

Régulation de charge

Jean-Claude König – Jean-Louis Roch

21 octobre 1999

Résumé

De manière générale, que ce soit dans un cadre dynamique ou statique, l'ordonnement d'un programme parallèle prend en entrée un graphe de tâches et un modèle de machine et fournit en sortie un ordonnancement à un facteur borné (aussi proche de 1 que possible) de l'optimal.

Dans le cadre dynamique (ou à la volée), le graphe est inconnu car construit au cours de l'exécution; la durée d'une tâche ou d'une communication n'est connue qu'après sa terminaison. Dans le cadre statique, le graphe est supposé connu avant l'exécution du programme. Dans cet article, nous introduisons les différents modèles de machines considérés et les techniques d'ordonnement qui leur sont associées dans les cadres statique et dynamique.

1 Introduction

Le placement des tâches sur l'architecture cible dans le but de minimiser la durée d'exécution est un problème complexe qui intègre un grand nombre de paramètres qui seront étudiés dans cette école et qui sont liés à des choix algorithmiques (routage par exemple) ou des contraintes physiques (mode de communication, bande passante par exemple).

Le placement et l'ordonnement des tâches doivent donc être décomposés en plusieurs étapes pour limiter à la fois la complexité et augmenter la portabilité et la généricité. La démarche classique est donc de placer et d'ordonner les tâches sur une machine virtuelle et d'émuler ensuite la machine virtuelle par la machine cible. Or il semble difficile (impossible) de concevoir une machine virtuelle simple et universelle qui soit acceptable (comme l'est la machine de von Neumann pour le calcul séquentiel) c'est-à-dire qui puisse servir de référence pour mesurer la complexité. Une raison clé de cet échec est l'impossibilité de faire abstraction des transferts de données (de coûts non uniformes de part la hiérarchie mémoire) pour obtenir des résultats de complexité ayant un sens expérimental.

Ainsi les modèles (de type PRAM synchrone [34] ou asynchrone [14, 4]) permettent d'analyser le coût d'un algorithme à partir de paramètres caractéristiques des graphes associés à ses exécutions : nombre d'opérations (les tâches sont séquentielles), volume d'accès en mémoire partagée (les arêtes de dépendance de données correspondent potentiellement à des accès distants), . . . Mais pour que les résultats ainsi obtenus puissent être exploitables et validés expérimentalement, il est nécessaire de montrer comment l'ordonnement du

graphe (indépendant du nombre de processeurs) peut être réalisé. Il s'agit d'une part de montrer qu'un ordonnancement optimal du programme apporte l'efficacité et d'autre part d'exhiber un algorithme permettant de calculer un ordonnancement optimal ou à un facteur de l'optimal (α -compétitif). Ces deux études sont complémentaires; ce chapitre est surtout centré sur la présentation d'algorithmes α -compétitifs. Dans ce cadre, un ordonnancement optimal pour un programme donné peut éventuellement être séquentiel.

L'ordonnancement trivial sur une machine PRAM théorique (principe de Brent) ne peut fournir une émulation raisonnable sur une machine, même simple, car il repose sur des hypothèses fausses comme l'indépendance entre les temps de latence pour accéder à la mémoire partagée et la possibilité de répartir équitablement sans surcoût de temps l'ensemble des instructions d'une étape sur les processeurs.

Une façon de rendre le modèle plus fiable est de montrer la possibilité d'ordonner des classes de graphes spécifiques sur des modèles de machine plus réalistes, comme :

- granularité globale : le nombre d'accès potentiels en mémoire distante est négligeable devant le volume de calcul;
- granularité grossière : toutes les tâches sont de même durée et ont chacune un volume d'accès distants négligeable devant leur nombre d'opérations;
- topologie particulière : le graphe de dépendance est un arbre dont les tâches sont de même coût.

De tels graphes sont suffisamment génériques pour être effectivement rencontrés dans des applications irrégulières : les classes précédentes englobent en effet de nombreux cas où par exemple durées de tâches et volumes d'accès à des données distantes sont inconnus ou variables, ou encore le graphe est inconnu et construit au cours de l'exécution par la création de nouvelles tâches.

