseau. Ces complexités sont données avec une probabilité supérieure à $(1-\epsilon)$.

3.4.2.3 L'élection avec détection de la terminaison

L'algorithme ELECT a été modifié de façon à permettre la détection de la terminaison. En effet, une hypothèse supplémentaire(H) a été introduite:

H: la valeur n de la taille du réseau est supposée être bornée inférieurement par L , et supérieurement par $2L$: ($L < n \leq 2L$).

Après la terminaison de l'algorithme ELECT, Max a la même valeur au niveau de tous les noeuds, et les liaisons parentes des noeuds constituent une forêt couvrante du réseau avec un ou plusieurs arbres. L'algorithme a réussi l'élection d'un leader si la forêt obtenue ne contient qu'un seul arbre.

Donc si $L < n \leq 2L$, l'arbre dont la racine est le noeud élu doit avoir une taille supérieure à L . De plus, s'il existe un arbre dont la taille est supérieure à L dans la forêt, on est certain qu'il est le seul à avoir cette taille, par conséquent, l'élu est unique.

3) Degré de répartition(ou symétrie)

Tous les algorithmes étudiés possèdent la propriété de symétrie de texte au sens défini dans [RAY 851.

4) Connaissance de l'état global ou local

- E1. Les processus connaissent la taille n du réseau ou une borne sur n .
- E2. Les processus ne connaissent ni la taille n du réseau ni une borne sur n .
- E3. Les processus connaissent la topologie du réseau.
- E4. Les processus ne connaissent rien sur la topologie du réseau.
- E5. Les processus connaissent les identités de leurs voisins immédiats.
- E6. Les processus connaissent les identités des ports d'entrée/sortie auxquels ils sont connectés.

5) Topologie en anneau

| Design.1 algo. | Hypo. I communic. | Hypo. I réseau | I info. pannes | glob./loc. | complexités messages | I temps. |
|----------------|---------------------------------|----------------|----------------|------------|---|----------|
| [CHA 7 8 1 | C2, C4, C5 C6, C9 | S1 | C13, S3 | E2, E4, E6 | $n(n-1)/2$ IO(n) $O(n \lg n)$ | $O(n)$ |
| [DOL 8 2 1 | C2, c4, c5 C6, C7, C9 | S1 | C13, S3 | E2, E4, E6 | $2n \lg n + n$ pire | / |
| [HIR 8 0 1 | C1, C4, C5 C6, C9 | S1 | C13, S3 | E2, E4, E6 | $O(n \lg n)$ pire | / |
| [KOR 8 1 1 | C1, C4, C5, C6 C7, C9, C12 | S1 | c13, S3 | E2, E4, E6 | $n^2/2$ $n \lg n$ | / |
| [BOD 8 6 1 | mêmes que Ici-dessus | S1 | c13, S3 | E2, E4, E6 | $\sqrt{306}n^2/4$ $n \lg n$ | / |
| [SPI 8 9 1 | C2, c4, C5 C6, C9 | S2 | c13, S3 | E1, E4, E5 | $O(n)$ moy. | / |
| [FRA 8 2 1 | C1, c3, C5 C6, C9 | S1 | C13, S3 | E2, E4, E6 | $2n \lg n + 3n$ pire | / |
| [LEE 8 7 1 | C1, c4, C5 C6, C7, C9 | S1 | C13, S3 | E2, E4, E6 | $4,44 \lg n + O(n)$ pire | / |
| [GOL 8 7 1 | C1, c4, C5 C6, C10 | S1 | C14, S3 | E1, E4, E6 | $O(n \lg n)$ pire | / |
| [FRE 8 7 b] | C2, c3, C5 C6, C9 | S1 | C13, S3 | E2, E4, E6 | $O(4n)$ pire | $O(n^2)$ |

6) Topologie complete

| Design. algo. | Hypo. communic. | Hypo. réseau | pannes | info. glob./loc. | complexites messages | temps |
|------------------|---------------------------|--------------|---------|------------------|--------------------------------------|--------|
| [GAR 82] Brutal | C1, C4, C5 C6, C7, C9 | S1 | C13, S4 | E1, E3, E5 | $O(n^2)$ pire | / |
| [GAR 82] Invit. | C1, C4, C5 C6, C7, C10 | S1 | C14, S4 | E1 | $O(n^2)$ pire | / |
| [KOR 84] | C1, C4, C5 C6, C7, C13 | S1 | C13, S3 | E2, E4, E6 | $5n \log k + O(n)$ $k = nb$ initl | / |
| [AFE 85] Algo. A | C1, C4, C5 C6, C7, C9 | S1 | C13, S3 | E2, E4, E6 | $2,773n \log n$ | $O(n)$ |
| [AFE 85] B | même | même | même | même | $2n \log n + 2n$ | / |
| [AFE 85] C | même | même | même | même | $2n \log n + \log^2 n + 2$ | / |
| [LOU 86] orient. | même + | S1 | C13, S3 | E2, E4, E6 | $3.62n$ | |
| [ATT 89] | I même que | [LOU 86] | / | / | / | $O(n)$ |

71 Topologie quelconque

| Design. algo. | Hypo. communic. | Hypo. réseau | pannes | info. glob./loc. | complexites messages | temps |
|-------------------|--------------------------|--------------|---------|------------------|---|--------|
| [HEL 86] Profond. | C1, C4, C6 C8, C9 | S1 | C13, S3 | E1, E4, E5 | $3(n-1)$ $(n-1)(n-3)$ | $O(n)$ |
| [HEL 86] Largeur | même que ci-dessus | même | même | même | $O(n^2)$ | / |
| [SKY 87] Invit. | C1, C5, C6 C9 | S1 | C13, S4 | E2, E4, E6 | $2E + 2n \log n + O(n)$ | |
| [LAV 89] I | C1, C4, C5 C6, C7, C9 | S1 | C13, S3 | E2, E4, E5 | $n \log k + O(n)$ $3n \log n + O(n)$ | |
| [LAV 87] I | C1, C4, C6 C7, C9 | S2 | C13, S3 | E1, E2, E4, E6 | $n + O(n)$ $n^2 + O(n)$ | |
| [MAT 89] I | C1, C4, C6 C9 | S2 | C13, S3 | E2, E4, E6 | | $O(D)$ |

(+): D délai d'acheminement des messages entre deux processus quelconques.

VII Les techniques **utilisées**

Dans ce paragraphe sont resumés les techniques et les mécanismes utilisés dans la conception des algorithmes d'élection. Le **problème** d'élection peut se ramener à deux **sous-problèmes**:

1. construction d'un algorithme de parcours,
 2. construction d'un algorithme d'élection utilisant l'algorithme défini en (1.) afin de transmettre les messages vers **tous** les processus **impliqués** dans l'élection.
- 1 Election quand les processus sont **dotés d'identités** distinctes

Dans le cas d'un **réseau** en anneau, le parcours **consiste**, pour un processus initiateur à transmettre un message vers un de ses deux voisins. Ce message est transmis le long de l'anneau jusqu'à atteindre son initiateur. Les algorithmes d'élection qui ont été étudiés sont basés sur le principe d'extinction sélective.

Extinction **sélective**

Chang et Roberts [CHA 79] ont proposé, en 1979, le principe de l'extinction sélective des messages pour un anneau unidirectionnel.

L'algorithme de Chang et Roberts comprend les trois phases suivantes:

1. La recherche du processus initiateur de plus fort numéro et l'arrêt des autres candidatures par extinction sélective. Pour cela chaque initiateur diffuse un message de reconnaissance contenant son numéro i . Si un tel message arrive sur un processus initiateur j de plus fort numéro, il n'est pas transmis, ce qui **annule** la candidature de i . Si $i > j$, le message est **propagé** par j . Il en est de même si le processus j n'est pas initiateur.
2. La fin de l'élection: elle est obtenue quand le message d'un site initiateur i a parcouru tout l'anneau et est revenu à son initiateur: $i=j$. **Seul** le message du processus de **numéro** maximal **peut** faire le tour de l'anneau **sans être** annulé.
3. L'annonce de l'élection: le processus **élu** fait **circuler** un message de fin de reconnaissance, **avec** son **numéro**, sur tout l'anneau; il informe ainsi **tous** les processus, de son **identité** et de la fin de l'élection.

Franklin, Hirshberg et Sinclair [HIR 80, FRA 82] ont adopté le principe de l'extinction sélective au cas de l'anneau bidirectionnel.

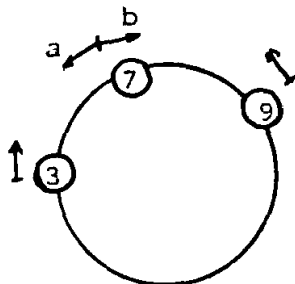
De nombreux auteurs ont propose des ameliorations à ces algorithmes initiaux et ont **abouti** parfois a des algorithmes très complexes [DOL 82, PET 82, GOL 87, ROT 87, LEE 87].

La **spécificité** des methodes d'extinction selective les unes par rapport aux autres repose [KAI 88] dans la phase 1. Il faut éliminer au plus vite les **petits** numeros. Dans **tous** les cas, l'unicite de **l'élú résulte** de l'unicite de la numérotation des processus, de l'unicite du processus dont le message de reconnaissance fait le tour de l'anneau et de l'unicite du message de fin d'élection. Les performances varient d'une election a une autre puisqu'elles dependent de la **numérotation** relative des processus.

Pour **éliminer** au plus vite les **petits** numeros il est possible d'appliquer une des methodes suivantes:

Méthode 1: Dans la variante probabiliste de l'algorithme de Chang et Roberts **proposée** dans [KOR 81], l'idée qui a été utilisée sur un anneau bidirectionnel consistait a modifier la phase d'initialisation en choisissant une direction $d \in \{\text{droite, gauche}\}$ avec la probabilité $1/2$. Le message contenant l'identite du processus sera **émis** dans la direction d .

La **complexité** dans le pire des cas est, pratiquement, la même que celle de l'algorithme unidirectionnel, **mais** un message sera **éliminé** plus rapidement s'il se dirige vers un message se **déplaçant** dans la direction opposée. La figure suivante [MAT 88] illustre une telle situation:



Si dans l'anneau unidirectionnel le message **émis** par le processus 3 est **éliminé** par le 7, alors sur l'anneau bidirectionnel il est possible de s'attendre a ce que **les** deux messages se rencontrent a mi-chemin si le 7 decide d'envoyer son message dans la direction **opposée** (cas a). **Donc, avec** la probabilité $1/2$ le message 3 se **déplace** seulement la **moitié** de la distance, ce qui réduit le nombre total de messages moyen. Cependant, si 7 envoie son message dans la **même direction** que 3 (**cas b**), il existe une probabilité supérieure a zéro tel qu'un autre **processus** (9 dans la figure ci-dessus) rencontre le message **émis** par 3 avant d'atteindre le processus 7.

Méthode 2: Il est possible d'obtenir une version **déterministe** en choisissant la direction d **égale** au **sens** des aiguilles d'une montre si l'identite est **paire**, et **égale** au sens contraire si

l'identité est impaire. Cependant, si les IDs des voisins sont connues, d'autres variantes **déterministes** sont possibles. Elles diffèrent par la phase d'initialisation.

En comparant son ID avec celles de ses voisins, un processus initiateur peut agir selon une ou plusieurs stratégies.

La figure suivante énumère quelques unes de ces stratégies.

| I stratégie | 21 | I 1 2 3 | I 1 3 2 |
|-------------|----|---------|---------|
| max | → | → | → |
| min | ← | ← | ← |
| min* | / | / | ← |

En général, ($n \geq 3$) il existe trois cas: le noeud initiateur peut être plus petit que les deux voisins, il peut être plus grand, et il peut être entre les deux. Pour tous les cas, la **stratégie** détermine le voisin auquel le message sera adressé.

Min* est la version proposée dans [BOD 861, elle suppose que tous les processus sont initiateurs. Un processus n'initialise pas l'algorithme et ne propage pas son message si il sait qu'un de ses voisins a une identité supérieure à la sienne.

Méthode 3: Une autre stratégie permettant d'accélérer le processus d'extinction selective de messages(processus) consiste pour un processus donné à adresser directement son IDi au processus dont l'identité IDj est inférieure à IDi. Cette **stratégie** utilise la notion de sens d'orientation global et nécessite la connexion des processus par d'autres liaisons. La topologie en anneau ne suffit pas à l'application d'une telle stratégie. Cette technique est appliquée, par exemple, dans [LOU 861 sur une topologie complète.

Ce qui différencie les algorithmes d'élection d'une topologie à une autre c'est l'algorithme de parcours sous-jacent à l'algorithme d'élection proprement dit. En faisant abstraction du problème de parcours, les mécanismes utilisés pour élire un processus unique sont semblables à ceux appliqués à la topologie en anneau. L'utilisation de la technique d'extinction selective a été, en effet, généralisée aux topologies complète et quelconque. Cela s'est fait de deux manières:

1. La comparaison **porte** directement sur les identites des processus participant a l'election.
2. La comparaison **porte** sur des couples (PH_i, i) où i designe l'identite du processus P_i , et PH_i represente la phase de P_i . Dans ce cas PH_i peut avoir plusieurs definitions:
 - a) PH_i represente la **taille** de l'arborescence en construction, **b)** elle represente le nombre de **racines battues**(puisque plusieurs arborescences peuvent être en construction simultanee), **c)** dans [SKY 871 la phase est definie autrement. Lorsque deux **messages(jetons)** (PH_i, i) et (PH_j, j) relatifs à deux processus initiateurs P_i et P_j , respectivement, se rencontrent au niveau d'un processus P_k , tel que $PH_i = PH_j$ alors, un nouveau **jeton** $(PH_k = PH_i + 1, k)$ est **créé**. Les jetons relatifs à P_i et a P_j ne sont pas transmis(ils sont **éteints**)

2 Election dans les **réseaux** anonymes

Ce type de reseaux est caracterise par l'absence d'**identités**(ou d'identites uniques) des processus. Ils sont bases sur une procedure permettant de choisir une identite. Cela **constitue** la premiere **tâche** que **doit** assumer chaque processus **dès** son activation. Il s'agit de choisir un **entier** S_i au hasard dans l'intervalle d'entiers $[1...d]$, ou un couple **ID** (S_i, T_i) , tel que S_i soit determine comme precedemment. T_i est la valeur d'une **certaine** variable aleatoire, par exemple le nombre de **jetés** d'une piece avant l'apparition de face.

Le but **principale** que l'on essaie d'atteindre c'est de **réaliser** une procedure assurant **le choix** d'identites uniques avec la plus grande **probabilité** possible.

A partir de la, les techniques habituelles d'election sont exploitees pour **élire** un processus. Parfois, au **cours** de l'algorithme, le **choix** d'autres identites peut **s'avérer** necessaire lorsqu'un ou plusieurs processus concurrents ont été dotes de la **même** identite.

Ce type d'algorithmes **réalisant** une election avec une **probabilité** d'erreur **même** infime ne sont pas interessants dans la pratique. L'election est utilisee dans le but de **résoudre** un probleme mettant en **défaut** le fonctionnement d'un systeme, et il ne faudrait pas que la solution **apportée** conduise elle aussi a un autre **problème**.

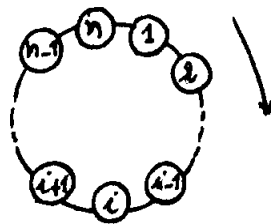
CHAPITRE 4

EVALUATION ET ANALYSE

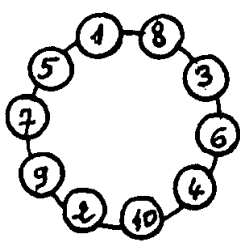
Dans ce chapitre est **présentée** l'analyse des résultats obtenus **par** l'évaluation d'un échantillon d'algorithmes. Les performances en nombre de messages et en temps **d'élection** sont données pour chaque algorithmes. Ce chapitre **comporte** également une étude comparative des différents algorithmes **évalués**.

L'évaluation a **été réalisée** sur les configurations suivantes:

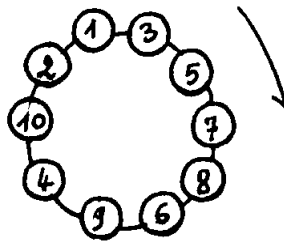
- anneau croissant: les processus sont placés de façon ordonnée selon leurs identités **croissantes**(ordre naturel).



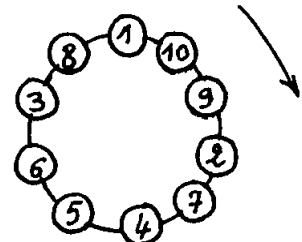
- anneau quelconque: le placement des processus est quelconque.



config. 1

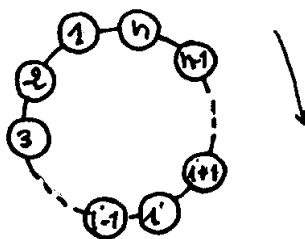


config. 2

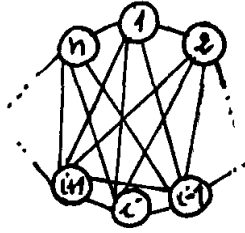


config. 3

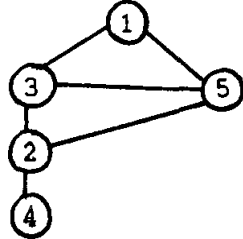
- anneau décroissant: les processus sont placés de façon ordonnée selon leurs identités **décroissantes**.



- reseau complet:



- reseau quelconque:



La representation des resultats a été faite comme suit:

1. L'election collective, en faisant varier la taille n du reseau, **permet** de représenter la variation du nombre de messages ou **le** temps d'election en fonction de n .
2. L'election collective **avec** $n-10$, **permet** de représenter le nombre de messages **par** site, **et/ou** le nombre total de messages ainsi que **le** temps d'election par configuration de reseau.
3. L'election individuelle **avec** $n-10$, **permet** de représenter **le** nombre de messages **par** site, **et/ou** le nombre de messages total et le temps d'election en fonction du processus initiateur.

I EVALUATION DES PERFORMANCES

Les resultats relatifs à chaque algorithme sont **présentés** et analyses dans ce paragraphe.

Les valeurs donnees ici **concernent** le nombre de messages necessaires à la determination de l'elu, **mais** chaque processus **émet** un message de plus (au total $n-1$ messages de **plus**) pour informer son voisin de l'identite de l'elu. De **même**, le temps **donné** est celui **mesuré** à partir du lancement de **l'élection** jusqu'à ce que l'identite de **l'elu** soit connue par un des processus (le premier a conclure).

1 Topologie en anneau

1.1 Algorithmes De Chang Et Roberts [CHA 79]

a) Complexité en messages:

Les bornes minimale (quand Pmax lance l'élection) et maximale prévues par l'analyse théorique, respectivement n [figures IV.6a, IV.6b, IV.5a, IV.5b], et $n(n+1)/2$ (figures IV.4a, IV.4b) sont les mêmes que les complexités obtenues.

| In | Topologie | mode | mode | lmsges | lmsges | I temps | I Figures |
|------|------------------|-------------|-------------|---------------|---------------------------|----------------|-----------|
| | réseau | exéc. | élect | théorique | obtenus | lobtenu | |
| n | Anneau croissant | Asyn. Sync. | col. | $O(n \log n)$ | $\leq 2n-3$ $2n-1$ | $3n-2$ n | IV.1 |
| 10 | même | Asyn. Sync. | lindivl P10 | | 10 10 | 16 10 | IV.6 |
| 10 | même | Asyn. Sync. | lindivl Pi | | 11-19 idem | 17-39 12-20 | IV.5 |
| In | Anneau décrois. | Asyn. Sync. | col. | $n(n+1)/2$ | $n(0.3n+2)$ $n(n+1)/2$ | $2n+9$ n | IV.4 |
| 1101 | même | Asyn. Sync. | lindivl P10 | $O(n \log n)$ | 10 10 | 15 10 | xv.7 |
| 10 | même | Asyn. Sync. | lindivl Pi | | 11-19 idem | 22-30 12-20 | IV.7 |
| 10 | Config 1 | Asyn. Sync. | col. | $O(n \log n)$ | 30 33 | 21 10 | IV.2 |
| | | Asyn. Sync. | lindivl P10 | | 10 10 | 17 10 | IV.6 |
| 10 | Config 2 | Asyn. Sync. | col. | $O(n \log n)$ | 19 22 | 19 10 | IV.2 |
| | | Asyn. Sync. | lindivl P10 | | 10 10 | 18 10 | IV.6 |
| 10 | Config 3 | Asyn. Sync. | col. | $O(n \log n)$ | 30 33 | 25 10 | IV.2 |
| | | Asyn. Sync. | lindivl P10 | | 10 10 | 21 10 | IV.6 |

Table IV.1

Pour une election collective sur un anneau croissant, les figures IV.1a et IV.1b donnent une **complexité** de $(2n-1)$ et de $(2n-3)$ respectivement. Pour n égal a 10 (figures IV.10a, IV.10b) l'election individuelle **génère** entre 10 et 19 messages.

Lorsque **Pmax(P10)** lance l'**élection** (figures IV.6a, IV.6b), la distribution des messages sur les différents processus est uniforme. **Tous** les processus emettent le **même** nombre de **messages** (1 message de type election: au total 10) pour toutes les configurations. Le message tmis par Pmax est le seul qui circulera sur l'anneau.

Comme on peut le **constater** (table VI.11, lorsque l'election est individuelle, des comportements identiques sont obtenus pour les deux modes synchrone et asynchrone, quelque **soit** la configuration et le processus initiateur.

Dans le cas defavorable oh l'election est collective sur un anneau décroissant, le comportement **asynchrone** (figure IV.4a) est **légèrement** meilleur par rapport au comportement **synchrone** (figure IV.4b), les **complexités** obtenues Ctant respectivement $n(0.3n+2)$ et $n(n+1)/2$.

b) Complexite en temps:

$O(n)$ est la **complexité évaluée** (figures IV.1c, IV.1d, IV.4c, IV.4d).

En mode asynchrone le temps enregitre est plus **élevé** qu'en mode synchrone. Dans ce dernier cas, n unites de temps suffisent pour que le processus Pmax **sache** qu'il est **élu**, lorsque l'election est **collective** (figures IV.1d, IV.4d) ou individuelle **initialisée** par Pmax (figures IV.6d). Si par **contre** l'initiateur n'est **pas** Pmax, ce temps varie entre 12 et 20 unites (figures IV.5d, IV.7d). En mode asynchrone, et toujours pour l'election individuelle, il varie entre 22 et 30 si l'anneau est décroissant (figure IV.7c), et entre 17 et 39 s'il est **croissant** (figure IV.5c).

c) Conclusion:

L'élection individuelle est plus performante que **l'élection** collective. Ce **résultat** est logique car dans le cas d'une election individuelle, seul le processus initiateur Cmet un message et a tout instant, l'extinction selective laissera un seul message en transit sur l'anneau. Dans le cas collectif, initialement **tous** les processus **génèrent** un message.

Le comportement asynchrone est meilleur en **élection** collective, **par** rapport au comportement synchrone pour ce qui est du nombre de messages **générés**. Cela peut s'expliquer par le fait qu'en mode asynchrone, les processus evoluent **à** des vitesses variables, et que les plus lents recevront des messages avant d'en emettre. Les informations **transportées** par ces messages eviteront aux processus recepteurs d'en emettre

inutilement.

Par contre, le temps **d'élection** obtenu est meilleur en mode synchrone, et **cela** pour les élections collective et individuelle. En effet, en mode synchrone les processus **évoluent** à la **même** vitesse ce qui **permet** de gagner du temps vu que plusieurs messages sont **traités** en parallèle. Ce **résultat** sera **rencontré**, par la suite, au niveau de **tous** les algorithmes.

