

Communications
dans les réseaux de processeurs

Opération RUMEUR du PRC C³

22 août 1994

Ouvrage rédigé avec le concours du Ministère de l'Enseignement Supérieur et de la Recherche (Direction de l'Information Scientifique et Technique et des Bibliothèques).

Préface

En moins de dix ans, le calcul parallèle est passé des laboratoires de recherche aux centres de calculs. Il s'agit pour l'informatique d'une véritable révolution, puisqu'elle touche profondément l'activité centrale de cette discipline : la programmation.

Ceux qui ont déjà eu l'occasion de programmer une machine parallèle savent à quel point les habitudes de la programmation séquentielle sont largement remises en cause. Pour obtenir les résultats espérés — un programme performant, ou tout simplement correct —, il faut changer de mode de pensée. Cette impression est largement confirmée par le développement de la théorie du calcul parallèle : les modèles changent, et leur étude conduit à des problèmes nouveaux.

Parmi ceux-ci, celui des communications est central. Si les programmes parallèles étaient faits de calculs sans communication, il n'y aurait rien de plus à dire que ce qui concerne les programmes séquentiels, ce qui, en soi, n'est déjà pas mal. Mais les communications sont l'essence même du calcul parallèle : la résolution d'un problème unique par plusieurs processeurs passe forcément par la distribution des données, l'échange et la comparaison de résultats intermédiaires, ou encore la diffusion des solutions de sous-problèmes ; toutes opérations qui nécessitent des communications.

L'étude des communications dans le calcul parallèle n'est pas récente, mais ce domaine a connu un essor important à l'arrivée des machines dites *massivement parallèles*. La raison en est simple. Il est physiquement possible d'interconnecter complètement un tout petit nombre de processeurs, mais dès que ce nombre croît, cela devient irréaliste. Il faut donc trouver des réseaux de communications efficaces, ainsi qu'une façon de les utiliser qui ne réduisent pas à néant les effets de l'augmentation du nombre des processeurs. Or, les machines massivement parallèles ont l'ambition de faire coopérer des centaines, voire des milliers de processeurs. On conçoit donc que ce sujet ait de l'importance.

En outre, les progrès de la technologie des processeurs n'améliorent en rien la situation. Même si les réseaux d'interconnexion, les algorithmes de routage et les dispositifs matériels de communication ont fait des progrès énormes en quelques années, le rapport entre vitesse de calcul et vitesse de communication des processeurs les plus récents est aujourd'hui plus que jamais défavorable aux communications. Il s'agit donc d'un problème fondamental, qui ne se résoudra pas de lui-même avec l'évolution de la technologie des circuits intégrés.

En France, le Programme de Recherches Coordonnées C^3 (Communication, Coopération, Concurrence) a reconnu très tôt la nécessité d'étudier ce sujet, puisqu'il l'a inscrit dans son titre lors de sa création en 1981, avant même que les premiers ordinateurs massivement parallèles ne soient conçus. Mais c'est la pratique de la programmation sur des machines parallèles qui a conduit la recherche sur ce sujet à se structurer. Dès 1986, C^3 a soutenu l'acquisition de machines parallèles à mémoire distribuée : l'Intel iPSC/1 de Rennes et le FPS T40 de Grenoble. Ces deux machines, prototypes de laboratoire plus que véritables calculateurs, ont permis à une génération de chercheurs de pratiquer le calcul parallèle et de déterminer ainsi, par l'expérience, les sujets de recherche critiques.

Les communications ont très naturellement été l'un de ces sujets. La coordination des équipes de recherche intéressées par ce problème a donné naissance en 1991 à une opération de C^3 , très joliment baptisée RUMEUR par ses animateurs. RUMEUR regroupe une vingtaine de chercheurs provenant de laboratoires répartis sur toute la France — principalement Bordeaux, Grenoble, Lyon, Orsay et Sophia-Antipolis ; mais aussi Lille, Rennes et Toulouse.

L'ambition de RUMEUR est de formaliser les modèles, de conceptualiser les idées et enfin, de développer des outils. Ces derniers devraient servir aux nombreux chercheurs qui font face à des problèmes de communications. Ces outils peuvent être théoriques, en faisant appel à des techniques empruntées aux mathématiques discrètes et à la théorie des graphes, ou plus pratiques, comme des bibliothèques de communications développées pour différentes machines parallèles.