Nous nous proposons dans ce chapitre de présenter les résultats connus pour l'ordonnancement dynamique et statique de tels graphes sur un modèle simpliste et d'étudier leur évolution lorsqu'on intègre au modèle certains paramètres caractérisant l'architecture cible comme :

- le nombre de processeurs
- la latence liée aux communications
- l'hétérogénéité des processeurs

Il est clair qu'un modèle intégrant seulement ces trois paramètres pour être à peu près valide nécessite de nombreuses hypothèses du type : seule l'application considérée s'exécute sur le système, la bande passante est suffisante pour ne pas entraîner une saturation des débits . . . Malgré cela un tel modèle devrait permettre de mieux analyser le parallélisme exploitable d'une application et de donner un ordonnancement plus efficace à l'exécution.

Organisation. Après la définition d'un problème d'ordonnancement, la première section précise les mécanismes systèmes nécessaires à sa mise en œuvre ainsi que les critères permettant d'évaluer ses performances. Puis, les algorithmes classiques d'ordonnancement lorsque les processeurs sont supposés identiques et les communications non prises en compte, sont rappelés. L'analyse de leurs performances permet d'introduire les techniques de base utilisées pour l'évaluation théorique d'un algorithme d'ordonnancement. Dans les sections suivantes,

les principaux algorithmes permettant de prendre en compte les communications sont présentés : ordonnancements dynamiques (section 3) et ordonnancements statiques (section 4). Enfin, les principaux résultats concernant l'ordonnement sur des machines hétérogènes sont présentés.

2 Le problème d'ordonnement

Dans cette partie, nous définissons le formalisme de graphe le plus général pris en entrée d'un algorithme d'ordonnement. L'exécution d'un programme parallèle décrit (de manière dynamique dans le cadre le plus général) un graphe de dépendance. Ce graphe est orienté sans cycle : les nœuds correspondent à des blocs d'instructions (tâches) et les arêtes à des relations de dépendance sur l'accès à une donnée (transitions).

2.1 Caractérisation du problème : graphe de flot de données

De manière générale, l'exécution d'un programme parallèle peut-être décrite par un graphe de flots de données (macro data-flow) : DFG . Ce graphe est biparti, orienté et sans cycle avec comme ensemble de nœuds $J = \{j_1, \dots, j_n\}$ correspondant aux blocs d'instructions séquentielles (j signifiant *job*, nous dirons tâche) et $T = \{t_1, \dots, t_m\}$ correspondant aux opérations de lecture et écriture d'une donnée (vue comme un bloc de mots) en mémoire globale (t signifiant *transition*).

Tous les nœuds d'entrée (i.e. de degré entrant nul) et de sortie (i.e. de degré sortant 0) sont des transitions; les premiers correspondent aux données en entrée du programme (fichiers en lecture par exemple) alors que les derniers correspondent à ses sorties.

Le graphe définit un ordre partiel sur J , dit de précédence et noté \prec . Pour tout $1 \leq i, k \leq n$, $j_i \prec j_k$ si et seulement s'il existe un chemin orienté de j_i vers j_k .

N.B. Dans toute la suite, nous nous limitons à l'ordonnement de graphes *d'arité bornée* (qui est celui des algorithmes dits EREW). Dans le cas général (algorithme CRCW), nous supposons donc que le graphe (le programme) a été transformé en une description sous forme graphe de tâches à arité bornée. Une telle transformation peut augmenter d'un facteur quadratique au plus le nombre de sommets (tâches et transitions) du graphe.

Contraintes de placement. Généralement, des fonctions de contraintes $\phi_t : T \mapsto \mathcal{P}(\{M_i, 1 \leq i \leq p\})$ et $\phi_j : J \mapsto \mathcal{P}(\{P_i, 1 \leq i \leq p\})$ permettent de spécifier des contraintes de localisation pour certaines transitions ou pour certaines tâches. Une transition en entrée peut par exemple n'être disponible que sur un module mémoire, une tâche ne peut être exécutée que sur un sous-ensemble de processeurs. Dans le cadre de cet article, nous nous limitons au cas où aucune contrainte n'est spécifiée, i.e. $\phi_t(x) = (\{M_i, 1 \leq i \leq p\})$ et $\phi_j(y) = (\{P_i, 1 \leq i \leq p\})$ pour tout $x \in T$, $y \in J$.