1.2 Algorithme De Dolev, Klawe Et Rodeh [DOL 82]

| In | Topologie | mode | lmode | lmsges | lmsges | ltemps | ltemps | Figures |
|-----|-----------|-------|-------|-----------------|---------|-----------|--------|---------|
| I | réseau | exéc. | élect | théorique | obtenus | lobtenu | | |
| In | Anneau | Asyn. | col. | $2n \log n + n$ | $3n$ | $\geq 2n$ | I | IV.1 II |
| | croissant | Sync. | | | $3n$ | $n+1$ | I | IV.3 I |
| In | Anneau | Asyn. | col. | | $3n$ | $< 2n$ | | IV.4 I |
| | décrois. | Sync. | | | $3n$ | $n+1$ | | |
| 10 | Config 1 | Asyn. | col. | I | 50 | 30 | | IV.2 I |
| I I | | Sync. | | | 50 | 16 | | |
| 10 | Config 2 | Asyn. | col. | I | | 36 | | IV.2 I |
| I | | Sync. | | | idem | 20 | | |
| 10 | Config 3 | Asyn. | col. | I | | 33 | | IV.2 I |
| | | Sync. | | | idem | 17 | | |

Table IV.2

a) Complexité en messages:

Le nombre de messages obtenu est en $O(n)$ (figures IV.1a, IV.1b, IV.4a, IV.4b). Que l'anneau soit croissant ou décroissant, $3n$ messages sont **générés** en mode d'exécution synchrone ou asynchrone.

Pour n **fixé** à 10 processus, 30 messages sont obtenus sur un anneau croissant et **sur** un anneau décroissant. 50 messages ont **été générés** sur un anneau où les processus sont placés dans un ordre quelconque (figures IV.2a, IV.2b). Ceci montre que le nombre de messages n'est pas sensible au placement des processus sur l'anneau, lorsque l'ordre est quelconque.

Le nombre de messages **générés** par chaque processus est le **même** (figures IV.3a). Cet algorithme **permet** de **réaliser** un **équilibre** de charge.

b) Complexité en temps:

La complexité obtenue est $O(n)$ unités de temps; $(n+1)$ unités en mode **synchrone** (figures IV.1d, IV.4d) et plus de $2n$ unités en mode **asynchrone** (figures IV.1c, IV.1d). Le temps obtenu pour un anneau croissant est exactement le même que pour un anneau **décroissant** en mode synchrone.

Pour les autres configurations, sont enregistrés 16, 20 et 17 unités de temps en mode synchrone, et 30, 36 et 33 unités de temps en mode asynchrone (figure IV.3b).

c) Conclusion:

Le **comportement** de l'algorithme est le même pour les deux modes d'exécution synchrone et asynchrone du point de vue nombre de messages. C'est-à-dire que le nombre de messages ne dépend pas de la vitesse **d'exécution** des différents processus. Par **contre**, du point de vue temps, le mode synchrone est plus performant, comme c'est le cas pour l'algorithme de Chang [CHA 79] ci-dessus.

1.3 Algorithme De Hirshbeg Et Sinclair [HIR 80]

a) Complexité en messages:

Dans les figures IV.1a et IV.1b, la variation du nombre de messages en fonction du nombre de processus composant le réseau, est une fonction **croissante**. Elle est caractérisée par des augmentations brusques qui se traduisent au niveau des figures IV.1a, IV.1b et IV.1c par une augmentation nette du nombre de messages. Ce **phénomène** peut être expliqué par le fonctionnement de l'algorithme lui-même. En effet, l'algorithme **réitère** sur des chemins de longueur 2^i , $i=0,1,2,\dots,k$, jusqu'à atteindre une longueur 2^k pouvant contenir tout l'anneau.

Il est facile de vérifier que pour des réseaux dont la taille n est comprise entre 2^{k-1} et 2^k , c'est-à-dire $n \in \{2^{k-1}+1, 2^{k-1}+2, \dots, 2^k\}$ les itérations porteront sur des chemins de longueur $1, 2, 4, \dots, 2^k$. Lorsque la taille du réseau prend la valeur $2^k + j$, $j \in \{1, 2, 3, \dots\}$, tel que $2^k + j < 2^{k+1}$, les itérations iront jusqu'à 2^{k+1} . Autrement dit, pour $2^{k-1} < n < 2^k$ la longueur maximale atteinte sera toujours 2^k , et cela se traduit par une demi-droite au niveau des figures IV.1a, IV.1b et IV.1c. Dès que n prend une valeur supérieure à 2^k , l'algorithme itère sur un chemin de longueur maximale égale à 2^{k+1} , et cela se traduit par le passage à la demi-droite supérieure. Les demi-droites sont parallèles; elles ont donc le même coefficient directeur. Les équations qui les caractérisent sont de la forme $an + bi$, où b_i varie d'une demi-droite à l'autre (voir la table IV.3).

Rappelons que les processus sont placés dans l'ordre naturel de leurs **identités**, ce qui pourrait expliquer que la complexité trouvée soit différente de celle qui était estimée par l'analyse **théorique** (où la configuration n'a pas été précisée). D'ailleurs, le nombre de messages est plus élevé

quand les processus sont places dans un ordre quelconque (table IV.3, et les figures A.1a, A.1b).

Lorsque l'élection individuelle est lancée par Pmax, le nombre de messages engendrés est le même pour toutes les configurations de l'anneau (figures IV.6a, IV.6b). De plus, le mode d'exécution n'a pas influence cette complexité.

L'élection individuelle lancée par un processus différent de Pmax est représentée sur les figures IV.10a et IV.10b. La complexité dans ce cas, varie entre 84 et 111 messages, et est pratiquement la même en mode synchrone ou asynchrone.

b) Complexité en temps:

L'évolution du temps d'élection en fonction de la taille du réseau suit le même processus que le nombre de messages. Elle est liée directement au fonctionnement de l'algorithme.

| In | Topologie | mode | mode | lmsges | lmsges | temps | Figures |
|----|-----------|-------|---------|--------------------|-------------------------|--------|---------|
| I | réseau | exéc. | élect | théorique | obtenus | obtenu | |
| | | | | | $4n+4 \sum_{i=1}^n 2^i$ | | |
| In | Anneau | Asyn. | col. | $8n+8n \log n$ | $4n+4 \sum_{i=1}^n 2^i$ | même | IV.1 |
| | croissant | | | $\sim O(n \log n)$ | $6n+4 \sum_{i=1}^n 2^i$ | allure | |
| | | Sync. | I | | | | |
| 10 | même | Asyn. | l'indiv | | 79 | 81 | IV.6 |
| | | Sync. | P10 | | 80 | 40 | |
| 10 | même | Asyn. | l'indiv | | 84-111 | 81-98 | IV.5 |
| | | Sync. | Pi | | idem | 48-57 | |
| 10 | Config 1 | Asyn. | col. | $O(n \log n)$ | 143 | 84 | IV.2 |
| | | Sync. | I | | 116 | 50 | |
| | | Asyn. | l'indiv | | 80 | 84 | IV.6 |
| | | Sync. | P10 | | 80 | 48 | |
| 10 | Config 2 | Asyn. | col. | $O(n \log n)$ | 127 | 84 | IV.2 |
| | | Sync. | I | | 132 | 51 | |
| | | Asyn. | l'indiv | | 79 | 89 | IV.6 |
| | | Sync. | P10 | | 80 | 48 | |
| 10 | Config 3 | Asyn. | col. | $O(n \log n)$ | 128 | 76 | IV.2 |
| | | Sync. | I | | 140 | 49 | |
| | | Asyn. | l'indiv | | 79 | 76 | IV.6 |
| | | Sync. | P10 | | 80 | 48 | |

Table IV.3

Contrairement à ce qui a été observé pour le trafic, le comportement temporel en mode synchrone est meilleur par rapport au comportement en mode **asynchrone** (figures IV.1c et IV.1d). Pour un anneau de 10 processus de configuration quelconque (figure IV.2d) le temps d'élection obtenu varie entre 48 et 51 unités en mode synchrone, et entre 76 et 84 unités en mode asynchrone (figures IV.2c).

Le temps enregistré pour une élection individuelle est donné dans les figures IV.5c, IV.5d, IV.6c, IV.6d ainsi que la table IV.3.

Les figures IV.5c, IV.5d montrent que le temps d'élection ne varie pas **beaucoup** lorsque l'initiateur change; car les messages sont **émis** dans les deux directions, et le comportement est symétrique de part et d'autre de cet initiateur quel qu'il soit.

cl Conclusion:

L'élection individuelle engendre moins de messages que l'élection collective. Le nombre de messages est plus **élevé** en mode synchrone et **cela** pour l'élection collective. Le temps est plus **élevé** en mode asynchrone quel que soit le type d'élection. L'explication de ce **résultat** rejoint celle donnée pour l'algorithme [CHA 791].

1.4 Algorithme-P De Korach, **Rotem** Et Santoro [KOR 81]

L'algorithme **décrit** dans [KOR 81] suppose que tous les processus initialisent simultanément l'algorithme. Sinon, à la réception du premier message, un processus Cmet un message, contenant son **identité**, avant d'analyser le message **reçu**.

Cet algorithme a **été modifié** de façon à permettre à un processus d'analyser le **contenu** du message **reçu** avant de transmettre **son** propre message (la version **modifiée** de l'algorithme a **été** désignée par "Korach-2"). Cette modification a **été** introduite dans le but **d'éviter** l'émission **inutile** d'un message dans le cas où l'identité transportée par le message **reçu** est supérieure à celle du processus récepteur. Autrement dit, un processus **génère** son propre message seulement si son **identité** est supérieure à celle **contenue** dans le message **reçu**.

a) Complexité en messages:

En ce qui **concerne** l'algorithme de Korach [KOR 81], la **complexité** obtenue dans le cas d'une élection collective, est $O(n \log n + n)$. Pour n **égal** à 10, le nombre de messages obtenu est au plus **égal** à 51 pour **tous** les **types** d'élection et de configurations **confondus** (table IV.41). En mode synchrone, le nombre maximum de messages est **égal** à 40.

La comparaison des deux versions de l'algorithme, la version d'origine [KOR 81] et la version modifiée Korach-2 permet de conclure que: pour l'élection collective, les deux algorithmes ont presque les **mêmes complexités** en mode **asynchrone** (figure IV.1a), et exactement la **même** en execution synchrone (figure IV.1b). Ces observations se generalisent au cas d'un anneau quelconque (figures IV.2a, IV.2b, table IV.4).

En election individuelle, les figures IV.5a et IV.5b montrent que la version Korach-2 offre la meilleure **complexité** (cela se verifie également sur les figures IV.6a, IV.6b). Dans le cas des figures IV.5a et IV.5b, le nombre de messages varie entre 23 et 51 en asynchrone, et entre 22 et 40 en synchrone pour l'algorithme [KOR 81]. Il varie entre 11 et 20 et entre 11 et 22 pour Korach-2.

1. Algorithme de [KOR 81]

| In | Topologie | mode | mode | lmsgs | lmsgs | I temps | Figures |
|----|-----------|-----------|-------|-------|-------------------|---------|---------|
| I | II | reseau | exéc. | élect | théorique | obtenus | lobtenu |
| | | | | | | | |
| In | II | Anneau | Asyn. | col. | $n^2/2$ | $2n+2$ | IV.1 |
| | | croissant | Sync. | | $O(n \log n + n)$ | n | |
| 10 | | même | Asyn | indiv | $(3n - H_n)/4$ | 38 | IV.6 |
| | | | Sync. | P10 | | 34 | |
| 10 | | même | Asyn | indiv | | 23-51 | IV.5 |
| | | | Sync. | P1 | | 22-40 | |
| 10 | Config 1 | | Asyn. | col. | $O(n \log n)$ | 25 | IV.2 |
| | | | Sync. | | | 26 | |
| | | | Asyn | indiv | | 20 | IV.6 |
| | | | Sync. | P10 | | 20 | |
| 10 | Config 2 | | Asyn. | col. | $O(n \log n)$ | 31 | IV.2 |
| | | | Sync. | | | 22 | |
| | | | Asyn | indiv | | 34 | IV.6 |
| | | | Sync. | P10 | | 26 | |
| 10 | Config 3 | | Asyn. | col. | $O(n \log n)$ | 27 | IV.2 |
| | | | Sync. | | | 27 | |
| | | | Asyn | indiv | | 19 | IV.5 |
| | | | Sync. | P10 | | 25 | |

Table IV.4

b) Complexité en temps:

L'algorithme de Korach [KOR 811] a une complexité de $2n$ unités en mode asynchrone et environ n en mode synchrone, pour l'élection collective sur un anneau croissant (figures IV.1c, IV.1d). La version Korach-2 présente la même complexité que la version d'origine [KOR 811] en mode synchrone, et une légère différence en mode asynchrone.

Lorsque n est fixé à 10 processus, et en faisant varier le placement des processus sur l'anneau, les résultats obtenus sont:

- Dans l'élection collective en mode asynchrone (figures IV.2c), sont utilisées pour [KOR 811] entre 16 et 20 unités de temps, et pour Korach-2 entre 19 et 23 unités. En mode synchrone 11 unités de temps ont été enregistrées pour les deux algorithmes (figures IV.2d).

2. Algorithme de Korach modifié: Korach-2

| In | Topologie | mode | lmode | lmgcs | ltemps | Figures |
|------|-----------|-------|---------|-------------------|--------|---------|
| I I | réseau | exéc. | élect | obtenu | | |
| In | Anneau | Asyn. | col. | | $2n$ | IV.1 |
| I I | croissant | Sync. | I | $O(n \log n + n)$ | n | |
| I--I | | | | | | |
| 10 | même | Asyn | lindivl | 10 | 14 | IV.6 |
| I I | | Sync. | P10 | 10 | 10 | |
| I--I | | | | | | |
| 1101 | même | Asyn | lindivl | 11-20 | 17-32 | IV.5 |
| I I | | Sync. | Pi | 11-22 | 11-19 | |
| i--I | | | | | | |
| 10 | Config 1 | Asyn. | col. | I 25 I | 19 I | IV.2 I |
| I I | | Sync. | I | I 2 6 | I 11 | I I |
| I I | | | | | | |
| I I | | Asyn | lindivl | 10 | 18 | IV.6 |
| I I | | Sync. | P10 | 10 | 10 | |
| I--I | | | | | | |
| 1101 | Config 2 | Asyn. | col. | 34 | 16 | IV.2 |
| I I | | Sync. | | 22 | 11 | |
| I I | | | | | | |
| I I | | Asyn | lindivl | 10 | 21 | IV.6 |
| I I | | Sync. | P10 | 10 | 10 | |
| I--I | | | | | | |
| 10 | Config 3 | Asyn. | col. | I 27 I | 20 | IV.2 |
| I I | | Sync. | | 27 | 11 | |
| I I | | | | | | |
| I I | | Asyn | lindivl | 10 | 21 | IV.6 |
| I I | | Sync. | P10 | 10 | 10 | |

Table IV.4(suite)

* L' election individuelle initialisee par Pmax, donne pour [KOR 81] entre 21 et 31 unites de temps en mode asynchrone, et entre 18 et 21 unites pour Korach-2(figures IV.5c, IV.6c). En mode synchrone(figures IV.5d, IV.6d) entre 11 et 12 unites sont obtenues pour [KOR 81], et 10 unites pour Korach-2 quelque soit la configuration.

L'election individuelle initialisee par un processus Pi autre que Pmax donne pour Korach-2 de meilleurs résultats que pour [KOR 81] en mode d'execution synchrone et asynchrone(figures IV.5c, IV.5d).

c) Conclusion:

Les deux versions CKOR 811 et Korach-2 ont les mêmes comportements temporels et en nombre de messages quand l'élection est collective. L'algorithme Korach-2 est plus performant en election individuelle. Comme pour les algorithmes precedents, le mode synchrone offre le meilleur temps d'élection.

2 Topologie complete

2.1 Algorithmes De Afek Et Gafni [AFE 851

Les trois algorithmes proposes A, B et C tentent chacun d'améliorer les performances de l' autre. L'algorithme B est conçu pour ameliorer la complexite en messages de A, et l'algorithme C est conçu de façon a preserver la complexite en messages de B et a en ameliorer la complexite temporelle.

2.1.1 L'algorithme A

| In | Topologie | mode | mode | lmsges | lmsges | temps | Figures |
|----|----------------|-------|-------|----------------|----------|---------|---------|
| | réseau | exéc. | elect | theorique | lobtenus | lobtenu | |
| | Reseau complet | Asyn. | col. | $2.77n \log n$ | $3n+1$ | $3n+6$ | IV.8 |
| | | Sync. | | | $3n$ | $2n+1$ | |
| 10 | | Asyn. | Icol. | 92 | 32 | 33 | IV.8 |
| | | Sync. | I | | 31 | 21 | IV.9 |
| 10 | | Asyn. | indiv | | 20 | 36 | IV.10 |
| | | Sync. | P10 I | | 20 | 20 | IV.11 |
| 10 | | Asyn. | indiv | | 20 | 35-43 | IV.10 |
| | | Sync. | Pi | | 20 | 20-20 | |

Table IV.5

a) Complexité en messages:

L'élection collective a une complexité (figures IV.8a, IV.8b) de $3n$ messages en mode synchrone et de $3n+1$ en mode asynchrone. Avec $n=10$, 32 messages sont enregistrés pour le cas synchrone et 31 pour le cas asynchrone.

En faisant varier le processus initiateur, pour une election individuelle, les figures IV.10a et IV.10b donnent le même nombre de messages (20) émis au total dans chaque cas et pour les deux modes d'exécution.

b) Complexité en temps:

Elle est située autour de la valeur de $2n$ pour l'élection collective synchrone (Figure IV.8d) et de $3n$ pour l'élection collective asynchrone (figure IV.8c).

En exécution asynchrone 49 unités suffisent pour la terminaison de l'algorithme dans une election collective (figure IV.9c). Le temps mesuré varie entre 35 et 43 unités pour l'élection individuelle (figure IV.10c). En exécution synchrone 31 unités de temps sont nécessaires pour l'élection collective (figure IV.2c), et, 20 unités pour l'élection individuelle quel que soit le processus initiateur (IV.10d).

2.1.2 L'algorithme B

| In | Topologie | mode | lmode | lmsges | lmsges | temps | Figures |
|------|-----------|-------|--------|------------------|---------|-----------|---------|
| I I | réseau | exéc. | élect | théorique | obtenus | lobtenu | |
| In | Réseau | Asyn. | col. | $2n \log n + 2n$ | $3n+5$ | $2.8n+25$ | IV.8 |
| I | complet | Sync. | | | $3n+4$ | $2n+8$ | |
| I | | | | | | | |
| 1101 | | Asyn. | col. | 86 | 35 | 49 | IV.8 |
| I I | | Sync. | | | 34 | 31 | IV.9 |
| 1--1 | | | | | | | |
| 1101 | | Asyn. | indivl | | 20 | 36 | IV.10 |
| I I | | Sync. | P10 | | 20 | 20 | IV.11 |
| I--I | | | | | | | |
| 1101 | | Asyn. | indivl | | 20 | 35-43 | IV.10 |
| I I | | Sync. | Pi | | 20 | 20-20 | |

Table IV.6

a) Complexité en messages:

La complexité évaluée est $O(n)$ (figures IV.8a et IV.8b), elle est donc inférieure à la complexité théorique de $2n \log n + 2n$. Avec $n=10$, le nombre de messages obtenu est égal à 35 en mode asynchrone, et à 34 en mode synchrone (figure IV.91).

Comme l'algorithme A, l'élection individuelle (figures IV.10a, IV.10b) nécessite 20 messages quelque soit l'initiateur et le mode d'exécution.

b) Complexité en temps:

Le temps est mesuré par $O(n)$ unités. Environ $3n$ unités pour le cas asynchrone et $2n$ unités pour le cas synchrone ont été enregistrées (figures IV.8c et IV.8d) respectivement. La figure IV.8c en mode asynchrone montre que l'algorithme B est plus lent que l'algorithme A pour une élection collective. L'élection individuelle, par contre est la même pour les deux algorithmes A et B (figures IV.10c, IV.10d, IV.11c et IV.11d).

2.1.3 L'algorithme C

a) Complexité en messages:

L'évaluation donne une complexité d'environ $2.4n$ pour l'élection collective (figures IV.8a et IV.8b). Mais, les trois algorithmes A, B et C ont pratiquement les mêmes complexités.

| In | Topologie | mode | mode | lmsges | lmsges | ltemps | Figures |
|-----|-----------|-------|-------|--------------------|---------|-----------|---------|
| I I | reseau | exéc. | élect | théorique | obtenus | lobtenu | I |
| In | Réseau | Asyn. | col. | | $3n$ | $2.4n+32$ | IV.8 |
| I | complet | Sync. | | $\sim O(n \log n)$ | $3n+4$ | $2n+8$ | |
| I-- | | | | | | | |
| 110 | | Asyn. | col. | 141 | 30 | 28 | IV.8 |
| I | | Sync. | | | 34 | 31 | IV.9 |
| -- | | | | | | | |
| 110 | | Asyn. | indiv | | 20 | 36 | IV.10 |
| I | | Sync. | P10 | | 20 | 20 | IV.11 |
| I-- | | | | | | | |
| 110 | | Asyn. | indiv | | 20 | 35-43 | IV.10 |
| I I | | Sync. | Pi | | 20 | 20-20 | |

Table IV.7

Le nombre de messages total générés par l'algorithme dans une élection individuelle initialisée par l'un quelconque des processus du réseau est C_{gal} à 20 (même que A et B), quelque soit le processus initiateur (figures IV.10a et IV.10b).

b) Complexité en temps:

Les figures IV.8c et IV.8d donnent la variation du temps d'élection en fonction du nombre de processus du réseau. Il en ressort que la complexité de l'algorithme C est la même que celle de B, en mode synchrone. En mode asynchrone, la valeur moyenne du temps d'élection de l'algorithme C est légèrement moins élevée que celle de B.

Les figures **IV.10c** et **IV.10d** représentent le temps d'élection enregistré pour des élections individuelles initialisées par les processus du réseau. Elles montrent que dans le cas synchrone 20 unités sont nécessaires quelque soit le processus initiateur. Pour le mode asynchrone le temps varie entre 35 et 43 unités.

c) Conclusion:

Pour les trois algorithmes, la conclusion qui s'impose est que le mode d'exécution n'a pas d'effet sur les **complexités** en **trafic** engendré. Il faut noter que pour une election collective en mode asynchrone, ou synchrone les trois algorithmes A, B et C **génèrent** le **même** nombre de **messages**(figures **IV.8a**, **IV.8b**). L'élection **individuelle**(figures **IV.10a** et **IV.10b**) est la **même** pour les trois algorithmes(20 messages obtenus partout).

D'après ces résultats, le but fixe dans [AFE 851 et qui consiste à améliorer la complexité en messages de B par rapport à celle de A n'est pas atteint.

Comme cela a été prévu, du point de vue temps d'élection, l'algorithme A est le plus performant. Cela est visible dans le cas de l'élection collective **synchrone**(figures **IV.8c** et **IV.8d**). Les algorithmes B et C ont la **même** complexité. L'amélioration de C par rapport à B, n'est **réalisée** (figure **IV.8c**), mais très faiblement, que dans l'**élection** collective asynchrone.

2.2 Algorithme De **Loui**, Hatsushita Et West [LOU 861

a) Complexité en messages:

Le nombre de messages est **égal** à environ $3n$ en mode d'exécution synchrone ou asynchrone pour l'**élection collective**(figures **IV.8a** et **IV.8b**).