A l'initiative de RUMEUR, une école d'été a été organisée en août 1992 à Cargèse (Corse) pour des étudiants de troisième cycle, répondant ainsi à une forte demande de formation dans le domaine des communications. Le support de cours de cette école, dûment corrigé et complété, grâce aux critiques des participants, est devenu cet ouvrage.

Grâce à la compétence de ses auteurs et à leurs échanges scientifiques soutenus — entre eux, mais aussi avec la communauté scientifique internationale dans laquelle RUMEUR a pignon sur rue —, ce livre présente et synthétise les connaissances les plus récentes sur les communications dans les machines parallèles. De plus, en très peu de temps, ses auteurs ont réussi le tour de force de faire un ouvrage très pédagogique, auquel ne manquent même pas des exercices corrigés.

Unique en son genre, il intéressera donc tous ceux qui, étudiants, informaticiens chevronnés ou utilisateurs avertis de calculateurs parallèles, sont intéressés par les aspects théoriques et pratiques des communications.

Patrice Quinton

Directeur de Recherche CNRS

Directeur du Programme de Recherches Coordonnées C^3

Table des matières

| | |
|--|-----------|
| Préface | v |
| Table des matières | vii |
| Contents | xi |
| Liste des figures | xiii |
| Auteurs | xvi |
| Avant-propos | xvii |
| 1 Éléments de théorie des graphes | 1 |
| 1.1 Définitions et propriétés des graphes | 2 |
| 1.1.1 Quelques définitions de théorie des graphes | 2 |
| 1.1.2 Connexité | 7 |
| 1.2 Paramètres des réseaux classiques | 8 |
| 1.2.1 Le cycle | 8 |
| 1.2.2 Le graphe complet | 9 |
| 1.2.3 L'hypercube | 10 |
| 1.2.4 La grille | 11 |
| 1.2.5 La grille torique | 12 |
| 1.3 Problème (Δ, D) | 13 |
| 1.4 Graphes de Cayley | 17 |
| 1.4.1 Définitions | 17 |
| 1.4.2 Hamiltonisme des graphes de Cayley | 18 |
| 1.4.3 Graphes sommet-transitifs | 19 |
| 1.4.4 Cas particuliers | 19 |
| 1.5 Exercices | 20 |
| 1.6 Solutions des exercices | 21 |
| 2 Graphes pour les réseaux d'interconnexion | 31 |
| 2.1 Graphes de de Bruijn et graphes de Kautz | 32 |
| 2.1.1 Définition à partir d'un alphabet | 32 |
| 2.1.2 Définition à partir du graphe représentatif des arcs | 36 |
| 2.1.3 Connexité et hamiltonisme | 38 |
| 2.1.4 Définition à partir des congruences | 41 |
| 2.2 Autres graphes | 42 |
| 2.2.1 Le graphe <i>shuffle-exchange</i> | 43 |
| 2.2.2 Le graphe <i>cube-connected-cycles</i> | 44 |
| 2.2.3 Le graphe <i>butterfly</i> | 45 |
| 2.2.4 La grille d'arbres <i>d</i> -aires | 46 |
| 2.2.5 Graphes composés | 47 |