Déterminisme des exécutions. Le graphe de dépendance est directement associé à une exécution; il ne peut souvent être complètement déterminé qu'après exécution complète du

programme. Cependant, dans tous les cas, nous supposons que l'exécution reste *déterministe*: le graphe de dépendance peut être construit complètement dynamiquement, mais reste le même pour les mêmes valeurs en entrée quel que soit l'ordonnancement effectué (même séquentiel). Cette contrainte est liée à la difficulté (impossibilité si $P \neq NP$) de trouver un ordonnancement polynômial pour un algorithme de coût polynômial sur une machine non déterministe (i.e. dont toute exécution admet un ordonnancement polynômial).

2.2 Construction distribuée du graphe de flot de données

Séparer la régulation de charge de la programmation de l'application requiert alors un langage de programmation capable de décrire ce flot, à un niveau macroscopique. Différents langages ont été proposés qui permettent de décrire de manière plus ou moins implicite un tel graphe. Ainsi BSP [35] permet de décrire le flot par étapes synchrones; ce graphe est très partiel: lors d'une étape, seules les tâches de l'étape suivante, supposées toutes indépendantes, sont annoncées ainsi que leurs dépendances avec les tâches de l'étape courante.

Le langage Cilk [20] permet de décrire de manière des graphes de tâches série-parallèle: les dépendances de données entre tâches ne sont pas connues mais le graphe de contrôle est entièrement déterminé.

Les langages Jade [32] et Athapascan-1 [12] permettent quant à eux de décrire un graphe de flot de données à un niveau macroscopique et dans un cadre assez général. Dans ces deux langages, une tâche correspond à une séquence d'instructions séquentielles telle que les accès effectués sur des données globales sont annoncés en début de bloc (accès en lecture, ou accès en écriture ou accès en mis-à-jour). Ainsi, les dépendances de flot de données entre tâches sont déterminées avant même l'exécution des tâches, par interprétation abstraite.

Les différences entre ces deux langages se situent au niveau de la syntaxe et du modèle de coût. D'un point de vue syntaxique, une tâche correspond à un bloc d'instructions; en Athapascan-1, elle correspond à un appel de procédure, ce qui offre d'une part une meilleure lisibilité du code et une plus facile réutilisabilité. D'un point de vue complexité, le coût de construction (implicite dans le langage) de la construction du graphe de flot de données est prouvé proportionnel à la taille de ce graphe, même en environnement distribué. Cette propriété est fondamentale pour la définition d'un modèle de coût associé au langage [12]. Ce modèle permet d'exhiber des ordonnancements offrant des compromis en temps mais aussi en mémoire [7, 5, 12].

2.3 Ordonnancement du graphe sur une machine

Une architecture distribuée est vue comme un ensemble d'unités de calcul P_i chacune associée à un module de mémoire M_i . L'ordonnancement consiste alors à affecter à chaque transition t (qui correspond à une donnée "en assignation unique") un module mémoire $M(t)$ et à chaque job j le processeur $P(j)$ sur lequel il sera exécuté. L'ordonnancement spécifie en outre quels événements conditionnent le démarrage ou la préemption d'une tâche (généralement fin de communication et fin de tâche de calcul).

2.3.1 Ordonnancement global et ordonnancement système

Deux niveaux d'ordonnancement peuvent être distingués :

- *macroscopique* : ordonnancement global des nœuds du graphe de dépendance (tâches et données) sur les ressources (processeurs et modules mémoire).
- *microscopique* : ordonnancement local des processus concurrents sur le processeur et des instructions atomiques pour minimiser les temps d'attente (par exemple anticipation du chargement d'une ligne dans le cache ou avancement des communications).

Nous ne nous intéressons dans ce chapitre qu'à l'ordonnancement macroscopique, l'ordonnancement microscopique étant supposé réalisé localement sur chaque processeur par le système.

Calculer un ordonnancement consiste donc à allouer à chaque nœud du graphe (donnée ou tâche) une ressource en complétant la relation de précédence \prec initiale. Chaque exécution du programme après ordonnancement définit une relation d'ordre partiel \prec_i sur l'ensemble des différentes instructions atomiques exécutées. Dans tous les cas, \prec_i doit être compatible avec \prec : l'ordonnancement calculé est dit "admissible" si, pour toute exécution et pour tout couple d'instructions atomiques (i_1, i_2) appartenant respectivement aux jobs j_1 et j_2 alors :

$$(j_1 \prec j_2) \implies (i_1 \prec_i i_2).$$

2.3.2 Opérations utilisées pour réaliser l'ordonnancement

Pour réaliser et contrôler l'ordonnancement calculé, l'ordonnanceur utilise certaines primitives supposées fournies par le système. Sauf précisé, toutes ces primitives sont supposées de coût constant.