Quand P_{max} lance l'**élection**(figure **IV.10a**), le nombre de messages **émis** par chaque processus, sauf P_{max} , est le **même**. En faisant varier le processus **initiateur**(figure **IV.10a**, **IV.10b**) le nombre total de messages **générés** par l'algorithme semble en **être très peu affecté**; la valeur moyenne est $C_{gale} \approx 29$. Cela peut s'expliquer **par** l'utilisation du **sens** d'orientation global qui **permet** un adressage direct des messages.

b) Complexité en temps:

En mode synchrone, le temps d'élection en **fonction** de la taille du **réseau** est **égal** à $n+1$ (figure **IV.8d**). En mode asynchrone ce temps est presque constant, c'est-à-dire qu'il varie autour d'une moyenne de 22 **unités**(figure **IV.8c**).

L'élection individuelle initialisée **par** les différents processus, **aboutit** à un temps variant entre 9 et 12 unités pour le mode synchrone (figure **IV.10d**). **Par contre** en mode asynchrone ces temps varient en **fonction** de l'initiateur entre

| In | Topologie | mode | mode | lmsges | mges | temps | Figures |
|-----|-----------|-------|-------|-----------|---------|---------|---------|
| | reseau | exec. | elect | theorique | obtenus | lobtenu | |
| | | | | (total) | | | |
| In | Réseau | Asyn. | col. | 3.62n | 3n | <=25 | IV.8 |
| | complet | Sync. | | | 3n | n+1 | |
| 10 | | Asyn. | col. | 36 | 28 | 19 | IV.8 |
| | | Sync. | | | 29 | 11 | IV.9 |
| 10 | | Asyn. | indiv | | 32 | 25 | IV.10 |
| | | Sync. | P10 | | 29 | 12 | IV.11 |
| 110 | | Asyn. | indiv | | 27-32 | 15-30 | IV.10 |
| | | Sync. | Pi | | 28-29 | 9-12 | |

Table IV.8

15 et 30 unités (figure IV.10c).

c) Conclusion:

Cet l'algorithme est plus performant, du point de vue temps, que les algorithmes A, B et C précédents (figures IV.8c et IV.8d).

3 Topologie Quelconque

3.1 Algorithme De Helary, Naddi Et Raynal [HEL 86b] (exploration en profondeur)

A/ Topologie en anneau

a) Complexité en messages:

Le nombre de messages obtenu pour une election collective synchrone ou asynchrone, sur un anneau croissant ou décroissant, est de l'ordre de $O(n)$ (figures IV.12a, IV.12b). Par exemple, pour 10 processus 18 messages sont obtenus.

Pour une election individuelle, lorsque le processus initiateur varie, le nombre de messages se situe entre 10 et 17 (figures IV.13a et IV.13b).

Le placement des processus de façon quelconque sur le réseau (figures IV.2, IV.2) fait varier le nombre de messages entre 24 et 33 messages pour l'élection collective.

Dans les deux modes d'élection collective et individuelle, le comportement observe est le même pour une exécution synchrone et asynchrone.

en mode asynchrone et entre 10 et 18 unités en mode synchrone.

En election individuelle (figures IV.5c et IV.5d) le temps varie entre 15 et 29 unités selon l'initiateur en mode asynchrone et, entre 11 et 18 unités en mode synchrone.

Comme pour les autres algorithmes, le comportement synchrone est plus performant que le comportement asynchrone.

B/ Topologie complete

a) Complexite en messages:

L'évaluation a abouti a une complexite de $2n-5$ messages (figures IV.12a, IV.12b). C'est-a-dire $3(n-1)$ messages sont obtenus en comptant les messages de conclusion, contre $(n-1)(n-3)$ prévus par l'analyse theorique dans le cas d'une election collective ($n=10$, 27 messages, contre 631).

Dans le cas d'une election individuelle où Pmax est seul a lancer l'election, on prévoit $2(n-1)$ messages au total, et c'est exactement la valeur qui a été obtenue (figures IV.11a et IV.11b). Quand P10 initialise l'election, 9 messages sont générés (sans compter les messages de conclusion). Si un processus autre que P10 lance l'election, 10 messages sont générés.

Le mode d'execution n'a pas influenceé l'election: les resultats recueillis dans les deux cas sont semblables, cela est vrai pour les elections collective et individuelle.

b) Complexite en temps:

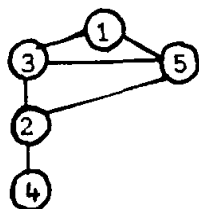
Lorsque l'election est collective la complexite obtenue pour le mode synchrone et le mode asynchrone est $O(n)$ (figures IV.12c et IV.12d). Pour $n=10$, 22 unités sont obtenues en mode asynchrone et 10 en mode synchrone (figure IV.91).

Si l'election est individuelle, le temps mesuré vaut 15 unités (pour la plupart des cas) ou 21 unités selon l'initiateur (figures IV.14c, VI.14d) en mode asynchrone. Il vaut 18 unités en mode synchrone avec Pmax comme seul initiateur (figure IV.11), et 19 unités avec les autres processus initiateurs.

Le mode synchrone donne le meilleur temps.

C/ Topologie quelconque

L'étude a porté sur le reseau suivant:



a) Complexité en messages:

Pour l'élection **collective**(figures IV.15a, IV.15b), le nombre de messages est égal à 11 en exécution synchrone, à 9 messages en exécution asynchrone. Un trafic plus important est observé au niveau des processus P2 et, P3 à cause de leur position **relais**(ils font suivre les messages en provenance des processus P4 et P5.

Si l'élection est individuelle, le nombre de messages varie entre 8 et 9 selon que c'est Pmax ou un autre processus qui lance l'élection(figures IV.16a, IV.16b). On prévoit 12 messages au maximum quand Pmax initialise l'élection, dans [HEL 86bl.

b) Complexité en temps:

Avec cette configuration de cinq pr. sus, les **résultats** sont:

- pour l'élection collective asynchrone 11 unités et pour le cas synchrone 7 **unités**(figur IV.15c) sont **nécessaires**.

- pour l'élection individuelle asynchrone 10 unités sont **enregistrées**(figures IV.16c). En mode synchrone le temps obtenu varie entre 8 et 9 unités, et est **égale** à 7 unités lorsque l'initiateur est Pmax(figures IV.16d).

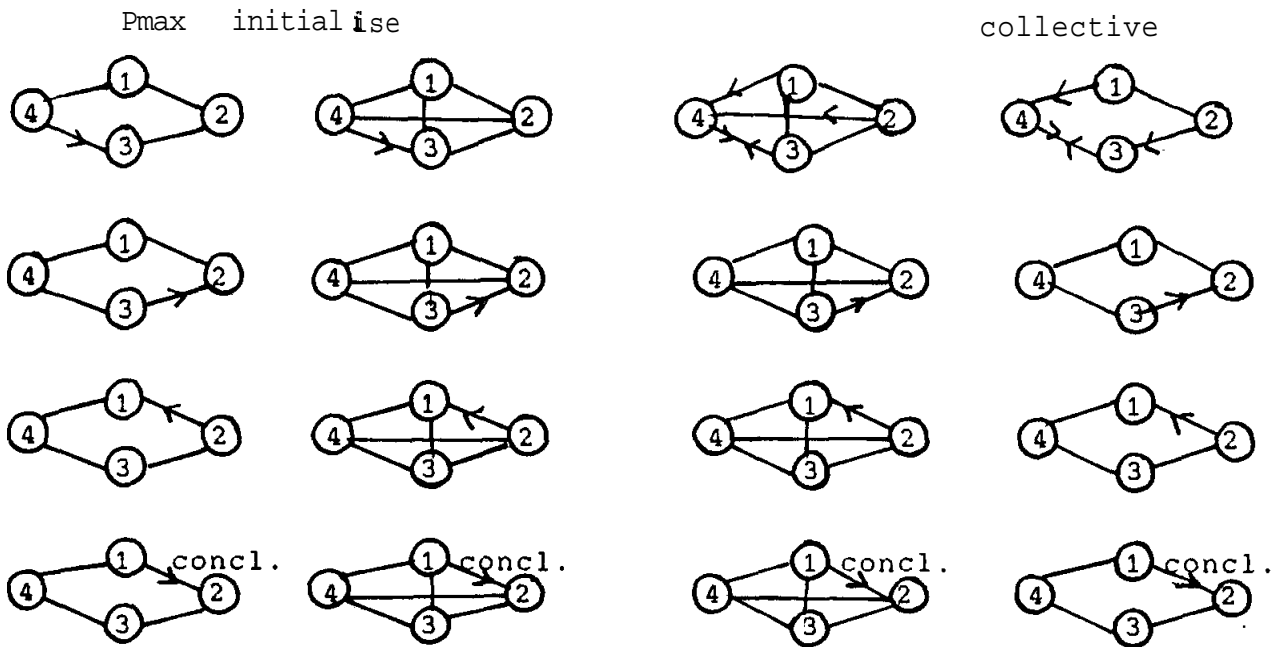
c) Conclusion

En résumé, pour les performances en trafic, le mode asynchrone et le mode synchrone donnent les **mêmes** complexités. **Mais**, en temps, le mode synchrone est mieux que le mode **asynchrone**(l'explication a été donnée précédemment pour l'algorithme [CHA 791).

Pour l'exploration en profondeur, que le **réseau soit** en anneau ou à maillage **complet**, les **complexités** obtenues en **messages**(figures IV.1a, IV.1b) et en **temps**(figures IV.1c, IV.1d) sont les **mêmes** pour les deux modes **d'exécution**, l'élection étant collective. Cette constatation est **également** vraie pour une élection individuelle **initialisée** par Pmax(figures IV.13a, IV.13b, IV.14a, IV.14b)

Ce **phénomène** peut être expliqué en regardant de plus **près** le fonctionnement de l'algorithme. En effet, comme le message est toujours **envoyé** vers le voisin dont l'identité est maximale, il sera toujours **amené à circuler** sur un anneau. ,

Soit l'exemple suivant:



Ces deux exemples montrent bien qu'il y a égalité de performances pour les deux topologies en anneau et a maillage complet.

3.2 Algorithme De Helary, Maddi Et Raynal [HEL 86b]

(exploration en largeur)

A/ Topologie en anneau

a1 **Complexité** en messages:

La **courbe** représentant l'évolution du nombre, de messages en fonction du nombre de **processus** (figures IV.12a; IV.12b) est de la forme $O(n \log n)$, pour les exécutions synchrone et asynchrone. Pour $n=10$, le... nombre de messages est égal a 49 en execution asynchrone et a 56 en execution **synchrone** (figures IV.2a, IV.2b).

Des configurations où le placement des processus est **quelconque** (figures IV.2a, IV.2b) ont été étudiées pour une election collective, et les **complexités** qui en ont, été déduites varient entre 45 et 47 messages en mode asynchrone et, entre 47 et 53 en mode synchrone.

En faisant varier l'initiateur, le nombre de messages varie entre 22 et 53 messages et entre 22 et 55 messages, pour respectivement un mode asynchrone ou synchrone. Ces **résultats** correspondent au cas où les processus sont **placés** dans l'ordre

croissant ou décroissant de leurs identités (figures IV.13a, IV.13b).

| In | Topologie | mode | mode | lmsges | lmsges | ltemps | Figures |
|----|-----------|-------|---------|-----------|-------------------|------------|---------|
| | réseau | exec. | élect | théorique | obtenus | obtenu | |
| n | Anneau | Asyn. | col. | | $O(n \log n + n)$ | $1.7n + 7$ | IV.12 |
| | croissant | Sync. | | | $O(n \log n + n)$ | $1.2n + 5$ | |
| | Idécrois. | | | | | | |
| 10 | | Asyn | lindivl | | 22 | 23 | IV.13 |
| | | Sync. | P10 | | 22 | 13 | |
| 10 | | Asyn | lindivl | | 22-53 | 29-37 | IV.13 |
| | | Sync. | Pi | | 22-55 | 14-18 | |
| 10 | Config 1 | Asyn. | col. | | 45 | 29 | IV.2 |
| | anneau | Sync. | | | 47 | 14 | |
| 10 | Config 2 | Asyn. | col. | | 47 | 22 | IV.2 |
| | anneau | Sync. | | | 53 | 14 | |
| 10 | Config 3 | Asyn. | col. | | 46 | 32 | IV.2 |
| | anneau | Sync. | | | 51 | 15 | |
| n | réseau | Asyn. | col. | $O(n^2)$ | $O(n^2)$ | $2n + 9$ | IV.8 |
| | complet | Sync. | | | $O(n^2)$ | $2n - 1$ | |
| 10 | | Asyn | lindivl | | 18 | 16 | IV.14 |
| | | Sync. | P10 | | 18 | 10 | |
| 10 | | Asyn | lindivl | | 18-80 | 15-28 | IV.14 |
| | | Sync. | Pi | | 18-108 | 11-18 | |
| 5 | réseau | Asyn. | col. | | 11 | 7 | IV.15 |
| | quelconq. | Sync. | | | 17 | 7 | |
| 5 | | Asyn | lindivl | | 8 | 6 | IV.16 |
| | | Sync. | P5 | | 8 | 5 | |
| 15 | | Asyn | lindivl | | 8-14 | 9-11 | IV.16 |
| | | Sync. | Pi | | 8-15 | 6-7 | |

Table IV.10

b) Complexité en temps:

L'évaluation donne une variation linéaire du temps d'élection en fonction du nombre de processus (figures IV.12c, IV.12d). Pour 10 processus, il faut 24 unités de temps pour que l'identité de l'élu soit connue, en exécution asynchrone (figure IV.11c), tandis que 14 unités sont nécessaires en mode

synchrone(figure IV.11d).

Pour l'élection individuelle, en faisant varier l'initiateur, le temps d'élection varie entre 29 et 37 unités, et entre 14 et 18 unités, le mode d'exécution étant respectivement, **asynchrone**(figure IV.13c) puis **synchrone**(IV.13d). A l'inverse de la complexité en messages qui est un peu plus élevée dans le mode synchrone, le temps d'élection est bien plus petit dans le mode synchrone par rapport au mode asynchrone.

B/ Topologie complete

a) Complexité en messages:

La complexité obtenue(figures IV.12a, IV.12b) est de l'ordre de n^2 , pour l'élection collective.

L'élection individuelle a donné 18 messages quand Pmax initialise l'élection(figure IV.111. Le nombre de messages varie entre 18 et 80 quand un processus P_i , autre que P_{10} , initialise l'élection pour le mode asynchrone (figure IV.14a), et entre 18 et 108 en mode **synchrone**(IV.14b).

b) Complexité en temps:

Pour une election collective(figures IV.12c, IV.12d), le temps mesuré est $O(n)$: 26 unités en mode asynchrone et 19 unités pour le cas synchrone dans un réseau de taille $n=10$ (figure IV.9).

Pour une election individuelle, le temps varie en fonction du processus lançant l'élection. Quand ce dernier est Pmax, 16 et 10 unités sont obtenues pour les modes **asynchrone**(figure IV.14c) et **synchrone**(figure IV.14d). Dans les autres cas le temps varie entre 15 et 28 unités, et entre 11 et 18 unités pour respectivement des exécutions **asynchrone**(figure IV.14c) et **synchrone**(figure IV.14d).

C/ Topologie quelconque

a) Complexité en messages:

Les figures IV.15a et IV.15b montrent que le nombre de messages est égal à 11 en mode asynchrone, et à 17 en mode synchrone pour une election collective.

Un trafic plus important est observé au niveau des processus 2, 3 et 5. Ces derniers ont un plus grand nombre de voisins par rapport aux processus 1 et 4. Pour une election individuelle initialisée par Pmax ou par les autres processus a tour de rôle(figures IV.16a, IV.16b), on obtient entre 8 et 14 messages en mode asynchrone, et entre 8 et 15 messages en mode synchrone.

b) Complexite en temps:

Pour **une** election collective 7 unites de temps sont obtenues en mode **asynchrone**(figure IV.15c), et en mode **synchrone**(figure IV.15d).

Pour l'election individuelle, entre 9 et 11 unites ont **été** enregistrees en mode **asynchrone**(figure IV.16c), 6 unites quand Pmax est l'initiateur. Dans le cas synchrone, le temps varie entre 6 et 7 **unités**(figure IV.16d), 5 unites sont **enregistrées** quand Pmax est l'initiateur.

3.3 Comparaison des explorations en profondeur et en largeur

3.3.1 Election collective

a) Comparaison des **trafics**

Reseau en anneau: L'exploration en profondeur donne de meilleurs **résultats** (figures IV.12a, IV.12b) en modes synchrone et asynchrone, par rapport à l'exploration en largeur. **Ceci** est previsible puisque dans le premier algorithme, un processus donne **un message** d'exploration vers un seul de ses **voisins**(celui dont l'identite est **maximale**), tandis que dans le second, il diffuse le **message** d'exploration a l'ensemble de ses voisins. Pour **n=10**, 27 messages sont **émis** par **une** exploration en **profondeur**(**synchrone** et **asynchrone**), 49 **messages**(**asynchrone**) et 56 **messages**(**synchrone**) par **une** exploration en largeur. Ce **résultat** est **vérifié** quelque **soit** la configuration de l'anneau(figures IV.2a, IV.2b).

Reseau complet: La **même** conclusion que le cas d'un reseau en **anneau**(figures IV.12a, IV.12b) peut **être répétée** ici.

Réseau quelconque: C'est aussi l'exploration en profondeur qui est la plus performante (figures IV.15a et IV.15b).

b) Comparaison des temps d'election

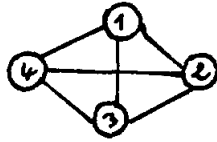
Dans le cas de l'election collective, le temps d'election **obtenu** dans une exploration en profondeur est plus petit que celui de l'exploration en largeur, et **cela** pour les topologies en anneau et à maillage complet, en execution synchrone ou asynchrone. Il faut, cependant, noter que la difference entre les deux explorations est plus importante dans le cas du maillage complet (figures IV.12C et IV.12d). La figure IV.2 **confirme** ce **résultat** sur les differentes configuration de l'anneau.

Pourtant, le lancement simultane de l'exploration **sur** plusieurs chemins, en parallele, pourrait laisser s'attendre à l'acceleration du **calcul** dans l'exploration en largeur **par** rapport à l'exploration en profondeur.

Ce **résultat** peut être expliqué par le fait que la multiplication excessive des messages se traduit par le temps pris dans le traitement de **certains** d'entre eux **véhiculant** des informations sans **intérêt** et ne permettant pas la conclusion.

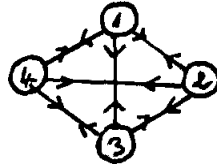
L'exemple de **l'élection** collective ci-dessous illustre cette situation:

Soit le graphe complet:

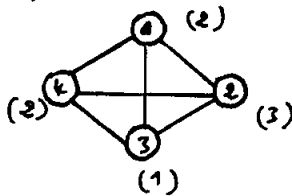


L'exploration en largeur donne les **étapes** de traitement suivantes:

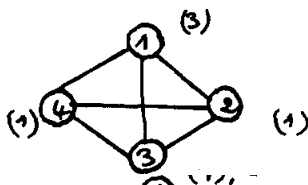
Etape 1: lancement de l'exploration: diffusion a l'ensemble des **voisins**.



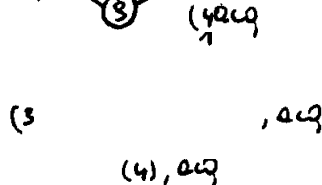
Etape 2: le processus 4 **reçoit** les messages (1), (2) ou (3). Les processus 1, 2 et 3 reçoivent des messages, autres que le message (4), qu'ils acquittent ou éteignent.



Etape 3: comme l'étape 2.



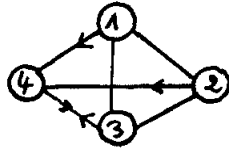
Etape 4: les processus 1, 2, 3 reçoivent chacun le message (4) qu'ils acquittent.



Etape 5: le processus 4 **reçoit** les acquittements et conclut.

L'exploration en profondeur donne les étapes suivantes:

Etape 1: chaque processus lance une exploration vers son plus grand voisin.



Etape 2: chaque processus reçoit le message qui lui est adressé: le processus 3 reçoit (4) qu'il transmet vers 2, le processus 4 reçoit (1), (2) ou (3) qu'il éteint.

Etape 3: le processus 2 reçoit (4) qu'il transmet vers 1.

Etape 4: le processus 1 reçoit les messages (4) et conclut.

Alors que dans l'exploration en profondeur le message émis par le processus P4 est traité par le processus P3 dès l'étape 2, dans l'exploration en largeur il n'est reçu par les destinataires qu'à l'étape 4.

Contrairement aux topologies complètes et à maillage complet, le cas quelconque étudié (figure IV.15c) montre que les deux explorations donnent le même temps d'élection pour le cas synchrone, cependant, c'est l'exploration en largeur qui offre le meilleur temps en mode asynchrone.

3.3.2 Election individuelle

a) Comparaison des trafics

Réseau en anneau: Les figures IV.13a et IV.13b indiquent que l'exploration en profondeur nécessite moins de messages par rapport à l'exploration en largeur.

Réseau complet: Même remarques que le cas de l'anneau (figures IV.14c et IV.14d).

Réseau quelconque: Les figures IV.16a, IV.16b, montrent également que l'exploration en profondeur est plus performante que l'exploration en largeur.

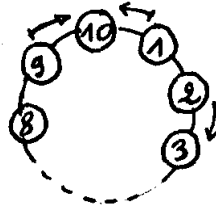
b) Comparaison des temps d'élection

Dans le cas de l'élection individuelle sur une topologie complète (figures IV.14c et IV.14d), l'écart (en temps) entre les deux explorations est d'autant plus réduit que le processus initiateur a une identité qui s'approche de la valeur Max(10). Les cas extrêmes sont:

a) P1 lance l'exploration en diffusant son message (1). A la reception du message (1), chacun des processus diffuse son propre message (P1 ayant l'identite minimale).

b) P10 lance l'exploration, en diffusant le messages (10). A la reception de ce message, **tous** les processus l'acquittent immediatement **sans** generer d'autre messages d'exploration.

Sur une topologie en anneau (figure IV.13c et IV.13d) c'est aussi l'exploration en profondeur qui est plus performante et plus particulierement en mode asynchrone. Dans ce cas, il faut remarquer **que** (surtout en mode synchrone) la difference est moins **élevée** lorsque l'initiateur est plus **près** de P10. Cela est **dû** au fait que l'exploration en profondeur est d'autant plus lente que le processus initiateur s'eloignent de P10. Le meilleur temps est **réalisé** quand P10, P9 ou P1 sont initiateurs.



La topologie quelconque **étudiée** (figure IV.16c et IV.16d) met en evidence la rapidite de l'exploration en profondeur en mode asynchrone, et la rapidite de l'exploration en largeur en mode synchrone.

4.3 Conclusion

La **multiplication excessive des messages** dans le **but d'accélérer la conclusion**, peut avoir un effet inverse à cause du temps pris dans le traitement de messages inutiles. Cela est **confirmé** par le fait que c'est en election collective, et plus particulierement dans la topologie à maillage **complet**, que l'exploration en profondeur est nettement plus rapide que l'exploration en largeur. Dans cette topologie la multiplication des messages est **maximale** du fait que chaque processus est **relié** à **tous** les autres processus du reseau. Par **contre**, la configuration quelconque **étudiée**, a montre que l'exploration en largeur a une meilleure performance temporelle que l'exploration en profondeur (figures IV.15C, IV.16d).

II ETUDE COMPARATIVE GENERALE

Une etude comparative des algorithmes évalués est présentée dans le present paragraphe. Une classification des algorithmes selon leurs performances est proposée.

1 Topologie En Anneau

Les algorithmes conçus pour la topologie en anneau [CHA 79, DOL 82, HIR 80, KOR 81, Korach-2] sont compares avec les algorithmes d'exploration en largeur et profondeur [HEL 86b] exécutés sur la même topologie.