| | | |
|----------|--|------------|
| 2.3 | Réseaux multiétages | 49 |
| 2.3.1 | Description préliminaire de deux réseaux simples | 50 |
| 2.3.2 | Définitions | 52 |
| 2.3.3 | Réseau de Clos | 54 |
| 2.3.4 | Réseau <i>butterfly</i> | 58 |
| 2.3.5 | Réseau de Beneš | 60 |
| 2.3.6 | Réseaux omega et <i>baseline</i> | 61 |
| 2.4 | Exercices | 64 |
| 2.5 | Solutions des exercices | 66 |
| 3 | Modèles de communication et routages | 75 |
| 3.1 | Machines parallèles à mémoire distribuée | 76 |
| 3.1.1 | Principes architecturaux des machines distribuées | 76 |
| 3.1.2 | Classification des machines | 78 |
| 3.2 | Modélisation du réseau d'interconnexion d'une machine distribuée | 82 |
| 3.2.1 | Modélisation d'un nœud du réseau | 82 |
| 3.2.2 | Les contraintes de communication | 84 |
| 3.3 | Différents modes de commutation | 85 |
| 3.3.1 | Commutation de circuits (<i>circuit-switching</i>) | 85 |
| 3.3.2 | Commutation de messages (<i>store-and-forward</i>) | 86 |
| 3.3.3 | Commutation de paquets (<i>packet-switching</i>) | 86 |
| 3.3.4 | Routage <i>wormhole</i> | 87 |
| 3.3.5 | <i>Virtual-cut-through</i> | 89 |
| 3.3.6 | Autres modes | 90 |
| 3.4 | Modélisation du temps de communication | 90 |
| 3.4.1 | Communication entre voisins | 91 |
| 3.4.2 | Communication à distance d | 94 |
| 3.5 | <i>Wormhole</i> versus <i>store-and-forward</i> | 94 |
| 3.5.1 | Technique du <i>pipeline</i> | 95 |
| 3.5.2 | Récapitulatif | 96 |
| 3.5.3 | Technique des chemins disjoints | 97 |
| 3.6 | Fonction de routage | 98 |
| 3.6.1 | Définition | 98 |
| 3.6.2 | Problèmes de congestion et routage adaptatif | 100 |
| 3.7 | Interblocage en routage statique | 101 |
| 3.7.1 | Exemple d'interblocage | 101 |
| 3.7.2 | Graphe de dépendance des canaux | 102 |
| 3.7.3 | Canaux virtuels | 104 |
| 3.7.4 | Fonctions de routage des réseaux classiques | 105 |
| 3.8 | Interblocage en routage adaptatif | 111 |
| 3.9 | Problèmes relatifs au routage | 113 |
| 3.9.1 | Indices de transmission | 114 |
| 3.9.2 | <i>Multicasting</i> et interblocage | 119 |
| 3.9.3 | Tolérance aux pannes | 120 |
| 3.10 | Exercices | 122 |
| 3.11 | Solutions des exercices | 122 |
| 4 | Communications globales | 127 |
| 4.1 | Les communications globales usuelles | 127 |
| 4.2 | Mode <i>store-and-forward</i> , temps constant | 129 |
| 4.2.1 | La diffusion | 130 |
| 4.2.2 | L'échange total | 140 |

| | | |
|----------|---|------------|
| 4.3 | Mode <i>store-and-forward</i> , temps linéaire | 145 |
| 4.3.1 | Bornes inférieures | 146 |
| 4.3.2 | La diffusion | 153 |
| 4.3.3 | Application : diffusion dans l'hypercube | 161 |
| 4.3.4 | L'échange total | 166 |
| 4.3.5 | La distribution | 169 |
| 4.3.6 | La multidistribution | 172 |
| 4.4 | Communications globales en <i>wormhole</i> | 173 |
| 4.4.1 | Critères de comparaison | 173 |
| 4.4.2 | Bornes inférieures | 174 |
| 4.4.3 | Diffusion dans l'anneau | 174 |
| 4.4.4 | Diffusion simple dans le tore | 175 |
| 4.4.5 | Diffusion optimale dans le tore pour les messages courts | 177 |
| 4.4.6 | <i>Store-and-forward</i> versus <i>wormhole</i> | 181 |
| 4.5 | Exercices | 181 |
| 4.6 | Solutions des exercices | 187 |
| 5 | Plongements | 197 |
| 5.1 | Plongements statiques | 199 |
| 5.1.1 | Définitions spécifiques | 199 |
| 5.1.2 | Résultats généraux | 201 |
| 5.1.3 | Complexité | 203 |
| 5.1.4 | Plongements de structures simples | 206 |
| 5.1.5 | Plongements dans l'hypercube | 207 |
| 5.1.6 | Autres réseaux apparentés à l'hypercube | 212 |
| 5.1.7 | Plongements dans les graphes de de Bruijn | 217 |
| 5.2 | Aspects dynamiques | 222 |
| 5.2.1 | Hypothèses pour l'émulation | 222 |
| 5.2.2 | Lien entre plongement et émulation statique | 223 |
| 5.2.3 | Emulation statique par les réseaux apparentés à l'hypercube | 225 |
| 5.2.4 | Emulation dynamique de l'hypercube par d'autres réseaux | 227 |
| 5.2.5 | Autres types d'émulations | 230 |
| 5.3 | Exercices | 231 |
| 5.4 | Solutions des exercices | 233 |
| 6 | Applications | 241 |
| 6.1 | Face à un problème à paralléliser | 242 |
| 6.1.1 | Quelle machine distribuée ? | 242 |
| 6.1.2 | Quel niveau de parallélisme ? | 243 |
| 6.1.3 | Parallélisme explicite ou parallélisme implicite | 243 |
| 6.1.4 | Paradigmes de programmation | 244 |
| 6.1.5 | Trouver le bon compromis | 244 |
| 6.2 | Utilisation de communications globales | 244 |
| 6.2.1 | Diffuser, distribuer, rassembler | 244 |
| 6.2.2 | Réductions | 247 |
| 6.2.3 | Produit matriciel | 249 |
| 6.3 | Adéquation algorithme/architecture | 253 |
| 6.3.1 | La transformée de Fourier discrète | 253 |
| 6.3.2 | La transformée de Fourier rapide (FFT) | 254 |
| 6.3.3 | Parallélisation à grains fins de la FFT | 257 |
| 6.3.4 | Parallélisation à grains moyens de la FFT | 259 |
| 6.3.5 | Parallélisation à gros grains de la FFT | 262 |
| 6.3.6 | Application | 263 |