Exécution d'une tâche . Une tâche (un nœud du DFG) peut être exécutée sur un processeur quelconque.

Communication. Un processeur peut lire ou écrire un objet dans n'importe quel module mémoire distant. Le coût de l'accès à n données est supposé borné par $h.n$ où h est une constante de la machine [34, 11]. Les accès mémoire sont asynchrones. Pour contrôler l'ordonnancement (via les transitions, une tâche accède en lecture à des données écrites par ses prédécesseurs), des outils de synchronisation globaux sont supposés disponibles (verrous, sémaphores).

Duplication. La duplication consiste à exécuter plusieurs fois une même tâche sur différents processeurs; ceci peut permettre de construire des ordonnancements minimisant les coûts de communication.

Migration. Certains algorithmes d'ordonnancement utilisent en outre une opération de préemption, qui permet de suspendre une tâche au cours de son exécution pour la relancer, éventuellement plus tard, sur un autre processeur [33, 3]. Deux types de préemption sont distingués :

- *migration par redémarrage* : lorsqu'une tâche est déplacée sur un autre processeur, son exécution reprend depuis son début.

- *migration* : lorsqu’une tâche est déplacée sur un autre processeur, son exécution reprend à partir de sa dernière instruction exécutée.

Un algorithme d’ordonnancement *non-préemptif* n’utilise aucune de ces opérations : il perd tout contrôle sur une tâche dès qu’il l’a démarrée sur un processeur, sauf information sur sa terminaison. Le coût d’une opération de migration est supposé constant.

Remarque. L’utilisation de la migration est très liée au modèle de programmation (comment sont implémentées les tâches) et au modèle de machine. Par exemple, si les tâches sont sans blocage (c’est le cas dans le graphe de dépendance), la migration est inutile lorsque la machine est composée de processeurs identiques. Par contre, elle peut s’avérer utile sur des machines hétérogènes [33].

2.4 Notation des problèmes d’ordonnancement

Nous allons utiliser pour classer les problèmes d’ordonnancement la notation introduite dans [17] et étendue dans [39] pour les problèmes avec communications. Nous l’étendons pour préciser le cas où le graphe de dépendance est construit en cours d’exécution, les coûts des tâches et des communications étant inconnus jusqu’à leur terminaison.

Cette notation est constituée de trois champs : α, β, γ . Ces trois paramètres permettent de régler le niveau d’abstraction de la machine virtuelle et les hypothèses sur le type de graphes considérés.

Le paramètre α définit le type des processeurs et leur nombre. Pour nous :

- $\alpha = P$: signifie que les processeurs sont identiques et que leur nombre m est une entrée du problème.
- $\alpha = \overline{P}$: signifie que les processeurs sont identiques, leur nombre étant suffisant ou non borné.
- $\alpha = P_m$: indique un nombre fixé de processeurs identiques.

Le paramètre β définit le type du graphe de précedence, les temps d’exécution des tâches, le coût de communication, l’autorisation ou pas de la duplication, l’autorisation ou pas de la préemption.

Donc $\beta = \beta_1\beta_2\beta_3\beta_4$ où :

- $\beta_1 \in \{prec?, prec, arbre, chaine, \dots, \bullet\}$: Type du graphe de précedence des tâches.
 - $\beta_1 = prec$: le graphe de précedence des tâches est quelconque mais de structure connue;
 - $\beta_1 = prec?$: le graphe de précedence des tâches est quelconque mais inconnu (construit en cours d’exécution);
 - $\beta_1 = arbre$: le graphe de précedence des tâches est un arbre;
 - $\beta_1 = \bullet$: les tâches sont indépendantes.
- $\beta_2 \in \{com, com?, c_{jk}, c, \bullet\}$: communication.
 - $\beta_2 = com$: temps de communication donnés par le graphe de précedence des tâches.