1.1 Election collective

a) Comparaison des trafics

La comparaison des algorithmes pour une election collective en mode asynchrone est representee dans la figure IV.1a, et en mode synchrone dans la figure IV.1b. L'algorithme base sur l'exploration en profondeur a le même comportement que l'algorithme de [CHA 79] conçu pour un anneau unidirectionnel. Quant a l'algorithme base sur une exploration en largeur il a une complexite plus élevée que celles des autres algorithmes a part celui l'algorithme de [HIR 81]. L'exploration en largeur a une complexite qui tend a dépasser celle de [HIR 80] au delà d'une taille de réseau egale a 40.

L' algorithme de [DOL 82] occupe la deuxieme place après les algorithmes de Chang [CHA 79] et l'exploration en profondeur. Les algorithmes de Korach et Korach modifié viennent en troisieme position.

Sur un anneau décroissant (figures IV.4a, IV.4b), l'exploration en largeur est mieux que l'algorithme de Chang [CHA 79]. L'exploration en profondeur est la plus performante. Elle est suivie de l'algorithme de Dolev [DOL 82].

b) Comparaison des temps d'lection

Cette comparaison est representee par la figure IV.1c pour le mode asynchrone, et la figure IV.1d pour le mode synchrone. L'algorithme de [HIR 80] donne la plus mauvaise performance dans les deux modes d'exécution, synchrone et asynchrone. L'exploration en largeur vient en deuxième place. Quant aux algorithmes d'exploration en profondeur et les algorithmes de Chang, Dolev et Korach [CHA 79, DOL 82, KOR 81, Korach-2], ils occupent la première place avec des complexités presque égales surtout en mode synchrone.

1.2 Election individuelle

a) Comparaison des trafics

La figure **IV.5a** montre les elections individuelles en mode asynchrone. L'exploration en profondeur ainsi **que** les algorithmes de Chang [**CHA** 791, de Korach modifie ont des complexites qui ne different pas **beaucoup**. La version initiale de l'algorithme de Korach [**KOR** 811 vient en deuxieme place. L'exploration en largeur a **une** complexite moins bonne, tandis que l'algorithme de Hirshberg [**HIR** 801 presente la complexite la plus mauvaise. La figure **IV.5b** montre le cas **synchrone** et les **mêmes** observations peuvent **être** faites.

b) Comparaison des temps d'élection

En mode **synchrone**(figure **IV.5d**), l'algorithme de Korach modifie est le plus performant, suivi de l'algorithme d'exploration en profondeur et de celui de Chang[**CHA** 791. **Ensuite** l'algorithme de [**KOR** 81] et l'exploration en **largeur** sont situes en troisieme place. **Enfin**, et en derniere place c'est encore l'algorithme de [**HIR** 80] qui presente la plus mauvaise performance.

Pour le mode **asynchrone**(figure **IV.5c**), l'algorithme de Korach modifie l'exploration en profondeur, et l'algorithme de **chang occupent** la premiere place. Il vient **ensuite** l'algorithme de Korach [**KOR** 811 ainsi **que** l'exploration en largeur. L'algorithme [**HIR** 801 **occupe** la derniere place **avec** un retard net par rapport **aux** algorithmes precedents.

2 Topologie complete

Les algorithmes compares pour cette topologie sont les algorithmes A, B et C de Y.Afek [**AFE** 851, l'algorithme de **Loui** [**LOU** 861, ainsi que les algorithmes de Helary [**HEL** 86b].

2.1 Election collective

a) Comparaison des trafics

Les figures **IV.8a** et **IV.8b** indiquent clairement **que** l'exploration en profondeur a une complexite en **messages** moins **élevée** que **celle** des algorithmes A, B, C [**AFE** 851 et celui de [**LOU** 861 qui ont les **mêmes** complexites. L'exploration en largeur donne la complexite la plus **élevée**.

b) Comparaison des temps d'élection

En mode **synchrone**(figure **IV.8d**), ce sont les algorithmes de **Loui** [**LOU** 861 et l'exploration en profondeur qui **offrent** le meilleur temps. En deuxi^eme place vient l'exploration en largeur puis l'algorithme A. Les algorithmes B et C sont les plus lents.

En mode **asynchrone**(figure **IV.8c**), c'est l'algorithme de **Loui** qui

est le plus performant. L'exploration en profondeur vient en deuxième place. L'exploration en largeur **occupe** la troisième position et la quatrième revient à l'algorithme A [AFE 85] suivi directement de l'algorithme C, puis B.

2.2 Election individuelle

a) Comparaison des **trafics**

Sur les figures **IV.10a**, **IV.10b**, il apparaît **que** l'exploration en profondeur est la plus performante en ce qui **concerne** la complexité en **messages**. **Après** cet algorithme, on **trouve** les trois algorithmes proposés dans [AFE 85], A, B, C. L'algorithme de [LOU 86] est **classé** en troisième place avant l'exploration en largeur qui donne la plus mauvaise complexité.

b) Comparaison des temps **d'élection**

En tenant **compte** des deux **diagrammes**(figures **IV.10c**, **IV.10d**) il est possible de **conclure** que l'exploration en profondeur est la plus rapide. En deuxième place vient l'algorithme de **Loui** suivi de l'exploration en largeur. La dernière place revient **aux** trois algorithmes A, B, C.

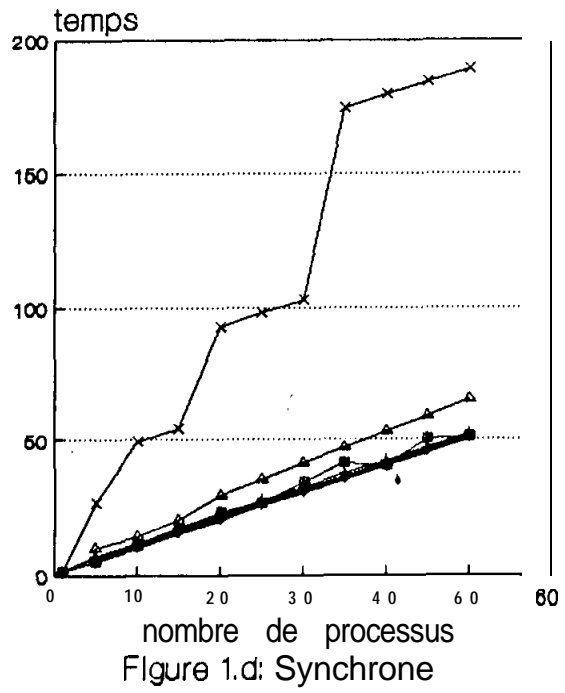
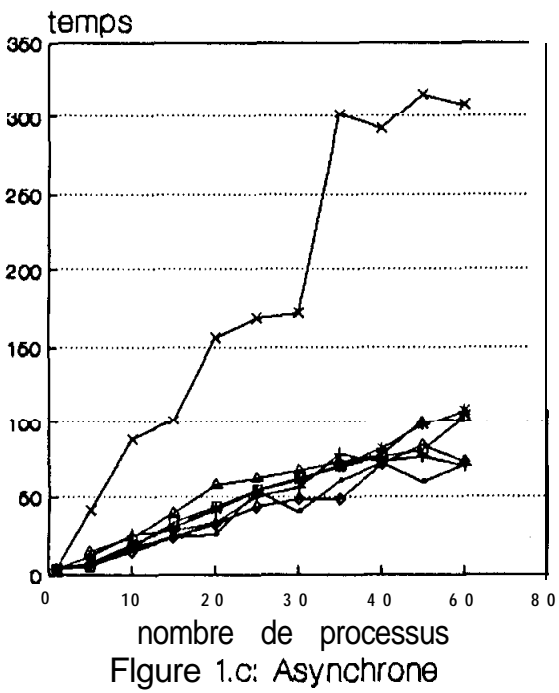
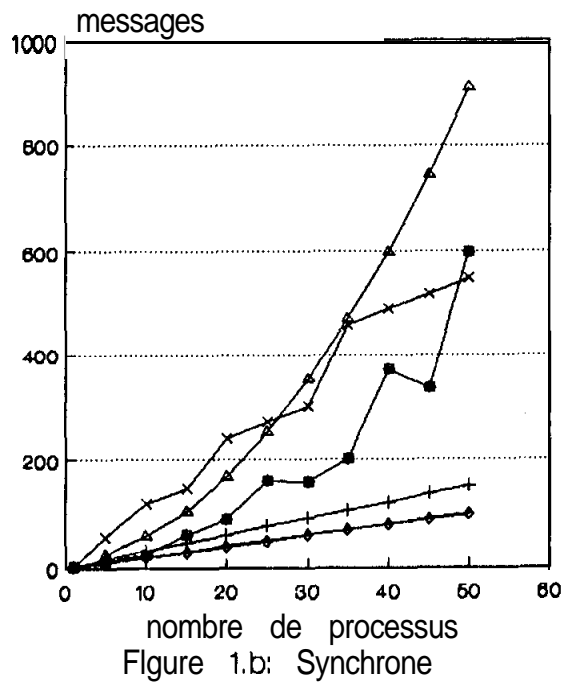
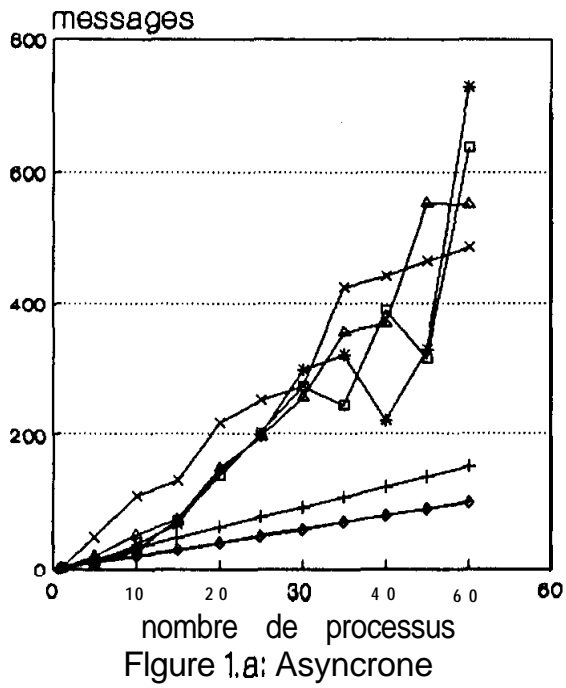


Figure IV.1: Collective (croissant)

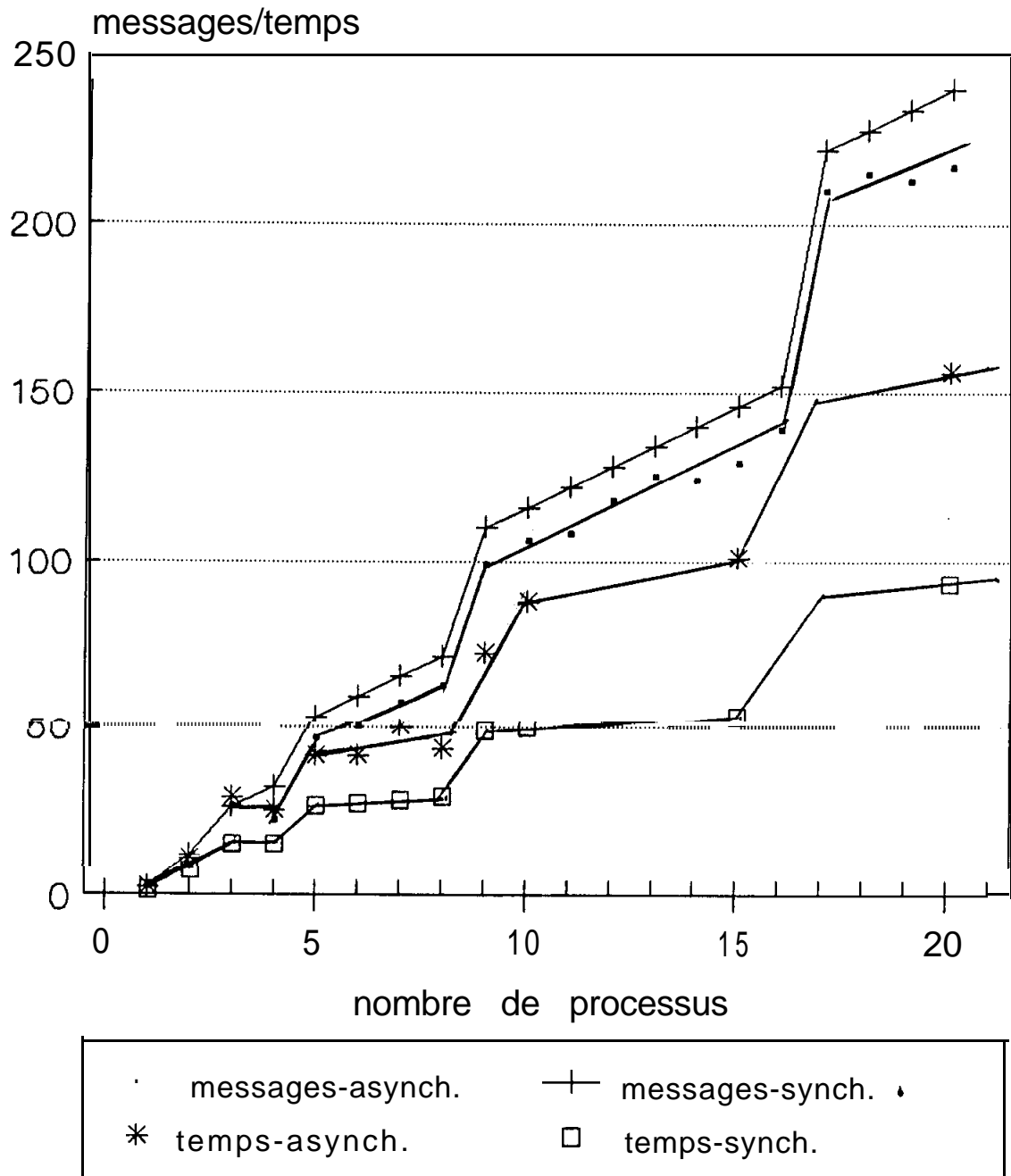


Figure 1.e: Algorithme de Hirshberg

Figure IV.1: Collective (croissant)

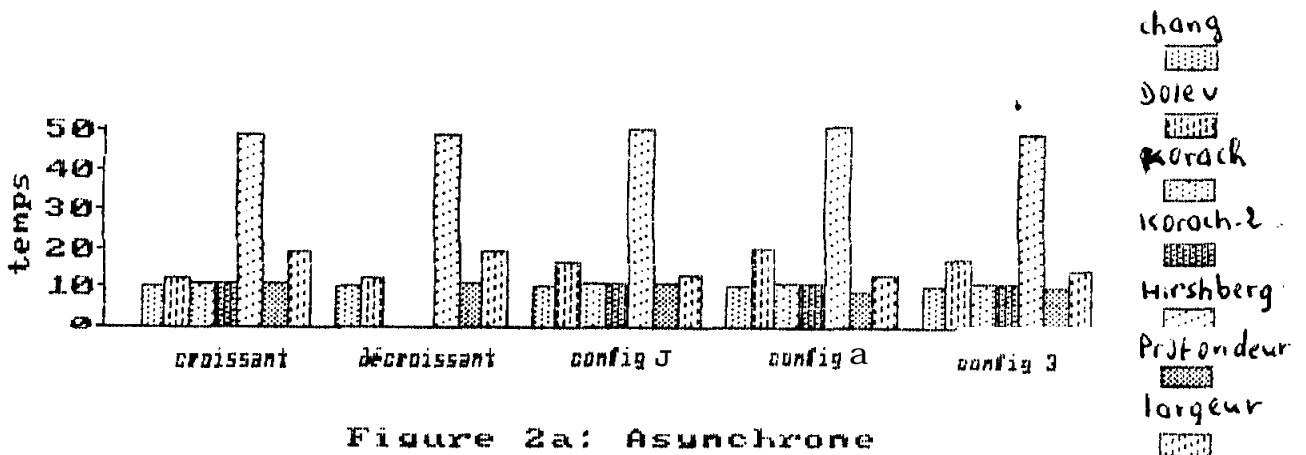
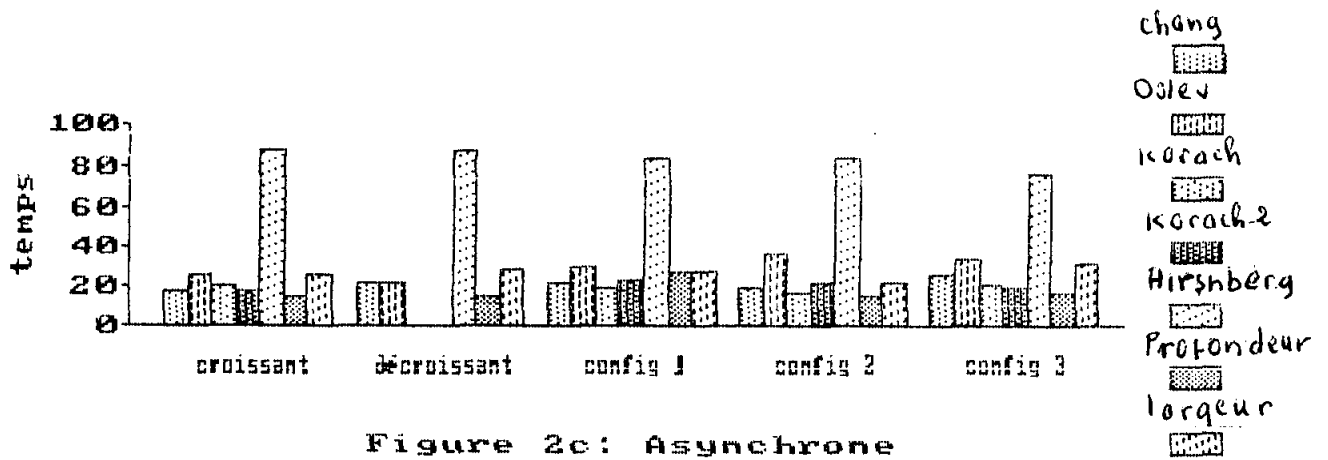
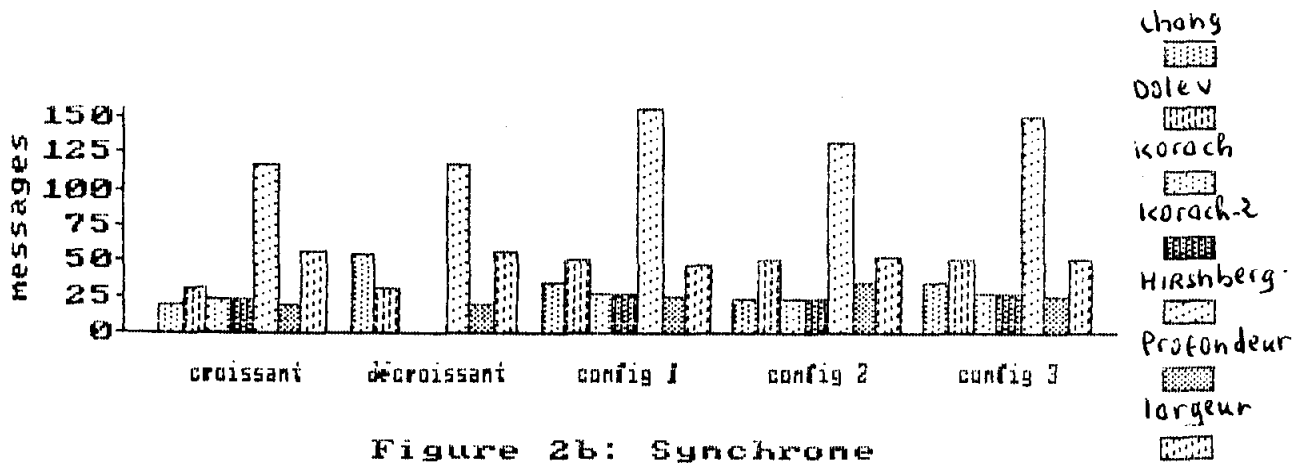
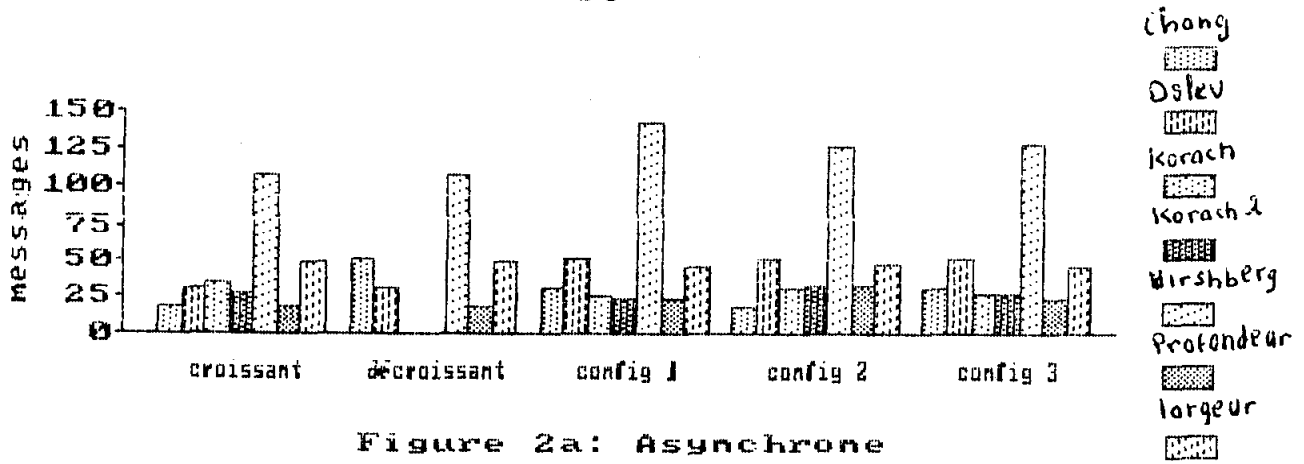


Figure IV.2: Collective (anneau)