| | | |
|-----|-----------------------------------|-----|
| 6.4 | Exercices | 264 |
| 6.5 | Solutions des exercices | 265 |

Annexes

| | | |
|----------|---|------------|
| A | Fiches descriptives de quelques réseaux | 269 |
| A.1 | Anneau ou cycle C_N | 270 |
| A.2 | Graphe complet K_N | 271 |
| A.3 | Hypercube $H(n)$ | 272 |
| A.4 | Grille $M(p_1, p_2, \dots, p_n)$ | 273 |
| A.5 | Grille torique $TM(p_1, p_2, \dots, p_n)$ | 274 |
| A.6 | Graphe de de Bruijn $B(d, D)$ | 275 |
| A.7 | Graphe de Kautz $K(d, D)$ | 276 |
| A.8 | Graphe <i>cube-connected-cycles</i> $CCC(n)$ | 277 |
| A.9 | Graphe <i>shuffle-exchange</i> $SE(n)$ | 278 |
| A.10 | Graphe <i>butterfly</i> $BF(n)$ | 279 |
| A.11 | <i>Star-graph</i> $S(n)$ | 280 |
| A.12 | Grille d'arbres d -aires $MT(d, h)$ | 281 |
| B | Des circuits VLSI aux machines parallèles | 283 |
| B.1 | Introduction | 283 |
| B.2 | Limites technologiques des VLSI | 283 |
| B.2.1 | Circuits intégrés | 284 |
| B.2.2 | Contraintes sur le nombre de liens | 285 |
| B.2.3 | Densité du réseau et bande passante globale | 287 |
| B.2.4 | Nécessité d'une mémoire distribuée | 287 |
| B.2.5 | Interface entre les processeurs et le réseau | 288 |
| B.3 | Architectures des processeurs | 288 |
| B.4 | Performances des machines | 290 |
| B.5 | Historique des machines parallèles | 294 |
| B.6 | Description des composants | 296 |
| B.6.1 | Les processeurs | 296 |
| B.6.2 | Les processeurs avec facilités de communication | 298 |
| B.7 | Les machines parallèles actuelles | 301 |
| B.7.1 | La dernière génération de machines MIMD | 301 |
| B.7.2 | La dernière génération de machines SIMD | 311 |
| B.7.3 | Les grandes lignes actuelles | 314 |
| B.8 | Perspectives | 316 |
| | Index | 321 |