- $\beta_2 = com?$: temps de communication inconnu avant fin de la communication effective mais constant.
- $\beta_2 = com? < \alpha \sum p_i$: temps de communication inconnu mais volume total de données communiquées borné par une constante fois le volume total de calcul.
- $\beta_2 = c_{jk}$: le temps de communication entre la tâche j et la tâche k est c_{jk} .
- $\beta_2 = c$: temps de communication partout égaux.
- $\beta_2 = \bullet$: temps de communications nuls.
- $\beta_3 \in \{dup, mig, \bullet\}$: la migration et la duplication sont autorisées ou non
- $\beta_4 \in \{p_i = 1, p_i = ?, \bullet\}$: durée des tâches;
 - $\beta_4 = (p_i = 1)$: toutes les tâches ont une durée d'exécution unitaire.
 - $\beta_4 = (p_i = ?)$: les tâches sont de durées quelconques et inconnue avant exécution.
 - $\beta_4 = \bullet$: les durées des tâches sont définies par le graphe de précédence des tâches.

Le paramètre γ définit l'objectif. Dans ce chapitre, le but sera toujours de minimiser la durée d'exécution, ce but est noté $\gamma = C_{max}$.

2.5 Évaluation du coût d'un algorithme d'ordonnancement

Dès que l'on considère des machines virtuelles non triviales, le problème de trouver un ordonnancement optimal d'un graphe devient NP-difficile. On s'intéresse donc à rechercher des heuristiques polynômiales (voire linéaire) pouvant fournir un ordonnancement proche de l'optimal.

L'efficacité d'une heuristique peut-être évaluée de diverses façons. Dans ce cours, nous nous contenterons de mesurer la qualité d'un algorithme d'ordonnancement par deux critères :

- **ratio** (ou compétitivité): le maximum du rapport sur tous les graphes entre la durée obtenue par l'ordonnancement de l'heuristique et la durée obtenue par un ordonnancement optimal.
- **surcoût** d'ordonnancement: le nombre d'opérations requises pour la mise en œuvre effective de l'ordonnancement.

2.5.1 Ratio de performance ou compétitivité

Définition 1 Soit h un algorithme d'ordonnancement. Le ratio de performance de h , noté $R(h)$, est le maximum du rapport sur tous les graphes entre la durée donnée par l'ordonnancement fourni par h et la durée de l'ordonnancement optimal, c'est-à-dire $R(h) = \text{Max} \frac{C_{max}(h)}{C_{max}^{opt}}$.

On dit qu'un algorithme h est ρ -approché (ou encore ρ -compétitif) si $R(h) \leq \rho$.

Un algorithme h est d'autant plus efficace que son ratio $R(h)$ est petit. La démarche dans le cas d'un problème NP-complet est donc de trouver la meilleure heuristique possible. En fait la plupart du temps nous ne sommes capables que de donner une borne supérieure du

ratio de la meilleure heuristique (en étudiant le ratio d'une heuristique connue) et une borne inférieure de ce ratio.

Une technique classique pour obtenir une borne inférieure du ratio est basée sur le résultat suivant, facile à démontrer et appelé "théorème de l'impossibilité" [25].

Théorème 1 *Etant donné un problème d'ordonnancement et un entier c , si la question de savoir si un graphe G peut-être ordonnancé en c unités de temps ou moins est NP-complet alors on ne peut pas espérer avoir une heuristique pour ce problème d'ordonnancement ayant un ratio inférieur à $\frac{c+1}{c}$.*

La valeur du ratio de performance permet de cerner la complexité pratique du problème et de traduire l'augmentation de la complexité du problème en fonction des nouveaux paramètres intégrés au modèle.

D'un point de vue pratique, l'heuristique de meilleur ratio n'est pas forcément celle à intégrer si le surcoût induit par sa mise en œuvre s'avère trop important.

2.5.2 Surcoût d'ordonnancement

Définition 2 *Soit h un algorithme d'ordonnancement et G un graphe de dépendance quelconque. Le surcoût d'ordonnancement de h , noté $\sigma(h)$, est le nombre d'opérations effectuées par h pour exécuter G .*

Il est souvent difficile de séparer a priori le surcoût lié à la mise en œuvre de l'ordonnancement de la longueur de l'ordonnancement lui-même (caractérisé par le ratio). Aussi, son étude est basée sur l'analyse du ratio et complétée par la prise en compte des opérations pour l'ordonnancement.