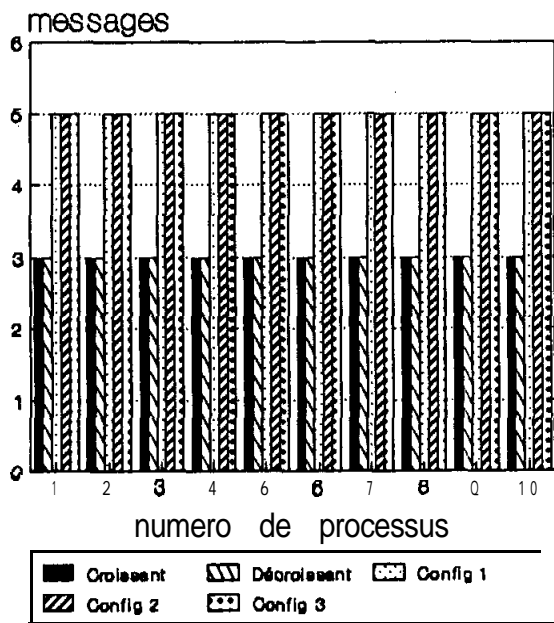


Figure 3a: Asynchrone et Synchronne

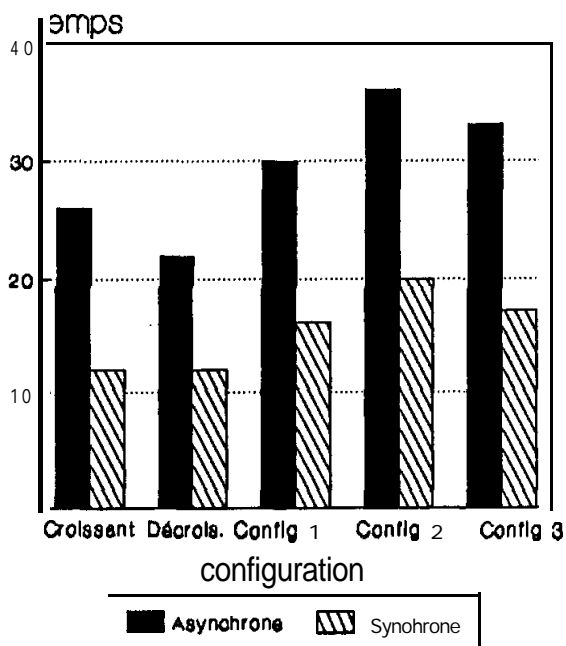


Figure 3.b: Asynchrone / Synchronne

Figure IV.3: Algorithme de Dolev

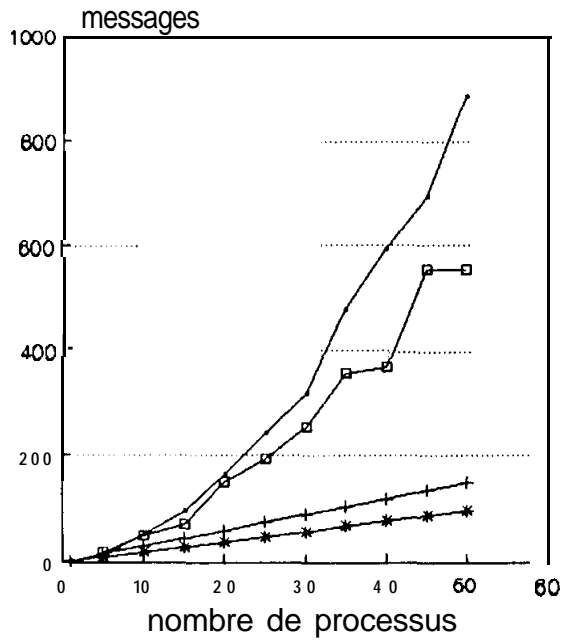


Figure 4.a: Asynchrone

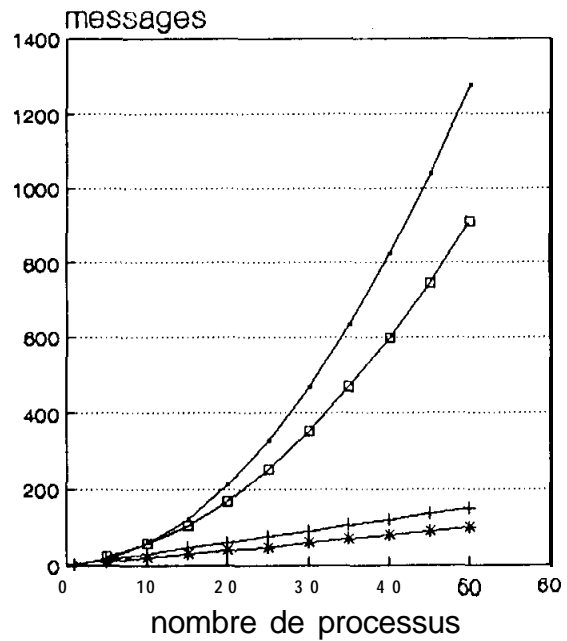


Figure 4.b: Synchrone

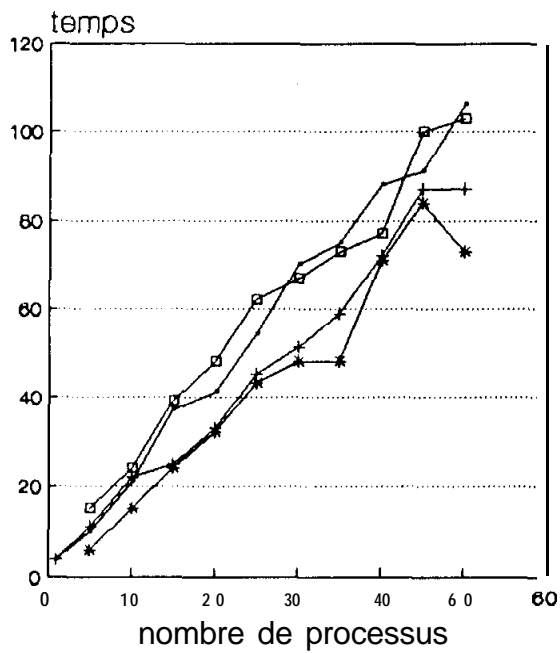


Figure 4.c: Asynchrone

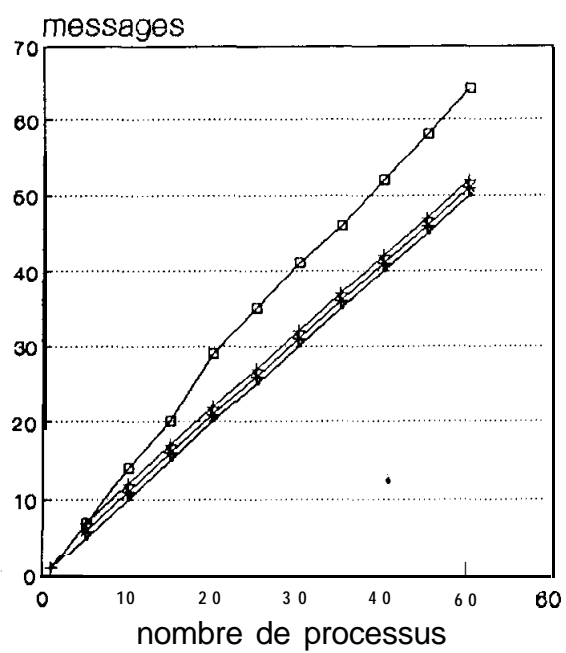


Figure 4.d: Synchrone

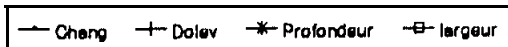


Figure IV.4: Collective (&croissant)

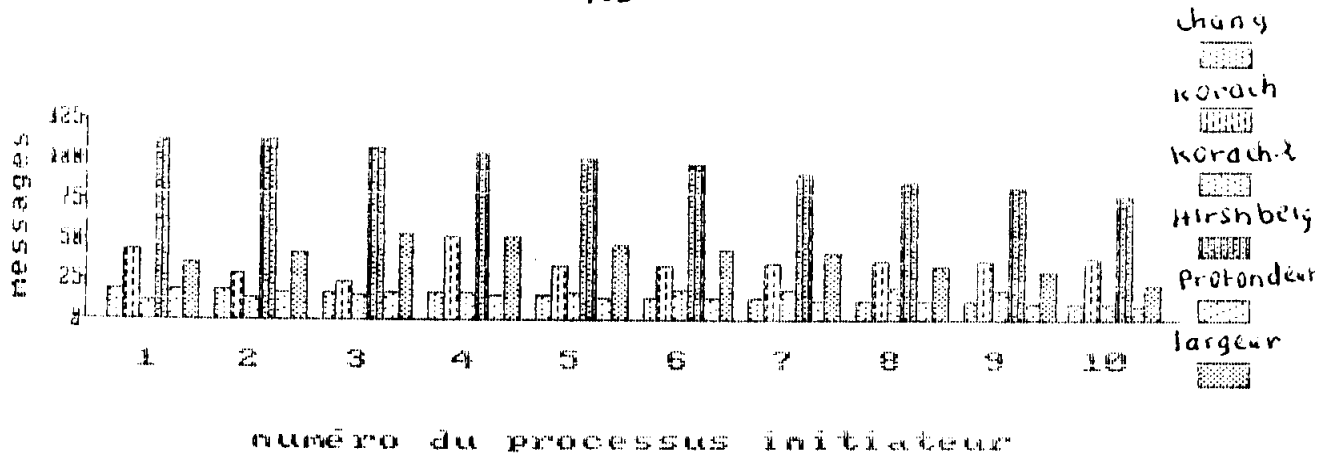


Figure 5a: Asynchrone

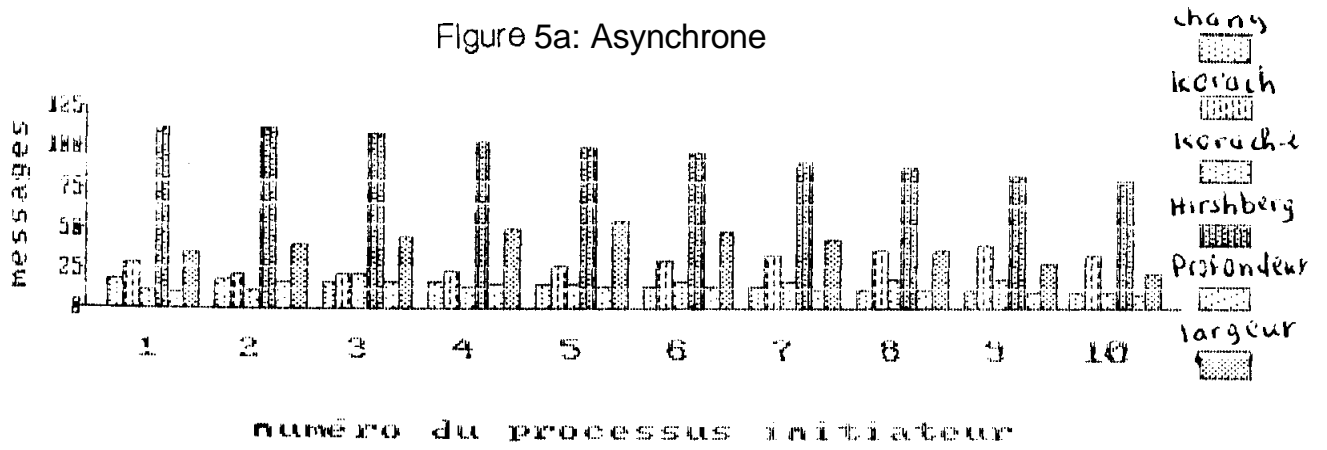


Figure 5b: Synchrone

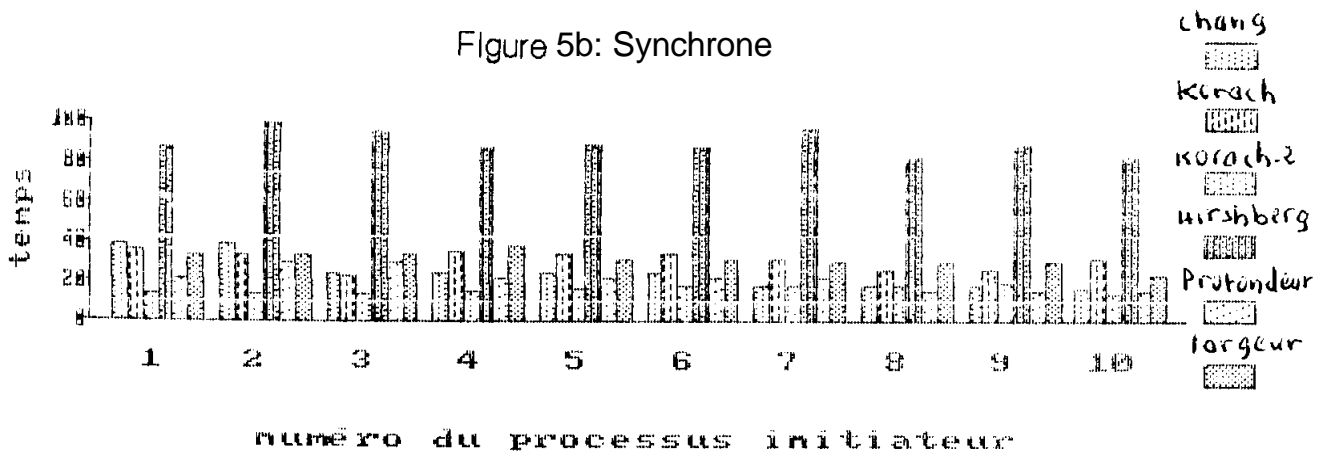


Figure 5c: Asynchrone

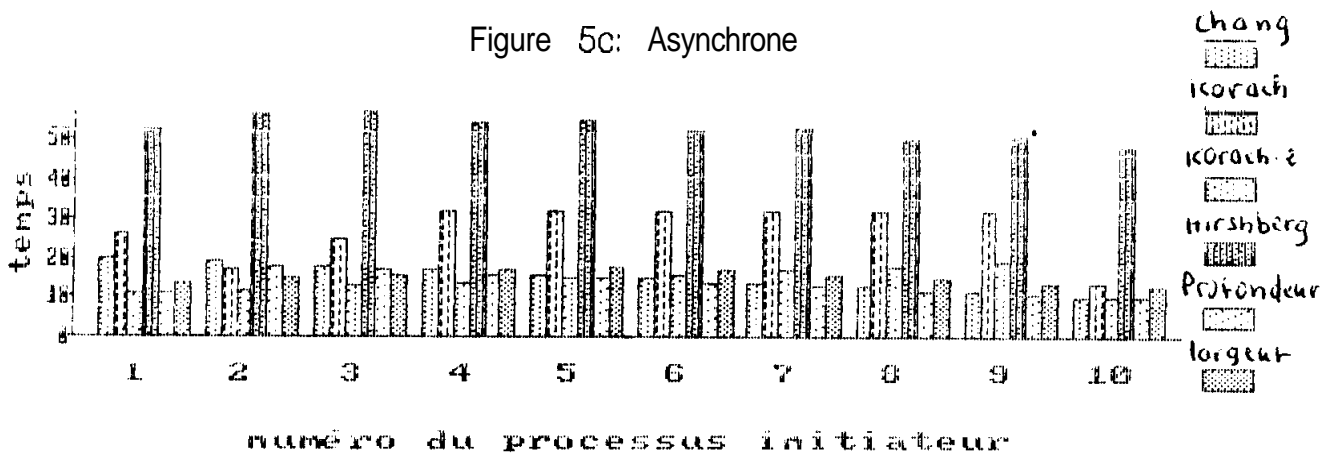


Figure 5d: Synchrone

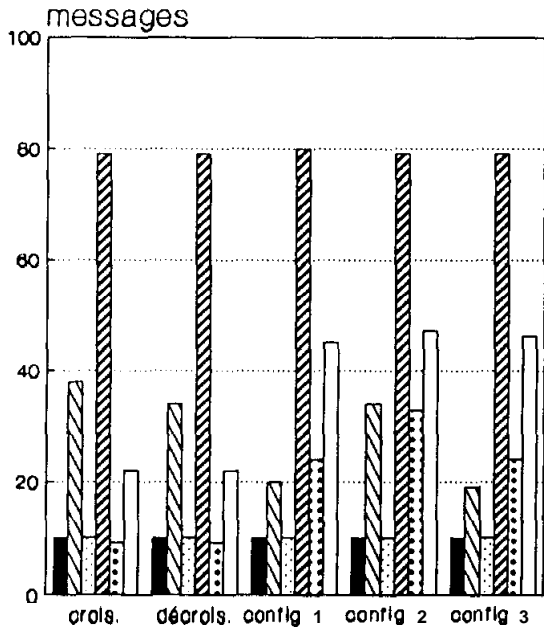


Figure 6.a: Asynchrone

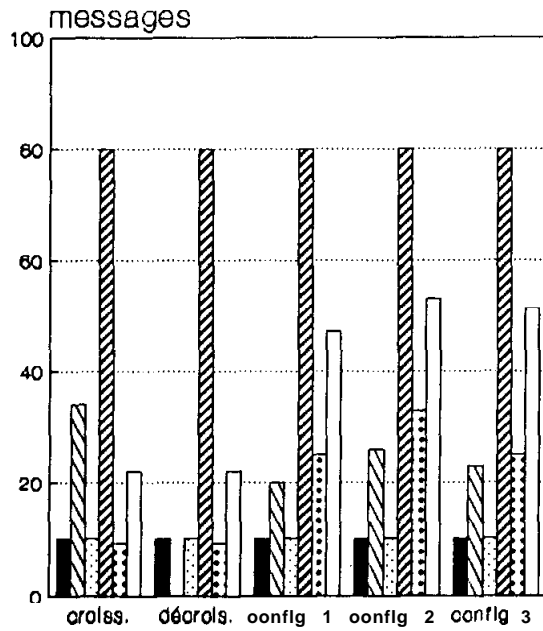


Figure 6.b: Synchrone

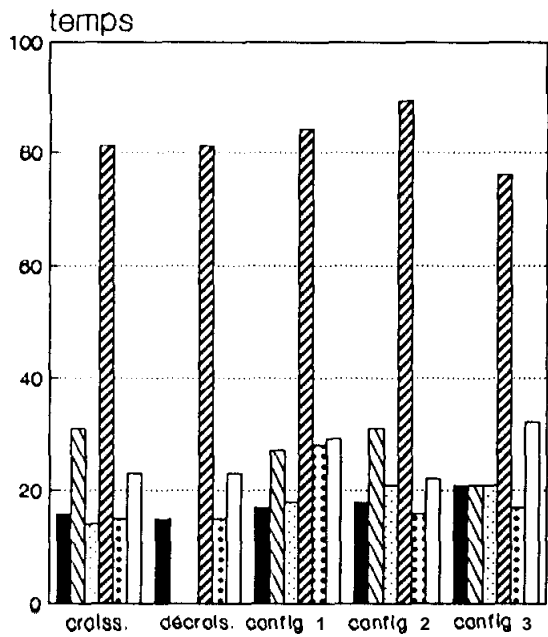


Figure 6.c: Asynchrone

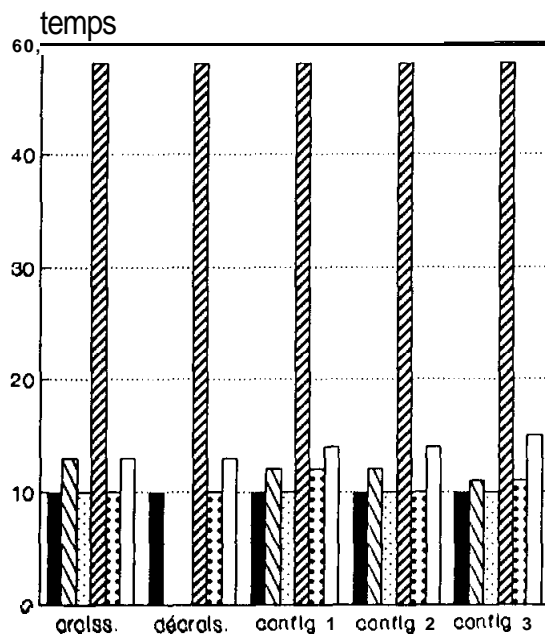
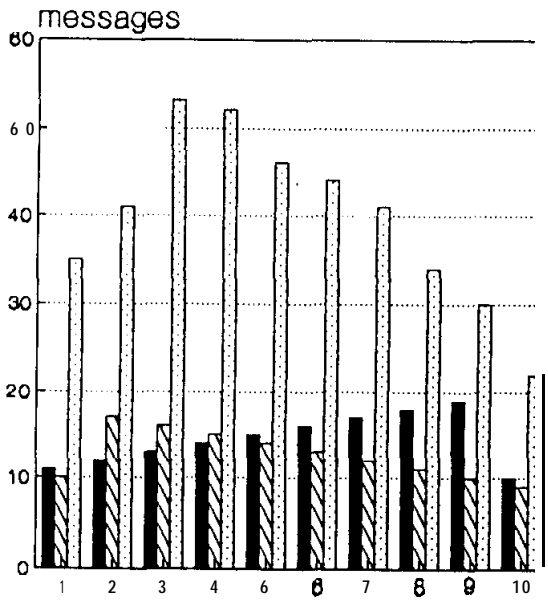
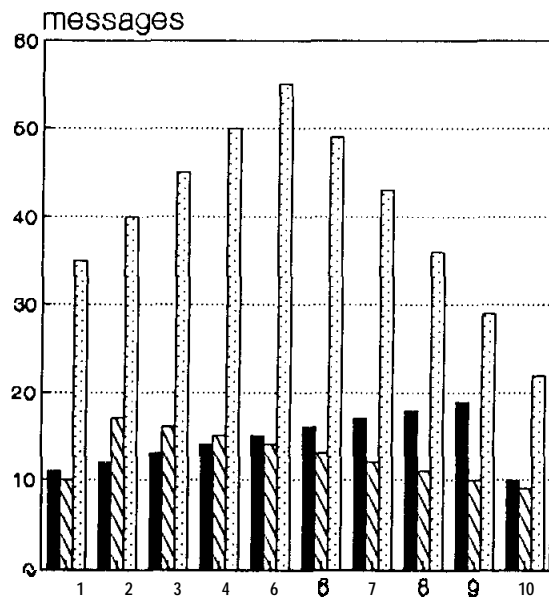


Figure 6.d: Synchrone

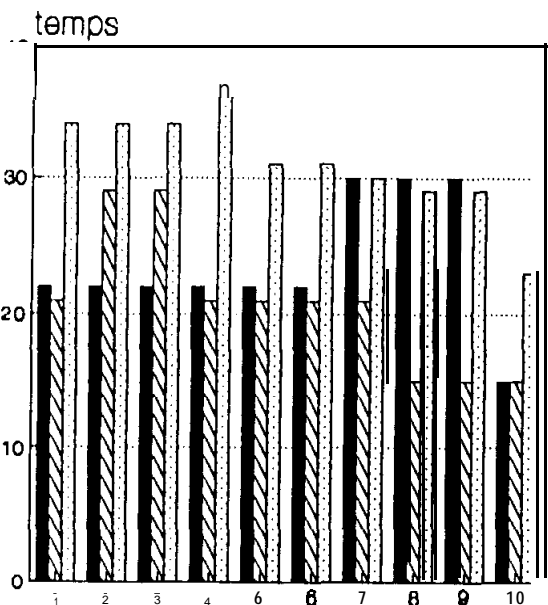
Figure IV.6: Individuelle (Pmax Mt.)