Liste des figures

| | | |
|------|---|----|
| 1.1 | Arbre binaire complet de profondeur 3 | 5 |
| 1.2 | Cycle de longueur 6 | 8 |
| 1.3 | Graphe complet d'ordre 6 | 9 |
| 1.4 | Hypercube de dimension 4 représenté par niveaux | 10 |
| 1.5 | Grille $M(4, 5)$ | 12 |
| 1.6 | Grille torique $TM(4, 5)$ | 12 |
| 1.7 | Deux représentations du graphe de Petersen | 15 |
| 1.8 | Le <i>star-graph</i> $S(4)$ sur l'ensemble $\{1, 2, 3, 4\}$ | 20 |
| 1.9 | Cycle hamiltonien et cycles de toutes longueurs paires dans la grille | 25 |
| | | |
| 2.1 | Graphe de de Bruijn $B(2, 3)$ | 32 |
| 2.2 | Graphe de de Bruijn $B(3, 2)$ | 33 |
| 2.3 | Graphe de Kautz $K(2, 2)$ | 34 |
| 2.4 | Un autre exemple de graphe de Kautz : $K(2, 3)$ | 35 |
| 2.5 | Graphes non orientés $UK(2, 2)$ et $UB(2, 3)$ | 36 |
| 2.6 | Graphe <i>shuffle-exchange</i> $SE(3)$ | 43 |
| 2.7 | Graphe <i>cube-connected-cycles</i> $CCC(3)$ | 44 |
| 2.8 | Graphe <i>butterfly</i> $BF(3)$ | 45 |
| 2.9 | Grille d'arbres ternaires $MT(3, 1)$ | 46 |
| 2.10 | Liaison par bus | 50 |
| 2.11 | <i>Crossbar</i> | 51 |
| 2.12 | Etats d'un commutateur $(2, 2)$ | 52 |
| 2.13 | Réseau de Clos $CL(p, q, r)$ | 54 |
| 2.14 | Réseau de Clos $CL(3, 5, 4)$ | 56 |
| 2.15 | Schéma de la bijection et multigraphe biparti | 57 |
| 2.16 | Coloration du réseau de Clos $CL(3, 5, 4)$ | 58 |
| 2.17 | Réseau <i>butterfly</i> de dimension 3 et son graphe associé | 59 |
| 2.18 | Réseau de Beneš de dimension 3 | 61 |
| 2.19 | Réseau omega de dimension 2 | 62 |
| 2.20 | Commande du réseau omega | 63 |
| 2.21 | Réseau <i>baseline</i> de dimension 3 et son graphe associé | 64 |
| | | |
| 3.1 | Architecture SIMD | 79 |
| 3.2 | Architecture MIMD | 80 |
| 3.3 | Architecture d'un nœud du réseau | 82 |
| 3.4 | Multiplexage des liens du réseau | 83 |
| 3.5 | Commutation de circuits | 86 |
| 3.6 | Commutation de paquets | 87 |
| 3.7 | Routage <i>wormhole</i> | 88 |
| 3.8 | Routage <i>wormhole</i> sur une grille 4×4 | 89 |
| 3.9 | Effet <i>pipeline</i> | 95 |
| 3.10 | Modèle pour le routage : $G = (V = P \cup M, E = C \cup I \cup O)$ | 99 |

| | | |
|------|--|-----|
| 3.11 | Interblocage sur la grille | 101 |
| 3.12 | $D(G, R)$ pour l'anneau orienté de quatre processeurs | 103 |
| 3.13 | Canaux virtuels sur un anneau de quatre processeurs | 104 |
| 3.14 | $D(G, R)$ pour l'anneau de quatre processeurs avec des canaux virtuels | 105 |
| 3.15 | Le cycle C_8 et son graphe de dépendance | 106 |
| 3.16 | Le cycle C_8 avec trois canaux virtuels | 107 |
| 3.17 | Deux canaux virtuels dans le graphe de de Bruijn binaire | 111 |
| 3.18 | Cycle avec cordes | 115 |
| 3.19 | Si on supprime l'arête $[x, y]$, le diamètre de G devient 4 | 123 |
| | | |
| 4.1 | Un arbre de recouvrement du 5-cube | 147 |
| 4.2 | Un anneau orienté | 152 |
| 4.3 | Etiquetage des liens d'un arbre | 156 |
| 4.4 | Etiquetage de délai 2 de deux chemins hamiltoniens pairs | 160 |
| 4.5 | Arbre de recouvrement binomial de racine 000 dans le 3-cube | 162 |
| 4.6 | Trois rotations de $SBT(000)$ dans $H(3)$ | 163 |
| 4.7 | Trois arbres de recouvrement arc-disjoints de racine 000 dans $H^*(3)$ | 164 |
| 4.8 | Echange total 1-port dans l'hypercube | 168 |
| 4.9 | Un arbre équilibré dans le 5-cube | 171 |
| 4.10 | Diffusion <i>wormhole</i> simple dans le tore | 176 |
| 4.11 | Un tore 5×5 dessiné comme un carré Π_1 | 177 |
| 4.12 | Une croix X_1 | 178 |
| 4.13 | Un tore 25×25 dessiné comme un carré Π_2 | 179 |
| 4.14 | Graphe de Petersen | 184 |
| 4.15 | Algorithmes de diffusion dans le graphe de Petersen | 190 |
| 4.16 | Arbres de recouvrement arc-disjoints dans le graphe de Petersen | 191 |
| | | |
| 5.1 | Plongement f du graphe G dans le graphe H | 200 |
| 5.2 | Arbre binaire complet double-racine $BDR(3)$ | 209 |
| 5.3 | Plongement d'un arbre binaire complet double-racine dans $H(5)$ | 209 |
| 5.4 | Graphe <i>Fast-Fourier-Transform</i> $FFT(3)$ | 213 |
| 5.5 | Plongement de $SE(3)$ dans $UB(2, 3)$ | 214 |
| 5.6 | Relations entre les réseaux apparentés à l'hypercube | 215 |
| 5.7 | Plongements de réseaux dans l'hypercube | 217 |
| 5.8 | Plongements de réseaux dans le graphe de de Bruijn | 220 |
| 5.9 | Résultats sur les plongements statiques avec expansion optimale | 221 |
| | | |
| 6.1 | Schéma logique de division et de calcul pour $n = 8$ | 255 |
| 6.2 | Calcul d'une FFT sur huit échantillons | 256 |
| 6.3 | Calcul d'une FFT inverse sur huit échantillons | 257 |
| 6.4 | FFT sur $\frac{n}{2} \log_2 n$ processeurs | 258 |
| 6.5 | FFT sur hypercube : plongement d'un arbre binaire | 259 |
| 6.6 | FFT sur hypercube : schéma d'exécution sur l'arbre binaire | 260 |
| 6.7 | FFT sur hypercube : liens utilisés par étape | 261 |
| 6.8 | FFT sur graphe de de Bruijn | 262 |
| | | |
| A.1 | Cycle de longueur 6 | 270 |
| A.2 | Graphe complet d'ordre 6 | 271 |
| A.3 | Hypercube de dimension 4 | 272 |
| A.4 | Grille $M(4, 5)$ | 273 |
| A.5 | Grille torique $TM(4, 5)$ | 274 |
| A.6 | Graphe de de Bruijn $B(2, 3)$ sur l'alphabet (binaire) $\{0, 1\}$ | 275 |
| A.7 | Graphe de Kautz $K(2, 2)$ sur l'alphabet $\{0, 1, 2\}$ | 276 |