Le surcoût peut être séparé en deux parties: d'une part les opérations requises pour le calcul de l'ordonnancement et d'autre part celles utilisées pour exécuter le graphe de tâches en garantissant le respect de la relation de précédence.

Cadre statique. Dans le cadre statique, le graphe étant connu à la compilation, le calcul de l'ordonnancement n'implique pas de surcoût. Par contre son contrôle nécessite un surcoût lié – au moins – à la taille du graphe. La réalisation d'un ordonnancement pré-calculé étant triviale, nous considérerons uniquement le ratio de performance pour l'analyse des algorithmes d'ordonnancement statique.

Cadre dynamique. Dans le cadre dynamique, au surcoût de contrôle s'ajoute aussi celui du calcul de l'ordonnancement. Dans tous les cas, chaque nœud et chaque transition devant être placés, il est borné inférieurement par la taille du graphe.

Notation. Pour exprimer le surcoût d'ordonnancement, dans la suite nous désignerons par N_a le nombre de tâches et par N_t le nombre de transitions. Comme tous les nœuds sont supposés d'arité bornée, nous désignerons par $n = N_a + N_t$ le nombre de nœuds dans G et par n_d le degré de G . Une borne sur la taille du graphe est donc nn_d .

3 Ordonnements sans communication

Nous présentons d'abord les résultats préliminaires concernant l'ordonnement (dynamique ou statique) sans communication.

3.1 Problème $P_m | prec, p_i = ? | C_{max}$: algorithme de liste

Dans le cas de la régulation dynamique de la charge, l'ordonnement le plus rencontré en pratique consiste basiquement, lorsqu'un processeur devient inactif, à lui allouer une tâche prête s'il en existe. Un tel ordonnancement est appelé "algorithme de liste" : il requiert de gérer une structure évoluant dynamiquement et permettant d'accéder aux tâches prêtes.

Dans cette section, seul le ratio de performance d'un tel algorithme pour un graphe de précedence quelconque est étudié. Le surcoût induit par sa mise en œuvre sur une architecture distribuée (gestion du graphe et contrôle de l'ordonnement d'une part, introduction des communications d'autre part) est étudié dans la section 4.

Théorème 2 [15] *Si les communications ne sont pas prises en compte, tout algorithme d'ordonnement de type liste a un ratio de performance borné par $(2 - \frac{1}{m})$ sur une machine constituée de m processeurs identiques.*

Nous rappelons la preuve de ce résultat car son schéma intervient dans de nombreuses démonstrations.

Soit T_m la longueur de l'ordonnement fourni par l'algorithme de liste sur m machines. Un algorithme de liste est tel que, à chaque instant, au moins un processeur exécute une tâche. Ainsi, si à un instant donné un processeur est inactif alors il existe au moins un processeur qui exécute une tâche. Soit t_{j_1} l'une des tâches qui s'est terminée à la date T_m et soit d_{j_1} la date de début d'exécution de t_{j_1} . Deux cas peuvent être distingués :

cas 1 : soit aucun processeur n'a été inactif avant d_{j_1} .

cas 2 : soit il existe au moins un processeur inactif à un certain instant avant d_{j_1} . Soit θ la plus grande date avant d_{j_1} à laquelle au moins un processeur est inactif. À θ , t_{j_1} n'est pas prête car sinon elle aurait été affectée à un processeur inactif. Donc il existe une tâche t_{j_2} telle que t_{j_2} est en cours d'exécution à θ et $t_{j_2} \prec t_{j_1}$. Soit d_{j_2} la date de début d'exécution de t_{j_2} .

En appliquant récursivement ce schéma jusqu'à ce que le cas 1 arrive, nous construisons une séquence de tâches $t_{j_k} \prec \dots \prec t_{j_2} \prec t_{j_1}$ telle que, à tout instant, soit tous les processeurs sont actifs soit un processeur exécute une tâche t_{j_i} ($1 \leq i \leq k$).