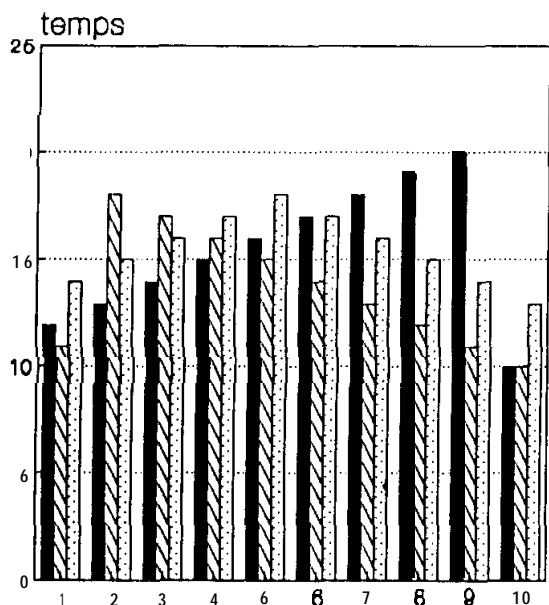
numéro du processus initiateur
Figure 7.a: Asynchrone



numéro du processus initiateur
Figure 7.b: Synchrone



numéro du processus initiateur
Figure 7.c: Asynchrone



numéro du processus initiateur
Figure 7.d: Synchrone

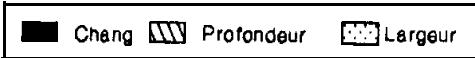


Figure IV.7: Individuelle (décroissant)

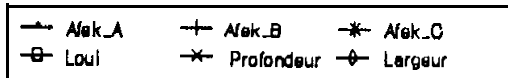
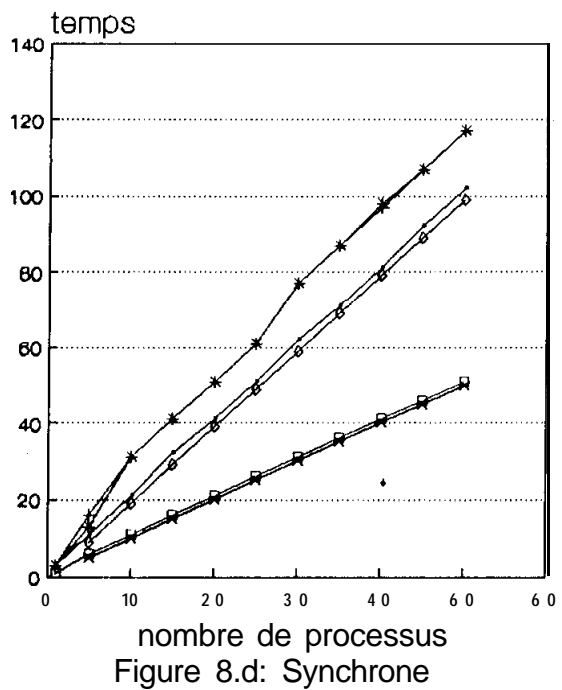
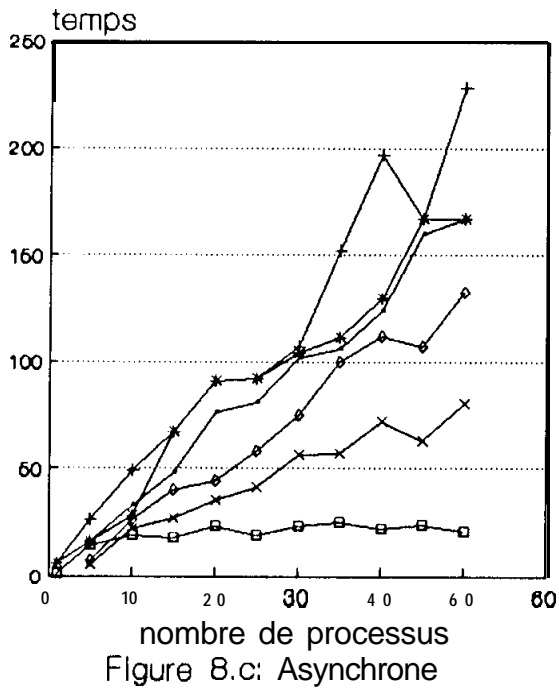
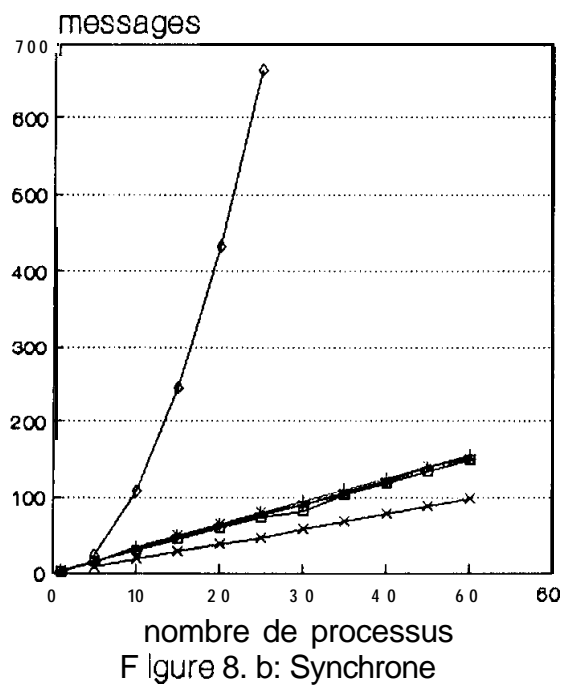
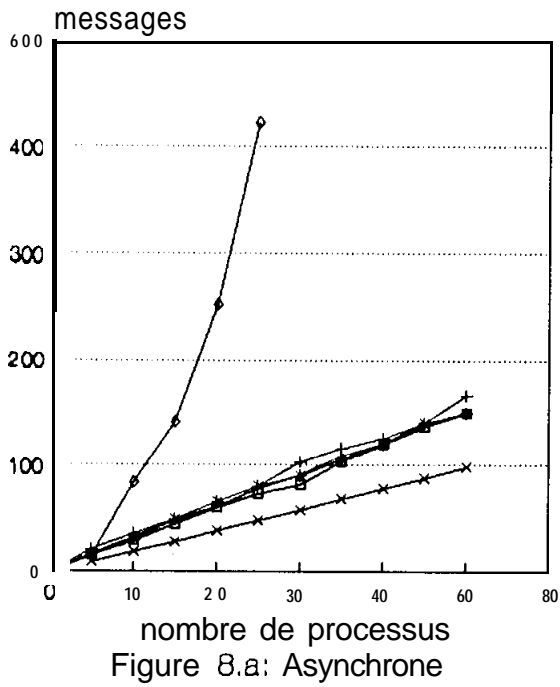


Figure IV.8: Collective (complet)

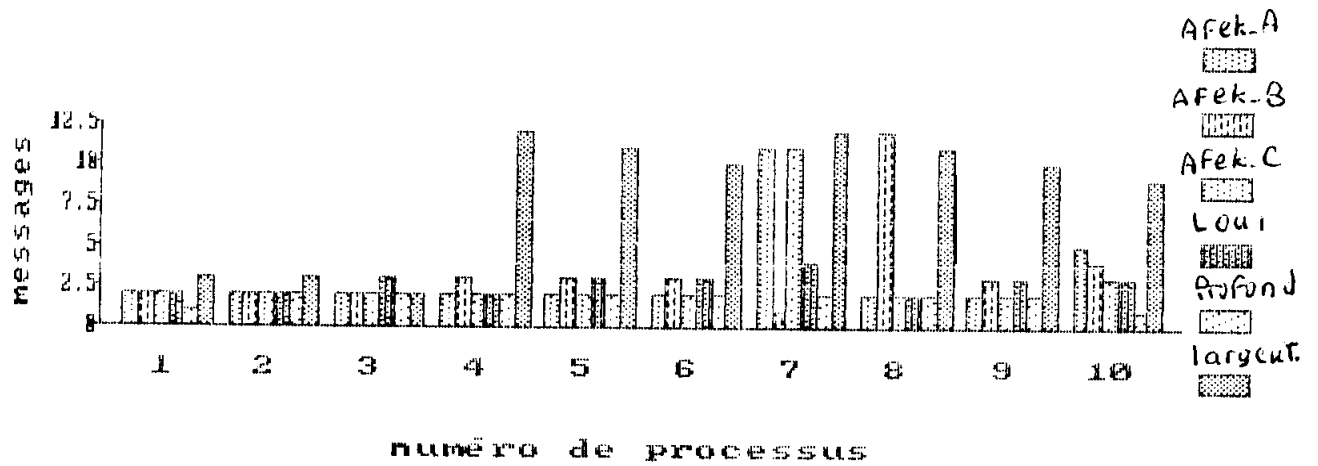


Figure 9a: Asynchrone

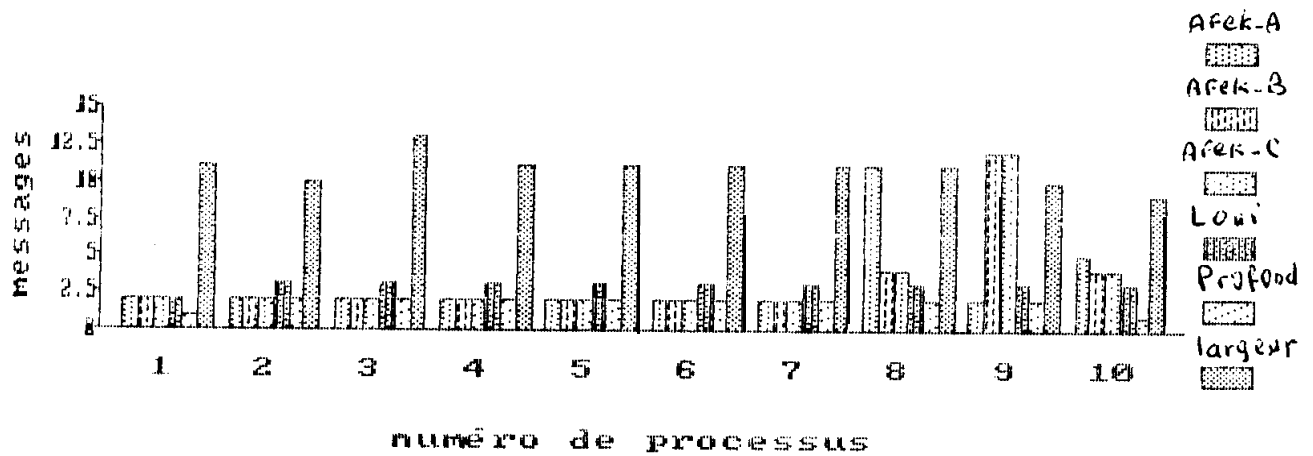


Figure 9b: Synchrone

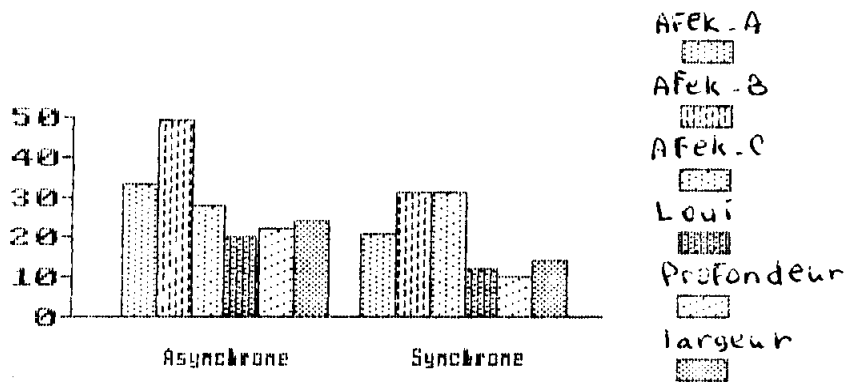


Figure 9c: Temps

Figure IV.9: Collective (complet)

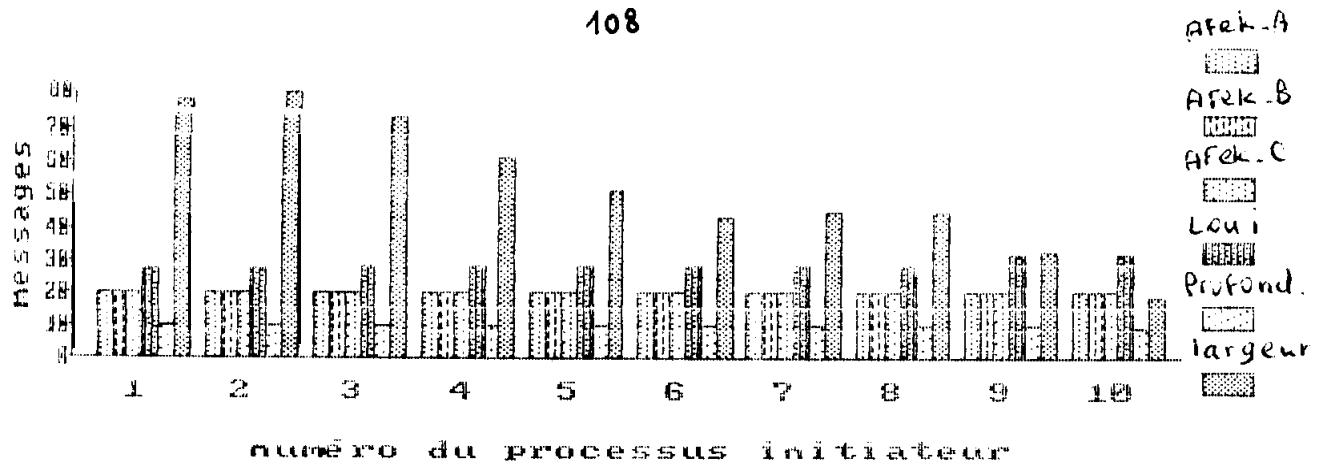


Figure 10a: Asynchrone

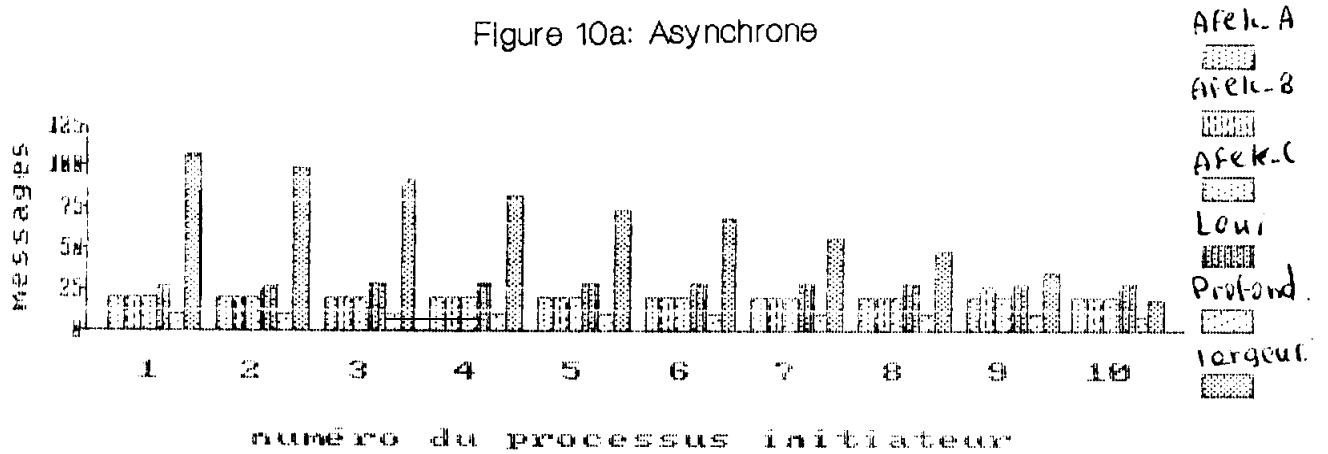


Figure 10b: Synchrone

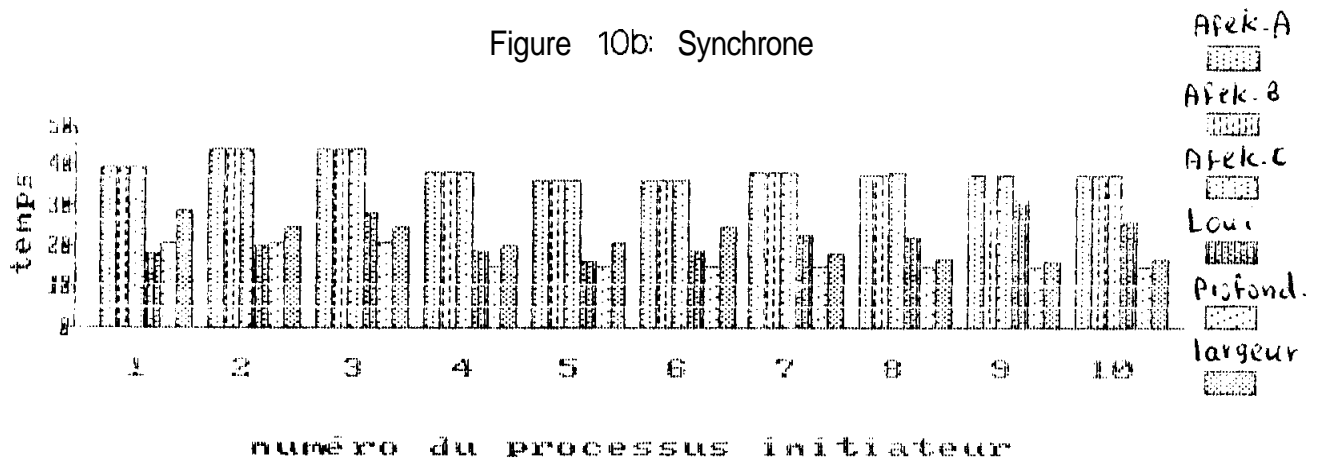


Figure 10c: Asynchrone

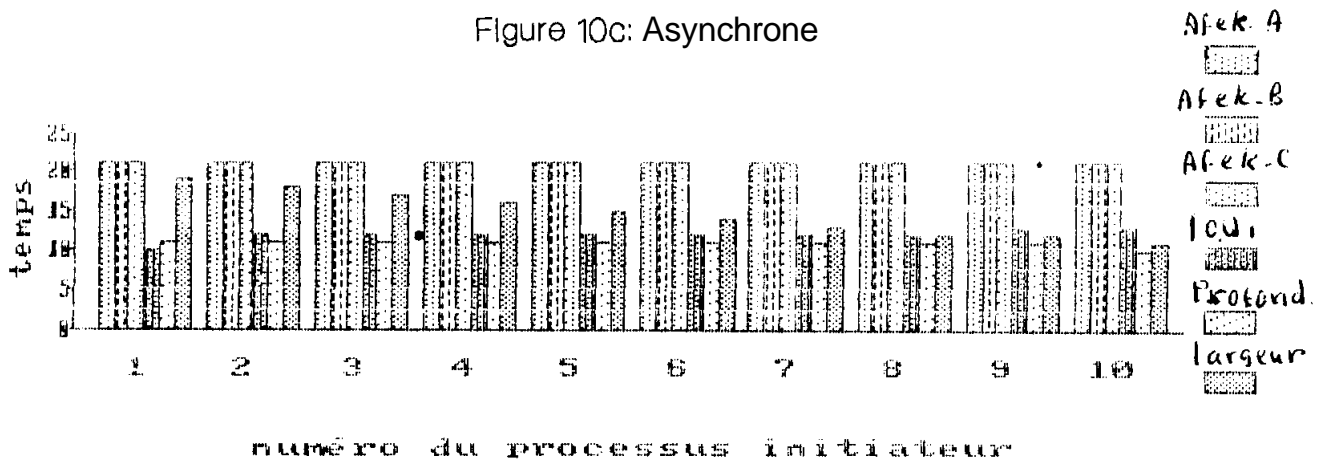


Figure 10d: Synchrone

Figure IV.10: Individuelle

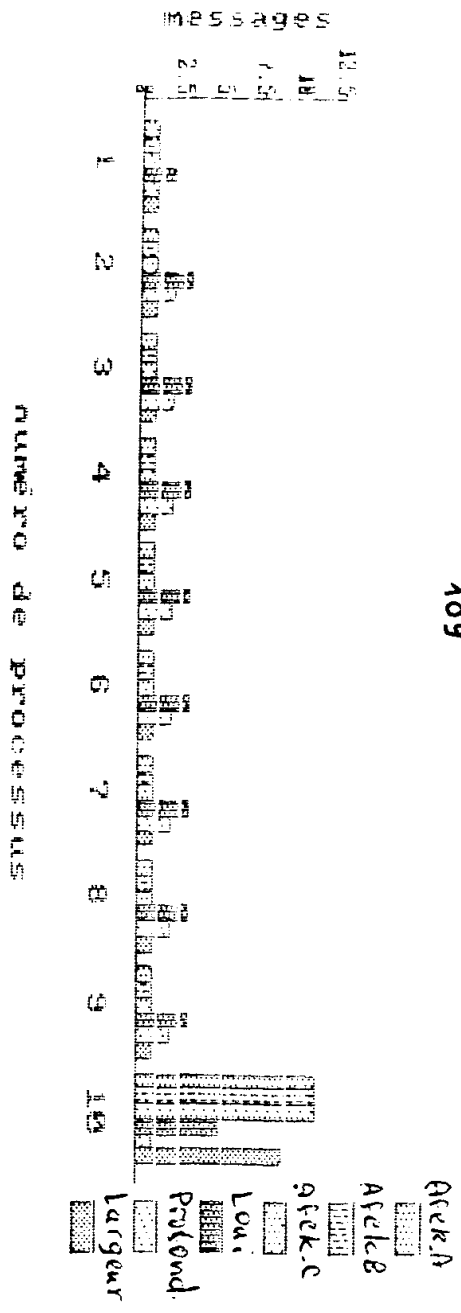


Figure 118: Asynchrone

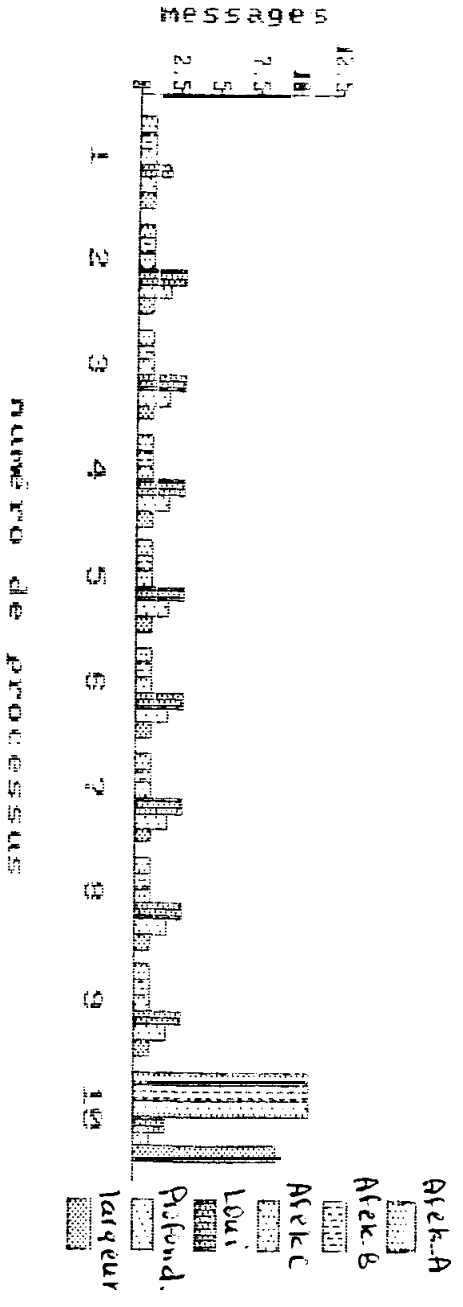


Figure 11b: Synchrone

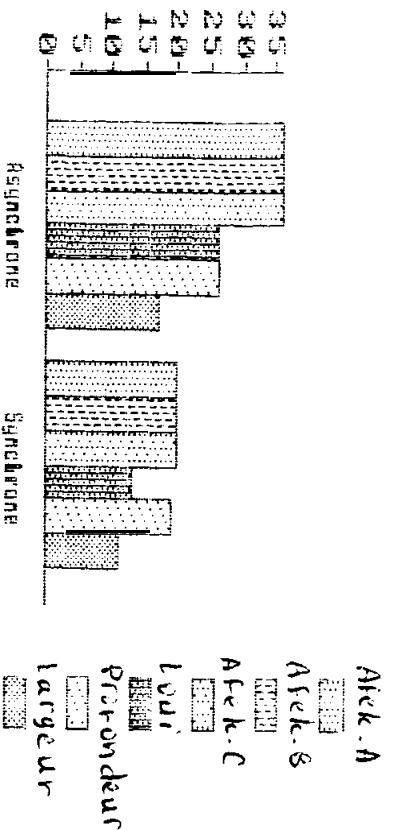


Figure 11c: Temps

Figure IV.11: Individuelle (Pmax Int.)