| | | |
|------|--|-----|
| A.8 | Graphe <i>cube-connected-cycles</i> $CCC(3)$ | 277 |
| A.9 | Graphe <i>shuffle-exchange</i> $SE(3)$ | 278 |
| A.10 | Graphe <i>butterfly</i> $BF(3)$ | 279 |
| A.11 | <i>Star-graph</i> $S(4)$ | 280 |
| A.12 | Grille d'arbres ternaires $MT(3, 1)$ | 281 |
| A.13 | Grille d'arbres binaires $MT(2, 2)$ | 282 |
| | | |
| B.1 | Différents niveaux d'encapsulation en VLSI | 284 |
| B.2 | <i>Fat-tree</i> du réseau de données de la CM-5 à 16 processeurs | 302 |
| B.3 | Un contrôleur du <i>fat-tree</i> | 303 |
| B.4 | Réseau d'interconnexion de la CS-2 à 32 processeurs | 305 |
| B.5 | Anneau d'anneaux de la KSR | 307 |
| B.6 | Réseau de la Paragon | 309 |
| B.7 | Réseau d'interconnexion du SP1 à 64 processeurs | 310 |

Auteurs

Jean-Claude BERMOND - Directeur de Recherches CNRS
I3S Sophia-Antipolis

Pierre FRAIGNIAUD - Chargé de Recherches CNRS
LIP-IMAG Lyon

Anne GERMA - Professeur
ENST Paris

Marie-Claude HEYDEMANN - Professeur
LRI Orsay

Emmanuel LAZARD¹ - Allocataire de Recherches
LRI Orsay

Philippe MICHALLON - Allocataire de Recherches
LMC-IMAG Grenoble

André RASPAUD² - Maître de Conférences
LaBRI Bordeaux

Dominique SOTTEAU - Directeur de Recherches CNRS
LRI Orsay

Michel SYSKA³ - Allocataire de Recherches
I3S Sophia-Antipolis

Denis TRYSTRAM - Professeur
LMC-IMAG Grenoble

avec la collaboration de

Eric DARROT - Allocataire de Recherches, *I3S Sophia-Antipolis*

Daniel LAFAYE DE MICHEAUX - Maître de Conférences, *I3S Sophia-Antipolis*

¹E. Lazard est maintenant Maître de Conférences au LAMSADE (Paris-Dauphine).

²A. Raspaud est aujourd'hui Professeur.

³M. Syska a obtenu un poste de Maître de Conférences.

Avant-propos

Calcul parallèle et communications

L'importance des machines massivement parallèles — systèmes pouvant comporter plusieurs milliers de processeurs — n'est plus à démontrer aujourd'hui ; ces machines répondent à un besoin sans cesse croissant de puissance de calcul.