Soit T_∞ le temps minimal sur un nombre infini de processeurs et T_1 le nombre total d'opérations exécutées (temps sur un processeur séquentiel, appelé aussi travail). Le temps total d'inactivité $\#I$ est défini par $\#I = mT_m - T_1$. Pour $1 \leq i \leq k$, soit l_i la durée de la tâche t_{j_i} . Nous avons $\#I \leq (m - 1) \sum_{i=1}^k l_i$, d'où :

$$pT_m \leq T_1 + (p - 1) \sum_{i=1}^k l_i.$$

En outre, les tâches t_{j_i} , $1 \leq i \leq k$ étant sur un chemin critique, $\sum_{j=1}^k l_i \leq T_\infty$. On a finalement :

$$T_m \leq \frac{T_1}{m} + \left(1 - \frac{1}{m}\right) T_\infty \quad (1)$$

Par ailleurs, soit T_m^* la longueur de l'ordonnancement optimal. On a $T_\infty \leq T_m^*$ et, comme T_1 opérations doivent être exécutées par tout ordonnancement, $mT_m^* \geq T_1$. Remplaçant dans 1, on obtient : $T_m \leq \left(2 - \frac{1}{m}\right) T_m^*$. \square

Le théorème précédent 2 est donné dans une version restreinte. [3]. De manière plus générale, la borne 1 reste inchangée lorsque la relation de précedence \prec' considérée par l'algorithme de liste est plus forte que celle \prec utilisée pour calculer l'ordonnancement optimal. La preuve est directe puisque on aura encore $t_{j_k} \prec \dots \prec t_{j_2} \prec t_{j_1}$. Clairement, la même remarque s'applique si la durée des tâches est augmentée artificiellement. Cela implique que ni l'ajout de contraintes de précedence telles que des barrières de synchronisation pour obtenir un graphe de dépendance bien structuré (par exemple semi-parallèle), ni l'insertion artificielle d'opérations de découpe pour que toutes les tâches aient la même durée ne peuvent améliorer le ratio d'un algorithme d'ordonnancement en ligne.

Remarque : principe de Brent constructif. En corollaire de l'encadrement 1, on obtient une démonstration constructive générale du principe de Brent pour des tâches de longueur arbitraire (et non forcément unitaires).

Corollaire 1 *Si un programme parallèle effectue T_1 opérations (temps séquentiel) et peut s'exécuter en temps parallèle minimal T_∞ sur un nombre infini de processeurs, alors il peut être exécuté sur m processeurs en temps*

$$\text{Max} \left\{ \left\lceil \frac{T_1}{m} \right\rceil, T_\infty \right\} \leq T_m \leq \left\lceil \frac{T_1}{m} + T_\infty \right\rceil \quad (2)$$

3.2 Bornes inférieures sur le ratio

Une question naturelle est alors de déterminer s'il est possible d'obtenir un ratio de performance inférieur à $\left(2 - \frac{1}{m}\right)$, soit sur le même problème soit en utilisant des opérations plus puissantes pour réaliser l'ordonnancement, comme le tirage aléatoire ou la migration. Ce problème est étudié dans [33] où la proposition suivante est prouvée.

Théorème 3 [33] *Si les temps de communication ne sont pas pris en compte, $\left(2 - \frac{1}{m}\right)$ est une borne inférieure pour le ratio d'un algorithme d'ordonnancement déterministe avec ou sans migration.*

Un algorithme probabiliste sans préemption ne peut avoir un ratio de performance inférieur à $\left(2 - \frac{1}{\sqrt{m}}\right)$, ratio qui est obtenu par un algorithme de liste dans lequel les tâches affectées aux processeurs inactifs sont tirées au hasard parmi l'ensemble des tâches prêtes.

Seule la preuve de la première partie est ici présentée car elle utilise la construction d'un adversaire, technique fréquemment utilisée pour obtenir des bornes inférieures sur un algorithme à la volée. Il s'agit de construire une instance du problème réalisant le pire cas et de montrer que ce pire cas ne peut être évité par aucun algorithme d'ordonnancement.

Le pire cas est obtenu pour l'instance suivante due à Graham [15]. G contient $1 + p(p - 1)$ tâches indépendantes; une tâche α est de longueur p tandis que les β_k , $1 \leq k \leq p(p - 1)$ autres sont de longueur 1. L'ordonnancement optimal est de longueur p ; il est obtenu en exécutant α_1 sur un processeur et les $p(p - 1)$ tâches β_k sur les $p - 1$ autres processeurs.