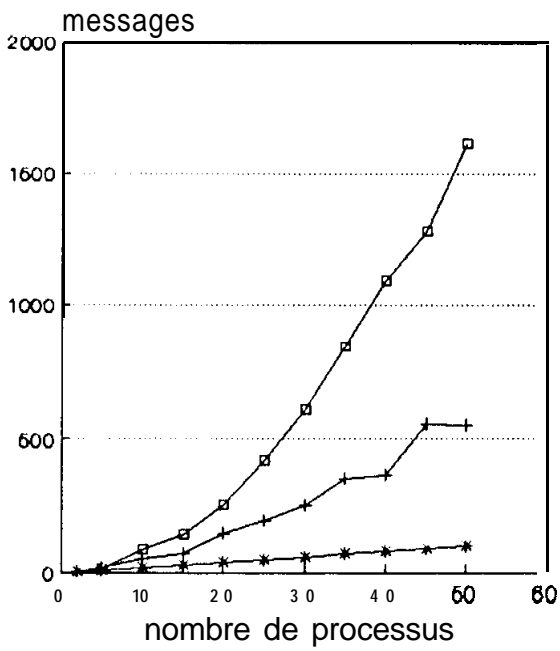


Figure 12.a: Asynchrone

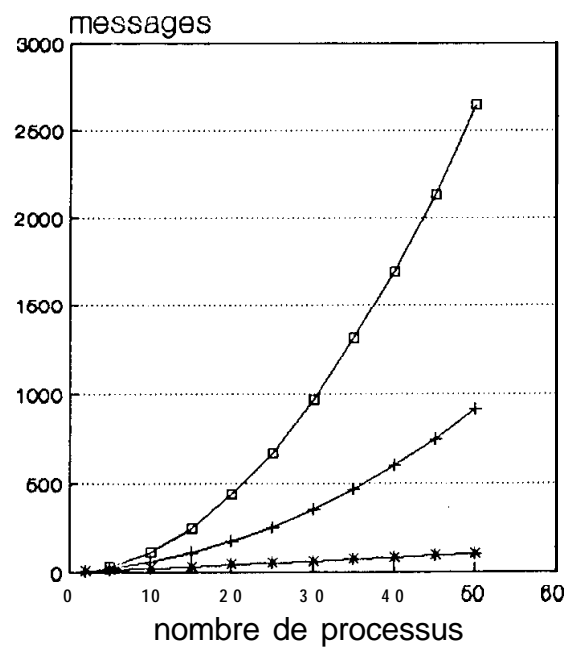


Figure 12.b: Synchrone

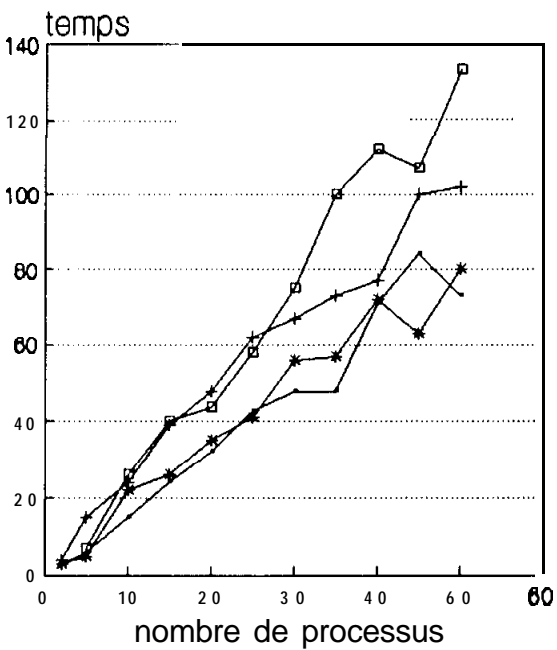


Figure 12.c: Asynchrone

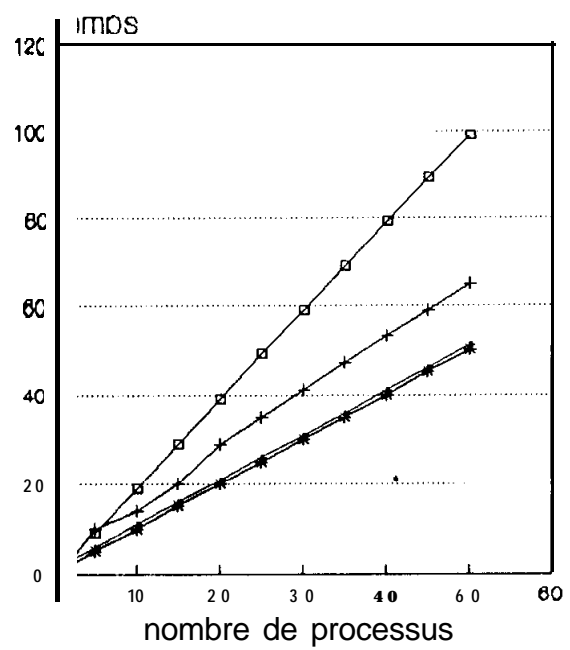


Figure 12.d: Synchrone

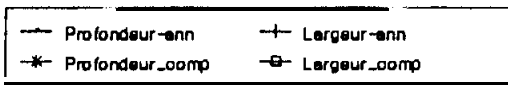


Figure IV.12: Collective

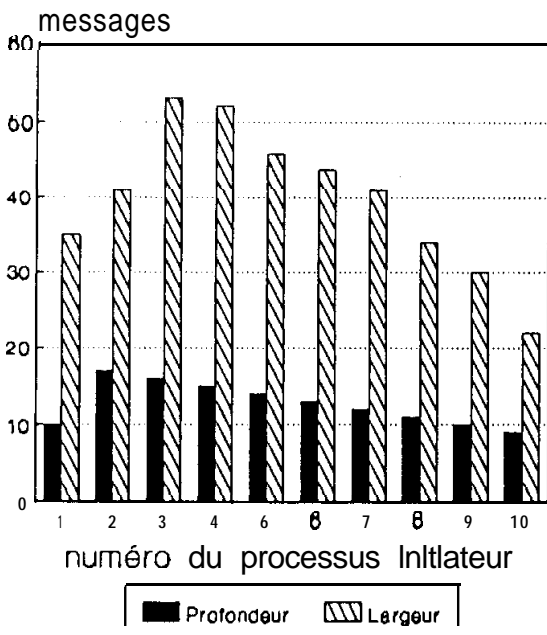


Figure 13.a: Asynchrone

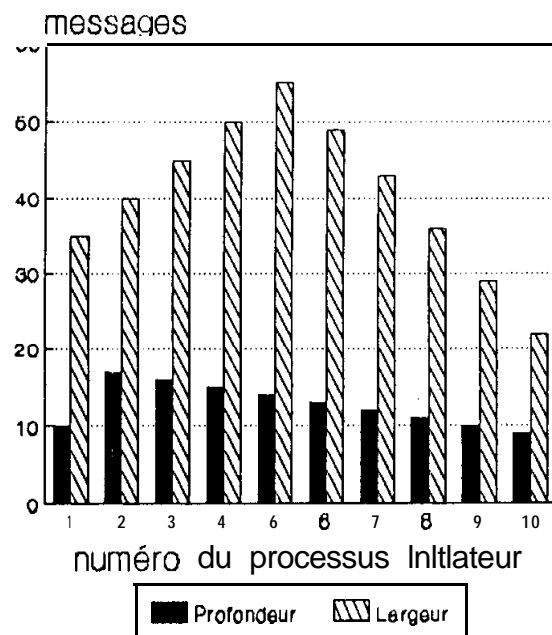


Figure 13.b: Synchrone

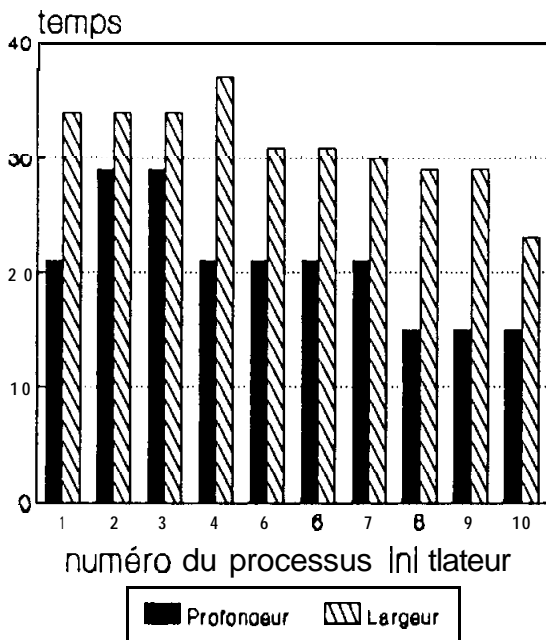


Figure 13.c: Asynchrone

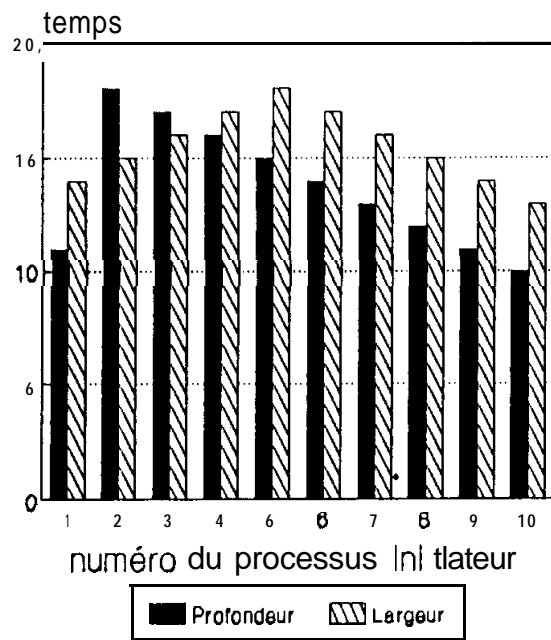


Figure 13.d: Synchrone

Figure IV.13: Individuelle (anneau)

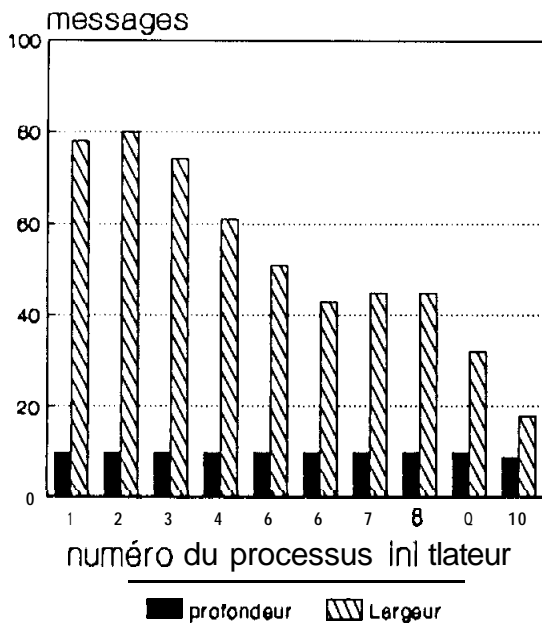


Figure 14.a: Asynchrone

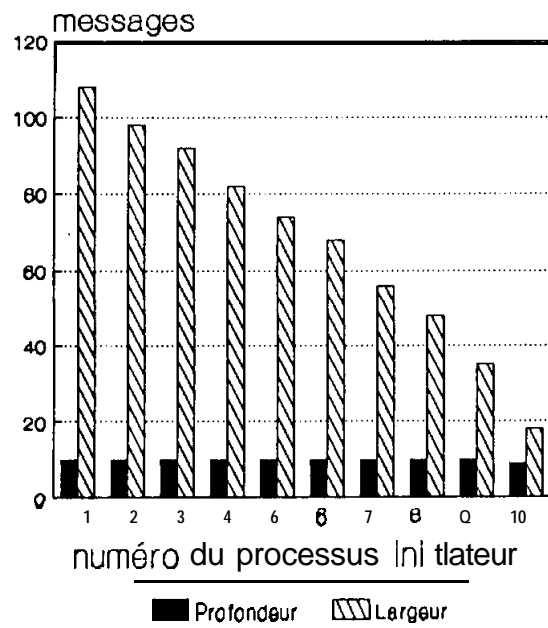


Figure 14.b: Synchrone

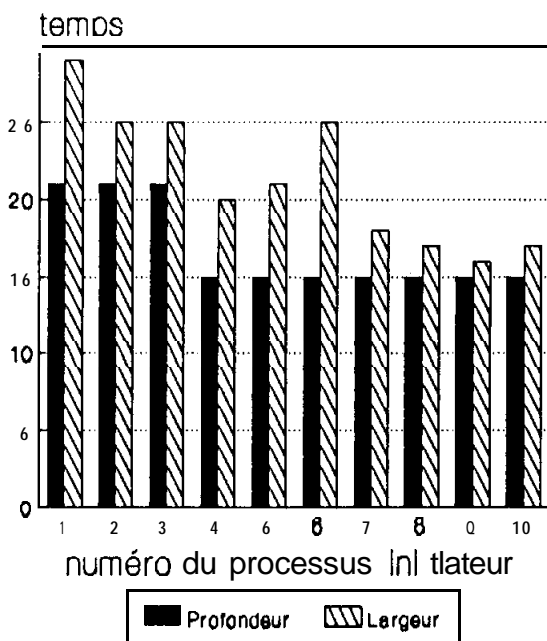


Figure 14.c: Asynchrone

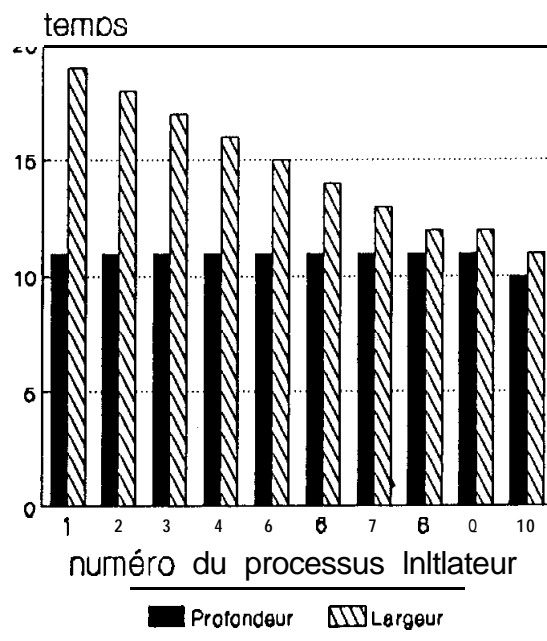


Figure 14.d: Synchrone

Figure IV.14: Individuelle (complet)

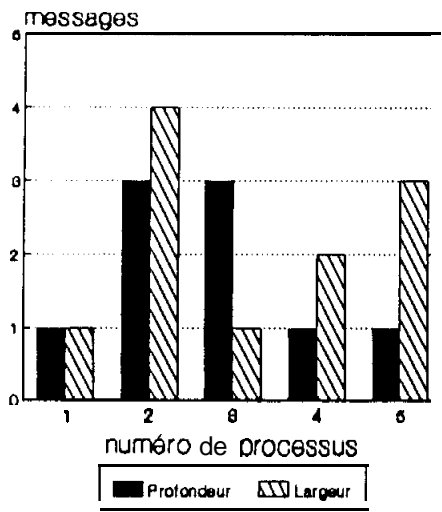


Figure 15.a: Asynchrone

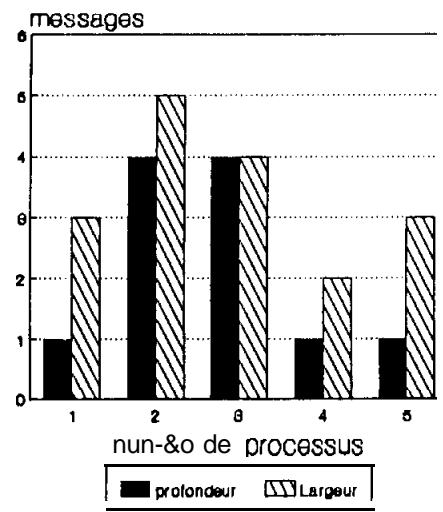


Figure 15.b: Synchrone

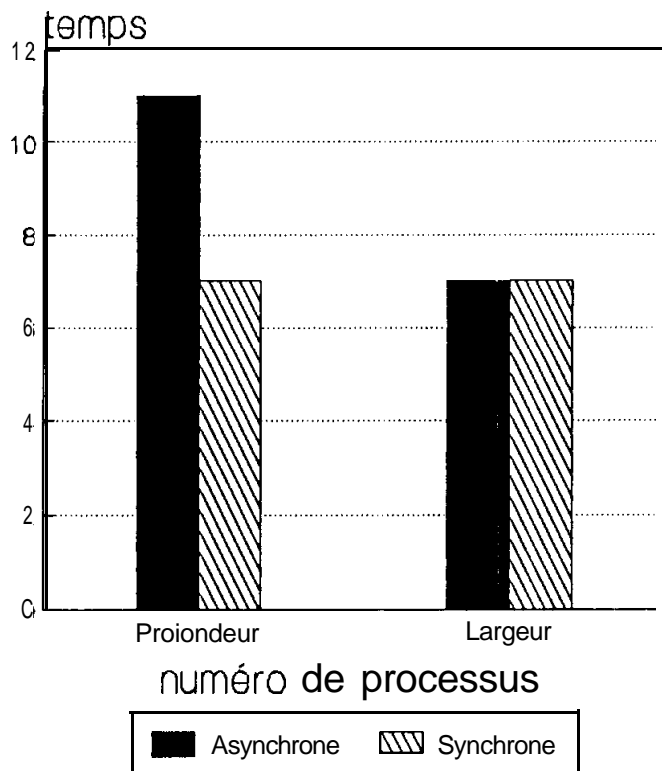


Figure 15.c: Asynchrone / Synchrone

Figure IV.15: Collective

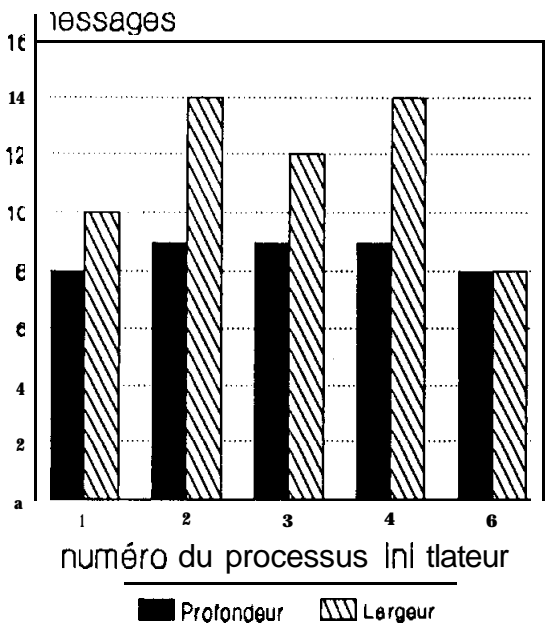


Figure 16.a: Asynchrone

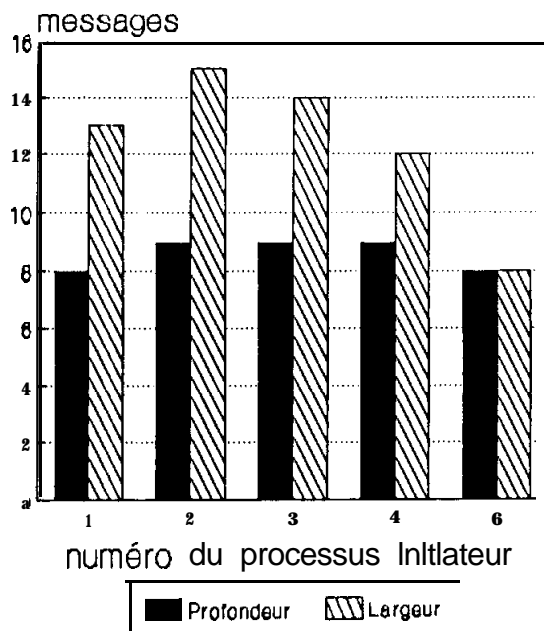


Figure 16.b: Synchrone

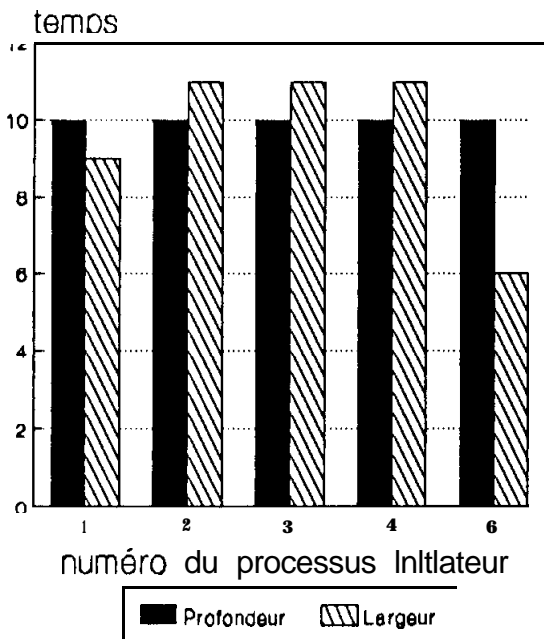


Figure 16.c: Asynchrone

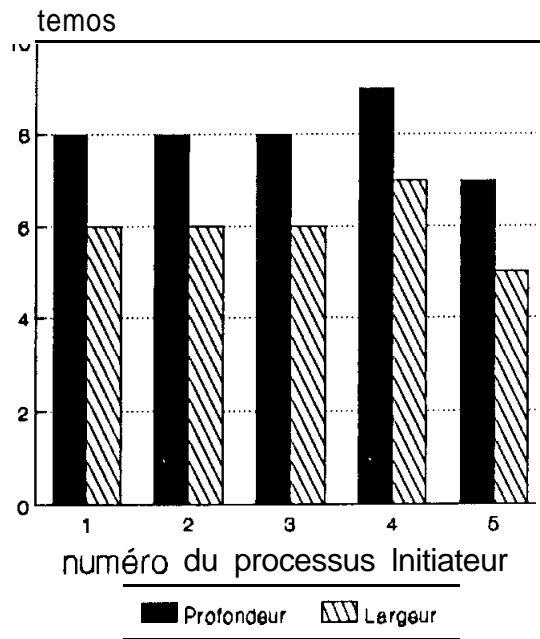


Figure 16.d: Synchrone

Figure IV.16: Individuelle

CONCLUSION

L'étude et l'évaluation des algorithmes distribués d'élection nous ont permis de mieux saisir le comportement de ce type d'algorithmes. Concernant les résultats obtenus, outre les performances évaluées, et l'étude comparative de celles-ci, un certain nombre de résultats ont été obtenus:

- Détection d'erreurs. Par exemple, l'algorithme de Loui [LOU 86] qui ne fonctionne pas sur certaines configurations de réseau (chapitre III).
- Confirmation de certaines complexités théoriques. Par exemple, le cas de l'algorithme de Chang [CHA 79].
- Obtention de complexités tout à fait différentes des complexités prévues. Par exemple, les algorithmes de Afek [AFE 85], et de Hirshberg [HIR 80].
- Certains comportements sont différents des prévisions théoriques. Afek [AFE 85] a proposé trois algorithmes A, B et C qui ont été conçus de façon à ce que l'algorithme B puisse améliorer la complexité en messages de A. Ce dernier, ayant une complexité temporelle de $O(n)$. L'algorithme C, était supposé réduire la complexité temporelle de B, en conservant sa complexité en messages. Or, l'évaluation a montré que les trois algorithmes ont pratiquement les mêmes comportements (chapitre IV).

L'étude des algorithmes a permis de comparer leurs performances en nombre de messages et en temps. Le fait que cette évaluation ne dépend pas d'un environnement d'exécution réel particulier permet d'affirmer que certains des résultats obtenus, tels que les résultats de la comparaison pourraient être généralisés à d'autres environnements.

Les meilleurs algorithmes sont:

1- Pour la topologie en anneau:

-- anneau croissant:

| collective | | individuelle | |
|------------|----------|--------------|-------------------|
| messages | temps | messages | temps |
| I Chang | Chang | Profondeur | 1. Korach-2 |
| Profondeur | Korach | | I 2. Profondeur |
| (+) | Korach-2 | | I 3. Chang (++) I |

(++): La numérotation signifie que les algorithmes sont classés dans cet ordre, mais la différence est minime.