De ce fait, de nombreux domaines de recherche nouveaux ont émergé ; par exemple, les langages parallèles, la parallélisation automatique, les architectures distribuées, les environnements de programmation, l'analyse de complexité parallèle et l'algorithmique parallèle.

La plupart des machines existantes reposent sur des architectures à mémoire distribuée. Une des différences essentielles entre les machines traditionnelles, dites de *Von Neumann*, et les réseaux de processeurs que sont les machines à mémoire distribuée, est la gestion des communications. Par exemple, la distribution des données, le placement des processus sur les différents processeurs, et la gestion des mouvements de données requis par les calculs, sont autant de problèmes qui font intervenir des communications. Les communications sont d'autant plus importantes que la granularité du parallélisme est fine. En effet, si l'on parallélise un programme grossièrement, c'est-à-dire en identifiant des procédures indépendantes, peu de communications sont nécessaires, et le traitement reste en majorité séquentiel. Au contraire, en distribuant le plus possible les traitements élémentaires, on obtient le parallélisme potentiellement le plus efficace, mais les communications induites peuvent alors être pénalisantes. Optimiser les communications est donc la clef de l'efficacité des algorithmes parallèles.

Les chercheurs en parallélisme étudient donc les problèmes liés aux communications afin d'accélérer les traitements, tant du point de vue du logiciel que de l'architecture du matériel. Il est indispensable d'aborder cette étude de façon globale. En effet, une connaissance approfondie des modèles et des mécanismes élémentaires de communication des machines est requise pour proposer des solutions efficaces, applicables aux systèmes réels.

Mais l'étude des communications fait appel à des connaissances variées : en architecture (circuits intégrés, cartes et liens de communication), en algorithmique (communications globales, routage, interblocage, mise en œuvre d'appli-

cations, etc.) et en théorie des graphes (coloration, plongement, cheminement, etc.). Cette grande diversité montre combien il est difficile d'avoir une vision synthétique du domaine. La synthèse des résultats actuels exposée dans ce livre, tend à offrir au lecteur un cadre général pour l'étude de divers aspects des systèmes parallèles : architecture, environnements, langages, algorithmique et programmation.

Jusqu'à présent, les communications ont été étudiées de manière dispersée par des chercheurs ne s'intéressant qu'à des problèmes très spécifiques (développement, dans un modèle donné, d'un schéma de communication donné, sur une architecture donnée). L'objet de ce livre est au contraire de présenter de manière unifiée les recherches sur les communications, sans pour autant en gommer les particularités.

Description de l'ouvrage

Ce livre est organisé en six chapitres, complétés par deux annexes.

Dans le chapitre 1 sont définis les termes de théorie des graphes utilisés dans cet ouvrage. On y trouve les définitions de familles de graphes classiques utilisés dans les architectures (anneau, graphe complet, hypercube, grille, etc.), ainsi que leurs propriétés élémentaires. Ce chapitre contient aussi une étude théorique du problème du nombre maximum de sommets d'un graphe de degré maximum et de diamètre fixés.

Le chapitre 2 présente d'autres graphes, bien adaptés pour servir de support aux réseaux (graphes de de Bruijn et de Kautz, graphe *shuffle-exchange*, etc.). Ces graphes ont un petit degré et un diamètre logarithmique en le nombre de sommets. Le chapitre 2 introduit aussi les réseaux multiétages qui, contrairement aux précédents, ont la particularité de ne pas être des réseaux statiques mais de pouvoir réaliser, grâce aux commutateurs qui en sont la base, un grand nombre de permutations des entrées sur les sorties.

Les modèles de communication et les différents types de routage sont abordés au chapitre 3. Une classification des machines parallèles et les principes architecturaux des machines à mémoire distribuée y sont donnés. De même, les principaux modes de commutation sont présentés (*circuit-switched*, *store-and-forward*, *wormhole*, *virtual cut-through*, etc.). Ceci conduit à une modélisation des réseaux et des temps de communication (temps constant et linéaire). Les modes de communications principaux sont comparés. Enfin, le problème de l'interblocage lors du routage est posé et quelques solutions y sont apportées.

Le chapitre 4 décrit les principaux types de communications globales (diffusion, échange total, distribution et rassemblement). Ce chapitre est découpé en trois parties, suivant les principales hypothèses décrites au chapitre précédent. Les deux premières parties traitent du mode *store-and-forward*, sous l'hypothèse temps constant, puis sous l'hypothèse temps linéaire. La troisième partie traite du

mode de routage *wormhole*. Dans chacune des parties sont présentés des bornes inférieures et des algorithmes déduits de méthodologies d'étude des différents protocoles. Ces résultats sont systématiquement illustrés : sur l'hypercube en mode *store-and-forward* et sur la grille en mode *wormhole*.