(+) : Ces algorithmes offrent pratiquement les mêmes performances.

-- anneau décroissant:

| collective | | | individuelle | | |
|------------|-----------|---------|--------------|-------|------------|
| messages | temps | | messages | temps | |
| | I Async | Sync | | | |
| Profondeur | 1.Profond | Chang | Profondeur | I | Profondeur |
| I Dolev | 2.Dolev | Profond | | | |

2- Pour la topologie à maillage complet:

| collective | | | individuelle | | |
|----------------|-------|------------|--------------|------------|----------------|
| messages | temps | | messages | temps | |
| I 1.Profondeur | Async | Sync | I | Profondeur | I 1.Profondeur |
| 2.Loui, A, | | | | | 2.Loui |
| I B, C | Loui | Profondeur | | | 3.Largeur |
| | | Loui | | | |

3- Pour le graphe quelconque:

| collective | | | individuelle | | |
|------------|---------|---------|--------------|------------|------------|
| messages | I | temps | messages | I | temps |
| Profondeur | Async | Sync | I | Profondeur | Async |
| | | | | | Sync |
| | Largeur | Profond | I | | Profondeur |
| | | Largeur | I | | Largeur |

En plus de la classification, l'étude comparative a permis de conclure que:

1. L'algorithme très simple de Chang et Roberts [CHA 791] qui était à l'origine de la technique d'extinction sélective, a une bonne performance comparée aux autres algorithmes qui ont été conçus plus tard et dont la mise en œuvre est bien plus délicate ou même parfois complexe ([DOL 821, [HIR 801. Mis à part le cas extrême où l'anneau est décroissant, cet algorithme est le plus intéressant de par sa simplicité.
2. Dans le cas de la topologie à maillage complet, l'algorithme (corrigé) de Loui [LOU 861] qui est basé sur le sens d'orientation global, donne de bons résultats par rapport aux algorithmes A, B et C de Afek [AFE 85] qui utilisent la notion de 'niveau' représentant le "travail" effectuée par les processus participant à l'élection.

3. En tenant **compte** de toutes les configurations de reseau **étudiées** il apparait **que** l'algorithme d'exploration en profondeur est le plus performant.
4. Concernant les deux algorithmes **conçus** pour des topologies quelconques, **nous** avons **observé que**, dans **certains** cas, le calcul d'election base **sur** une exploration en profondeur est plus rapide en temps d'election **que** celui base sur l'exploration parallele. Le **phénomène** a **été** explique par le fait que dans **le** calcul parallele, la generation excessive de messages influence de maniere negative le temps d'election.

Les resultats **obtenus** permettent d'affirmer que si les evaluations de **complexités** theoriques sont necessaires, elles ne suffisent pas, pour autant, a faire un choix de mise en oeuvre dans **un** cadre reel, et qu'elles doivent, par consequent, Ptre **complétées** par d'autres **méthodes** de validation-evaluation parmi lesquelles la simulation **occupe** une place importante.

Le choix d'un algorithme d'election depend, principalement, de la **nature** du **problème** a **résoudre**, et de l'importance que l'application donne **au trafic** ou **au** temps de reponse. La **simplicité** de mise en oeuvre peut, aussi, constituer **un critère** de choix non negligeable.

Le travail **réalisé** dans cette these est loin d'avoir **cerné** **tous** les problemes lies à ce type d'algorithmes. **D'autres** parametres pourraient **être étudiés**, tels que l'influence de la **longueur** des **messages** sur les performances.

Les algorithmes **synchrones** qui sont encore **peu nombreux**, **n'ont pas été évalués**. Une etude **appropriée** de **leurs propriétés** est encore a **réaliser**.

Il **serait, également,** interessant et instructif de **completer cette** etude par l' experimentation de ces algorithmes dans **un** environnement d'execution **réel**[BAD 911.

La mise en evidence des techniques de **base** sur lesquelles sont bases les algorithmes d'election pourrait constituer le point de depart d'une etude plus approfondie de ces algorithmes et d'autres algorithmes de **contrôle** distribue. **Elle** permettrait de degager les points **communs** a ces differents algorithmes, et de **mettre au point** **une** methodologie pour **leur** conception.

B I B L I O G R A P H I E

- [ADA 88a] Adam, M.
Exerimentation on parallel machines in helpful to
analyse distributed algorithms / Ph. Ingels, C. Jard,
J.M. Jezequel, M. Raynal. - Workshop on "parallel and
distributed algorithms", France, (September 1988).
- [ADA 88b] Adam, M.
Algorithmes distribues **synchrones** et systemes repartis
asynchrones: concepts, mise en oeuvre et
experimentation/P. Ingels, M. Raynal.-IRISA, PI n°411,
(Juin 1988).-p 28.
- [AFE 85] Afek, Y.
Distributed algorithms for election in unidirectional
and complete networks.- University of California: these
de docteur, Inovembre 1985).
- [ATT 87] Attiya, H.
Constructing efficient election algorithms from
efficient traversal algorithms.- (Juillet 1987).
- [ATT 89] Attiya, H.
Efficient election chordal ring networks/ J.V. Leeuwen,
N.Santoro, S. Zaks.- Algorithmica, (1989).- pp 437-446.
- [ANG 87] Angluin, D.
Local and global properties in networks of processes.
-Proc. 12th ACM symposium on theory of computing,
(April 1980).-pp 82-93.- In [LAV 87].
- [AUD 88] Audureau, E.
Theorie de la programmation et logique temporelle,
seconde partie: validation d'algorithmes paralleles/
L.F. Del Ceno, P. Enjalbert.-TSI, (1988).-pp 181-200.
- [BAD 91] Badache, N.
Une machine virtuelle Estelle **repartie/** S. Bouallag,
N. Semaoune.- Proc. 1st Maghrebin Symposium On
Programming and System.- Algiers, (October 21-23 1991).
-pp 115-131.
- [BAR 87] Baruch, A.
optimal distributed algorithms for minimum weight
spanning tree, counting leader election and related
problems.- ACM Proc. 9th symp. on theory of computing,
networks, (May 1987).-pp 230-240.
- [BER 81] Bernstein, P.A.
Concurrency control in distributed data base systems/
N. Goodman.-Computing Surveys, vol. 13, 2(June 1981).
-pp 185-221.- In [MAD 87].

- [BOD 861 Bodlaender, H.L.
New upperbounds for decentralized extrema-finding in a ring of **processors**.- STACS 86, Lecture notes in C.S. (Springer **Verlag**), (1986).- pp 119-129.
- [BOD 88] Bodlaender, H.L.
A better lower bound for distributed leader finding in bidirectional asynchronous rings of processors.- **Info. Proc. Let.**, vol 27, (1988).
- [BOU 88] Bouge, L.
On the existence of symmetric algorithms to find leaders in **Networks** of Communicating Sequential Process.- **Acta informatica**, vol 25, (1988), (Springer **Verlag**).
* pp 179-201.
- [BUD 871 Budkowski, S.
An introduction to Estelle: A specification language for distributed algorithms / P. Dembinski.- **Computer Networks and ISDN systems**, 14 (1987) 3-32.
- [CHA 791 Chang, E.J.
An improved algorithm for decentralized extrema-finding in circular configuration of processors / R. Roberts.
- **Comm. ACM**, vol 22, N05, (may 1979).- pp 281-283.
- [CHA 881 Chan, M.Y.
Distributed election in complete networks/ Chin, F.Y.L.
* **Distributed computer** (Springer **Verlag**), (1988).
- pp 19-22.
- [COR 811 Cornafion (nom collectif).
Systemes informatiques répartis, concepts et techniques.
- **DUNOD Informatique**, (Aout 1984).
- [COU 861 Courtiat, J.P.
Estelle: un langage ISO pour les algorithmes distribués et les protocoles / P. Dembinski, R. Groz, C. Jard.
- Publication interne n°326, **IRISA**, (Décembre 1986).
- [COU 871 Courtiat, J.P.
Estelle: un langage ISO pour les algorithmes distribués et les protocoles / P. Dembinski, R. Groz, C. Jard.
* **TSI, Technique et Sciences Informatiques**, (1987).
- [DIJ 801 Dijkstra, E. W.
Termination detection for diffusing computation / C.S. Scholten.-**Inf. Proc. Letters**, vol. 11,1 (August 1980).-pp 1-4.
- [DOL 821 Dolev, D.
An $O(n \log n)$ unidirectional distributed algorithm for extrema-finding in a circle/M. Klawe, M. Rodeh.-**Journal of algorithms**. vol 3. (1982).-pp 245-260.- in [RAY 851

- [DUD 871] Duda, A.
Estimating global time in distributed system/ G. Haddad, Y.Haddad, G. Bernard.- In **Proc.** 7th Int. Conf. on Distributed Computing Systems, Berlin, (1987).
- [FRA 821] Franklin. R.
On an improved algorithm for decentralized **extrema-** finding in circular configuration of processors.
- communication of the ACM. vol 25. n^o4.(May 1982).
-pp 336-337.
- CFRE 87a] Frederickson, G. N.
The impact of asynchronous **communication** on the problem of electing a leader in a ring / N. A. Lynch.-Proc. of 16th annual ACM on theory of computing, (1989).
- pp 493-503.
- [FRE 87b] Frederickson, G. N.
Electing a leader in a synchronous ring / Lynch A.
-**journal of** the association of computing machinery.
vol 34. n^o1. (January 1987).-pp 98-115.
- [GAF 85a] Gafni. E.
Improvements in the bit complexity of two **message-** optimal election algorithms.- ACM symp. on PODC, (1985).
- PP 175-185.
- [GAF 85b] Gafni. E.
Time and message Bounds for election in synchronous and asynchronous complete networks/ Afek. Y.
-Acm symposium on PODC, (1985).-pp 186-195.
- [GAR 82] Garcia-molina, H.
Elections in distributed computing system.- IEEE Transactions on Computer, vol-31, N^o1, (January 1982).
-pp 48-59.
- [GOL 87] Goldreich. O.
Electing a leader in a ring with link failures / Shrira L.-**Acta** informatica. vol 24. (1987).-pp 79-91.
- [GRO 87] Groz. R.
Peut on verifier des algorithmes distribues par simulation?.-Algorithmique distribuee, INRIA, la **colle-sur-loup** (Alpes maritimes), (6-10 avril 1987).
-pp 143-145.
- [GUS 84] Gusella, R.
A network time controller for a **distributed** Berkeley UNIX system / S. Zatti.- Distr. **Proc.** Tech. Comm. Newsletter, vol. 6, n^oSI.2, IEEE, (June 1984).-pp 7-15.
-In [GUS 86].

- [GUS 851] **Gusella, R.**
 The berkeley UNIX 4:3BSD time synchronization protocol:
 protocol **spécification** / S. Zatti.- University of
 California, Berkeley: Report n°UCB/CSD 85/250,(June 851).
 -In [GUS 861.
- [GUS 861] **Gusella, R.**
 An election algorithm for a distributed clock
 synchronization program / S. Zatti.-IEEE, (1986).
- [HEB 851] Hebuterne. G.
 Evaluation de performances des algorithmes distribués.
 -Actes du premier colloque, An **gouleme**, (16-18 sept 1985).
 pp 151-166.
- [HEL 86a] Helary. M.
 Controlling knowledge transfers in distributed algorithms.
 application to deadlock detection.-rapport de recherche
 NQ 493, INRIA, (Mars 1986).
- [HEL 86b] Helary. M.
Calcul distribue d'un **extremum et du routage** associe
 dans un **réseau** quelconque / A. Maddi, M. Raynal.
 - RR N°516, INRIA, (Avril 1986).
- [ISO 891] ISO/IEC JTC 1/SC 21.
 Working draft for a **tutorial** for **Estelle**.- Information
 retrieval transfer and managemenr for **OSI USA (ANSI)**,
 Boston, (JUN 19-20, 1989).
- [JAR 871] Jard, C.
 Introduction à la verification des algorithmes
distribués.-Algorithmique distribuee, INRIA, la colle-
sur-loup, (Alpes maritimes), (6-10 avril 871).
 -pp 119-138.
- [JAR 891] Jard, C.
 un compilateur Estelle multi-processeurs pour l'experi-
 mentation d'algorithmes distribués sur machines
 paralleles / J.M. Jezequel.-Publication interne nQ453,
 (Janvier 891).
- [JEZ 89a] Jezequel, J.M.
 Outils pour l'experimentation d'algorithmes distribués
sur machines paralleles. - These de docteur de
 l'universite de Rennes, (19 **Octobre 89**).
- [JEZ 89b] Jezequel, J.M.
 Un compilateur Estelle multi-processeurs **pour** l'experi-
mentation d'algorithmes distribués sur machines
 paralleles.- Rennes: **IRISA**. Publication interne NO 453,
 (Janvier 1989).

- [KAI 88] C. Kaiser,
Election sur un anneau par parcours de reconnaissance
Anneau unidirectionnel ou bidirectionnel.- Paris: CNAM,
dept. mathématique et informatique, (23 décembre 1988).
- [KOR 811] Korach, E.
A probabilistic algorithm for decentralized extrema
finding in a circular configuration of processors.- Res.
Rep. C. S.- University of Waterloo: Dept. of computer
science, (1981).- In [LAV 89C].
- [KOR 841] Korach, E.
Tight lower and upper bounds for some distributed
algorithms for a complete network of processors /
S. Moran, S. Zaks.- ACM, (1984). pp 199-207.
- [LAM 781] Lamport, L.
Time, clocks and the ordering of events in distributed
Systems.- ACM, (July 19781, vol 21, number 7.
- [LAV 871] Lavault, C.
Un algorithme d' election entièrement distribue et
décentralisé dans un **système** distribue quelconque.
- RR N0625, INRIA, (19871.
- [LAV 89C1] Lavault, C.
Average number of messages for distributed leader
finding in rings of processors.- Padesbour FRG: Lecture
notes on **theoretical** aspects of computer science,
(Fevrier 16-18, 1989).-pp 269-281.
- [LAV 8911] Lavallee, I.
Yet another distributed election and (minimum-weight)
spanning tree algorithms/ C. Lavault.-RR NQ1024, INRIA,
(Avril 19891.
- [LAV 901] Lavallee, I.
Efficient routing protocols in nameless networks /
C. Lavault.- RR nQ1254, **(Juin 19901.**
- [LEL 771] Le Lann, G.
Distributed systems. Towards a formal **approch.**-IFIP
Congress 1977, Holland.-pp **155-160.**-In [RAY 851.
- [LOU 861] Loui, M. C.
Election in complete network with a sense of direction/
T. A. Matsushita, D.XB. West.- Information **proc. let.**,
vol 22, **(1986).**- pp 185-187.
- [MAD 871] Maddi, A.
Contrôle des transferts d'informations dans **les systèmes**
distribués, application a la detection de l'interblocage
et à la l' election dans un reseau quelconque de
processus.- **Université** de Rennes 1: these de docteur en
informatique, **(Mai 1987).**

- [MAH 881 Mahmoud Bacha. R.
Etude et implementation d' algorithmes de contrôle distribue: exclusion mutuelle et election.-Alger: USTHB. these de magister.-(Septembre 1988).
- [MAS 891 Masuzawa, T.
Optimal fault-tolerant distributed algorithms for election in complete networks with a global sence of directives / N. Nishikoma, K. Haqihora, N. Tokuba.-proc 3rd Int. Wok., Nice, France (1989).- Lecture notes in computer Science.- p 11.
- [MAT 89] Matias, Y.
Simple and efficient election algorithms for anonymous networks / Y. Afek.-Lecture notes in C.S., Distr. alg., 3rd Int. Work., Nice, France, (Septembre 26-28, 1989).-pp 183-194.
- [NAV 881 Navaratnam, S.
Reliable group communication in distributed systems / S. Chanson, G. Neufeld.-Proc. 8th Int. Conf. Distributed Computing Systems. -CS Press, Los Alanitos, California, order nQ865, (June 88).-pp 439-446.
- [PAC 841 Pachl, J.
Lower bounds for distributed maximum-finding algorithms/ E. Korach, D. Rotem.- Journal of the association for computing machinery, vol 31, N^o4,, (Octobre 1984).-pp 905-918.
- [PET 821 Peterson. G. I.
An $O(n \log n)$ unidirectionnal algorithm for the circular extrema problem.-ACM Toplas. vol 4. n^o4, (octobre 1982).-pp 758-762.
- [PUJ 871 Pujolle, G.
Reseaux et télématique/ D.Secret, D. Dromard, E.Horlait. -Edition: Eyrolles, Tome I, Tome II, (1987).
- [RAY 851 Raynal. M
Algorithmes distribués et protocoles.-Eyrolles, (1985).
- [RAY 88] Raynal. M.
Synchronisation et contôle des systemes et des programmes repartis / J.M. Helary.- Edition: Eyrolles, (1988).-p 194.
- [RAZ 841 Razouk, R.R.
Performance Analysis using Timed Petri Nets/C.V. phelps. - 4th IFIP workshop on protocol specification, verification and testing. North holland ed: (Juin 1984). -In [HEB 85].

- [SAN 861 Santoro, N.
Guessing games and distributed computations in
anonymous networks / J. Van Leeuwen, J. Urrutia, S. Zaks.
-Preliminary draft, (March 86).
- [SHR 85] Shrira, L.
Methodological construction of reliable distributed
algorithms / M. Rodeh.-TR361, computer science dept.,
Technion, Haifa 32000, (February 1985).- in [GOL 871.
- [SIN 801 Sinclair. J. B.
Decentralized extrema-finding in circular configurations
of processors / Hirshberg D.S.-Comm. ACM. vol 23. n^o 11.
(November 1980).-pp 627-628.
- [SKY 871 Skym. S.
A distributed Election and Spanning tree algorithms based
on depth first search traversals.-Marhus university,
computer science department, Denmark, (August 1987).-p11.
- [SMI 791 Smith, R.G.
The contact net protocol: High level communication and
control in a distributed problem solver.- Huntsville,
AL.: in Proc. 1st Int. Conf. Distr. Comput. Syst.,
(October 1979).-pp 185-192.-In [GAR 821.
- [SPI 891 Spirakis. P.
Symmetry breaking in asynchronous rings with $O(n)$
messages / B. Tampakas, T.Athamasios.-Lecture notes in
computer science.-Distributed algorithms.- 3rd Int.
workshop. Nice. France. (Septembre 1989).-pp 235-241.
- [VER 871 Verjus, J.P.
construction et preuve d'algorithmes distribues
-algorithmique distribuée.-INRIA, la colle-sur-loup
(Alpes maritimes), (6-10 Avril 1987).-pp 147-168.
- [YAM 891 Yamashita. M
Electing a leader when processor Identity Number are
not distinct / T. Kameda. -Distributed algorithms 3rd
International workshop. proc.-Lecture notes in computer
science, Nice, France, (septembre 1989).-pp 304-314.

A N N E X E

```

specification S;

default individual queue;

const  nbsites=10;
       nbsites_1=9;

type   numsite=1..nbsites;

channel liaison(sortie,entree);
by sortie:
    elir(numero: integer) ;
    elu(numero: integer) ;
by entree: ;

module type-mod (mon_numero: integer) ;
ip    5: liaison (sortie) ;
      e: liaison(entree);
end;

body type-bod for type_mod;

state participant,non_participant,fin,elu;

var
coordonateur: integer ;

initialize to non_participant
begin
    coordonateur:=0;
end;

trans
from non_participant to participant
(* provided mon_no=10*)
priority 1
begin
    output 5. elir (mon_no) ;
end;

trans
from non_participant to participant
when e.elir(numero)
priority 3
begin
    if numero<mon_numero then output s.elir (mon_numero);
    else output s.elir(numero);
end;

```

```

trans
from participant to participant
when e. el ir (numero)
priority 3
begin
  if numero > mon_no then output s. el i r (numero) ;
end;

trans
from participant to fin
when e.elu(numero)
priority 4
begin
  coordonateur: =numero;
  if coordonateur < mon_no then output s. el u (numero) ;
end;

trans
from participant to elu
when e.elir(numero)
provided (numero = mon_numero)
priority 5
begin
  output a. el u (numero) ;
  coordonateur: =mon_numero;
end;

END; (*endbody*)

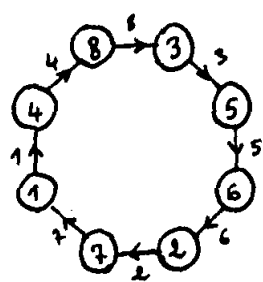
MODVAR noeud: array[1..nbsites] of type-mod;

initialize
BEGIN
  ALL i: numsite DO INIT noeud[i] WITH type_bod(i);
  ALL i: 1..nbsites_1 DO connect noeud[i].s TO noeud[i+1].e;
  connect noeud[nbsites].s TO noeud[1].e;
end;
END.

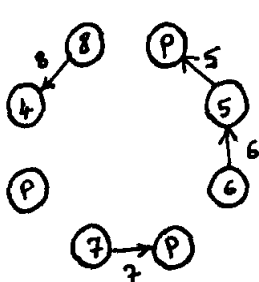
```

Exécution de l'algorithme de LOUI (Lou 86] non corrigé :

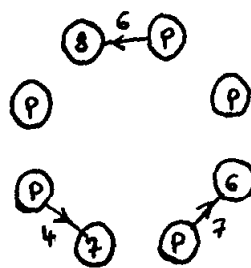
| Proc1 | Proc2 | Proc3 | Proc4 | Proc5 | Proc6 | Proc7 | Proc8 |
|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|----------------------|----------------------|
| D=71 E-7 | D=71 E-7 | D=71 E-7 | D=777 E- | D=777 E- | D=7714 E- | D=7714 E- | D=7714 E- |
| S(1,1,1) R(1,7) | S(1,1,2) R(1,6) | S(1,1,3) R(1,8) | S(1,1,4) R(1,1) | S(1,1,5) W. 31 | S(1,1,6) R(1,5) | S(1,1,7) R(1,2) | S(1,1,8) R(1,4) |
| Passé R(7,4) | Passé R(7,7) | Passé R(7,6) | S(7,7,4) R(7,8) | S(7,7,5) R(7,6) | S(7,7,6) | S(7,7,7) | S(7,7,8) |
| S(7,14,4) | S(7,14,7) | S(7,14,6) | Passé | Passé | R(14,7) | R(14,4) | R(14,6) |
| | | | | | Passé | S(-6,-6,7) erreur | S(-6,-6,8) erreur |



(α)



(β)



(γ)

Exécution de l'algorithme de LOUI [Lou 86] corrigé :

| Proc1 | Proc2 | Proc3 | Proc4 | Proc5 | Proc6 | Proc7 | Proc8 |
|------------------------------|------------------------------|------------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|
| D=21 E=7 | D=21 E=7 | D=21 E=7 | D=227 E=2 | D=227 E=2 | D=226 E=3 | D=2285 E=5 | D=2283 E= |
| S(1.1.1) R(1.7) passif | S(1.1.2) R(1.6) passif | S(1.1.3) R(1.8) passif | S(1.1.4) R(1.1) actif | S(1.1.5) R(1.3) actif | S(1.1.6) R(1.5) actif | S(1.1.7) R(1.2) actif | S(1.1.8) R(1.4) actif |
| | | | S(7.7.4) R(7.8) | S(7.7.5) R(7.6) | S(7.7.6) R(6.7) | S(7.7.7) R(6.4) | S(7.7.8) |
| | S(7.6.4) R(7.7) | S(7.6.5) R(7.5) | Passif | Passif | | | R(6.5) |
| | | | R(2.7) | R(1.8) | passif | S(2.2.7) | S(2.2.8) |
| | | | S(1.3.7) | S(1.3.8) | R(3.8) | | R(3.7) |
| | | | | | S(2.5.8) | B(5.8) passif | S(5.5.8) |
| | | | | | | R(5.8) S(3.0.8) | R(0.8) Etu |