Le problème du plongement de graphes dans le contexte des communications est motivé dans la première partie du chapitre 5. La deuxième partie, la plus importante du chapitre, traite de la notion classique de plongement dans l'aspect statique. Les principaux résultats connus sont donnés (complexité, paramètres des plongements de structures simples dans l'hypercube et les graphes de de Bruijn) et résumés dans le tableau final. La dernière partie aborde les aspects dynamiques de la notion de plongement.

Au chapitre 6 sont présentées deux applications typiques du calcul numérique utilisant des communications : le produit matriciel et le calcul de la transformée de Fourier discrète. Ces deux exemples illustrent la mise en œuvre de schémas de communication présentés dans les chapitres précédents.

Dans l'annexe A, une série de fiches synthétiques sur les réseaux usuels récapitule pour chacun ses propriétés principales (ordre, degré, nombre d'arêtes ou d'arcs, etc.).

L'annexe B commence par un rapide rappel sur les architectures des processeurs et une définition du vocabulaire nécessaire à la compréhension de tout document traitant des performances d'un processeur ou d'une machine parallèle. Puis, les processeurs les plus utilisés ainsi que les machines parallèles les plus récentes sont présentés. La description de ces machines est principalement axée sur leur réseau d'interconnexion et leurs propriétés liées aux communications.

Chaque chapitre contient une bibliographie et est illustré par des exercices dont énoncés et solutions sont regroupés à la fin du chapitre.

Quelques conventions et notations

Toutes les définitions, théorèmes, propositions, corollaires, etc. sont référencés indifféremment, dans leur ordre d'apparition au fil du texte. Un théorème est un résultat de portée générale (non nécessairement difficile à démontrer). Une proposition, au contraire, concerne des cas particuliers. Nous avons remplacé la plupart du temps les termes anglo-saxons par leur traduction française, cependant, les termes non traduits sont écrits en italiques. Nous rappelons ci-dessous des notations habituelles :

- $|X|$ désigne le cardinal de l'ensemble X .
- $\ln n$ désigne le logarithme népérien de n , et $\log_b n$ désigne le logarithme en base b de n .

- C_n^p et $\binom{n}{p}$ représentent le nombre de combinaisons de p éléments, pris parmi n éléments.
- Soient f et g deux fonctions réelles positives de variable entière :
 - La notation $f = O(g)$ signifie qu'il existe un réel positif c et un entier M tel que, pour tout x , $x \geq M$, on a $f(x) \leq cg(x)$.
 - La notation $f = \Omega(g)$ signifie qu'il existe un réel positif c et un entier M tel que, pour tout x , $x \geq M$, on a $f(x) \geq cg(x)$.
 - La notation $f = \Theta(g)$ est équivalente à $f = O(g)$ et $g = O(f)$.

Jean de Rumeur

Cet ouvrage est écrit par une équipe de chercheurs spécialistes de domaines différents mais complémentaires (théorie des graphes, placement et ordonnancement, algorithmique numérique, architecture, etc.). Ces auteurs ont uni leurs compétences, en tirant profit de plusieurs années de collaboration au sein de l'opération RUMEUR du PRC C³. Ce livre fait suite à un support de cours rédigé pour une école d'été, organisée en août 1992 à l'Institut d'Etudes Scientifiques de Cargèse, pour des étudiants de troisième cycle. Il a été modifié et complété grâce aux nombreuses remarques et critiques faites au cours de cette école, ainsi que dans les cours de DEA dispensés par les auteurs.

Remerciements

Les auteurs tiennent à remercier particulièrement :

- leurs laboratoires de recherche, qui ont mis à leur disposition les moyens matériels nécessaires à la réalisation de cet ouvrage ;
- le Programme de Recherches Coordonnées C³, qui a financé les réunions de RUMEUR qui ont permis aux auteurs de collaborer régulièrement ;
- le CIMPA, qui a aidé à organiser l'école d'été de Cargèse en 1992 ;
- les collègues et les étudiants qui les ont aidés dans la tâche ingrate de relecture, en particulier : Dominique Barth, Christophe Calvin, Charles Delorme, Odile Favaron, Rabah Harbane, Jean-Claude König, Stéphane Pérennes et Jean-François Saclé.