Aucun algorithme d'ordonnancement à la volée ne peut faire d'hypothèse sur la longueur d'une tâche qui est inconnue tant que la tâche n'est pas terminée. La technique de l'adversaire consiste alors à faire en sorte que la tâche α soit démarrée le plus tard possible. Chaque fois que l'ordonnanceur lance l'exécution d'une tâche sur un processeur, nous supposons donc que cette tâche est une tâche β_k , de longueur 1. Ainsi l'ordonnanceur ne peut demander l'exécution de la tâche α qu'après que toutes les tâches β_k aient été exécutées, donc au plus tôt au top $(p - 1)$. La longueur de l'ordonnancement délivré est alors de longueur au moins $(p - 1) + p$, ce qui montre la borne inférieure. Il est clair que la préemption ou la migration ne peuvent améliorer ce ratio. \square .

En conséquence, ni la migration ni l'introduction d'aléatoire ne peuvent améliorer le ratio de performance en comparaison d'un algorithme de liste.

3.3 Problème : $P|indep, p_j|C_{max}$

Si la durée des tâches est supposée connue, trier la liste (par exemple en ordonnant en priorité les tâches les plus longues) ne permet pas d'améliorer la borne de Graham. Par contre, il est possible d'obtenir une amélioration en restreignant la classe de graphes considérés.

Ainsi si les tâches sont indépendantes le problème reste NP-complet (réduction au problème de la partition) mais l'heuristique LPTF – *largest processing time first* – (algorithme de liste avec la liste des tâches triée en ordre décroissant de durée) permet d'obtenir un ratio de $\frac{4}{3} - \frac{1}{3m}$ [16].

Une restriction plus intéressante est celle où toutes les tâches ont la même durée, ce sous problème étudié dans la section suivante servira de base de comparaison lors de l'étude de l'impact des communications.

3.4 Problème : $P|prec, p_j = 1|C_{max}$

Théorème 4 *Il n'existe pas d'algorithme d'ordonnancement pour le problème $P|prec, p_j = 1|C_{max}$ ayant un rapport de performance inférieur à $\frac{4}{3}$ si $P \neq NP$.*

Ce résultat est déduit du résultat suivant par le théorème de l'impossibilité :

Lemme 1 [26] *Le problème de décider pour un graphe de précédence quelconque s'il existe un ordonnancement valide de longueur au plus 3 est NP-complet.*

Nous allons utiliser une réduction du problème CLIQUE qui est un problème connu NP-complet. On rappelle qu'une *clique* un graphe dont tous les sommets sont adjacents. Il est évident que si $G_1 = (V_1, E_1)$ est une clique et si $G_2 = (V_2, E_2)$ a le même nombre de sommets que G_1 ($|V_1| = |V_2|$) mais n'est pas une clique, alors $|E_1| > |E_2|$.

Le problème CLIQUE est le suivant :

- **Instance :** Un graphe non-orienté $G = (V, E)$ et un entier k .
- **Question :** Existe-t-il dans G une clique de taille k ?

Nous réduisons une instance de CLIQUE à une instance de $P|prec, p_j = 1|C_{max}$ de la manière suivante (la figure 1 illustre la réduction). Soit $l = \frac{(k-1)k}{2}$ le nombre d'arêtes dans une clique de taille k ; Soit $k' = |V| - k$ et $l' = |E| - l$.

On considère l'instance suivante de $P|prec, p_j = 1|C_{max}$:

- Le nombre de machines est $m = \max\{k, l + k', l'\} + 1$.
- Pour chaque sommet $v \in V$, on introduit une tâche T_v et pour chaque arête $e \in E$ une tâche L_e . On a $T_v \prec L_e$ si v est une extrémité de e .
- On ajoute $(3m - l - k - l' - k')$ tâches $X_x (x = 1, \dots, m - k)$, $Y_y (y = 1, \dots, m - l - k')$ et $Z_z (z = 1, \dots, m - l')$ avec les précédences $X_x \prec Y_y \prec Z_z, \forall x, y, z$.

Le nombre total de tâches est $3m$. Un ordonnancement valide avec $C_{max} = 3$ ne peut donc pas avoir de temps d'inactivité.

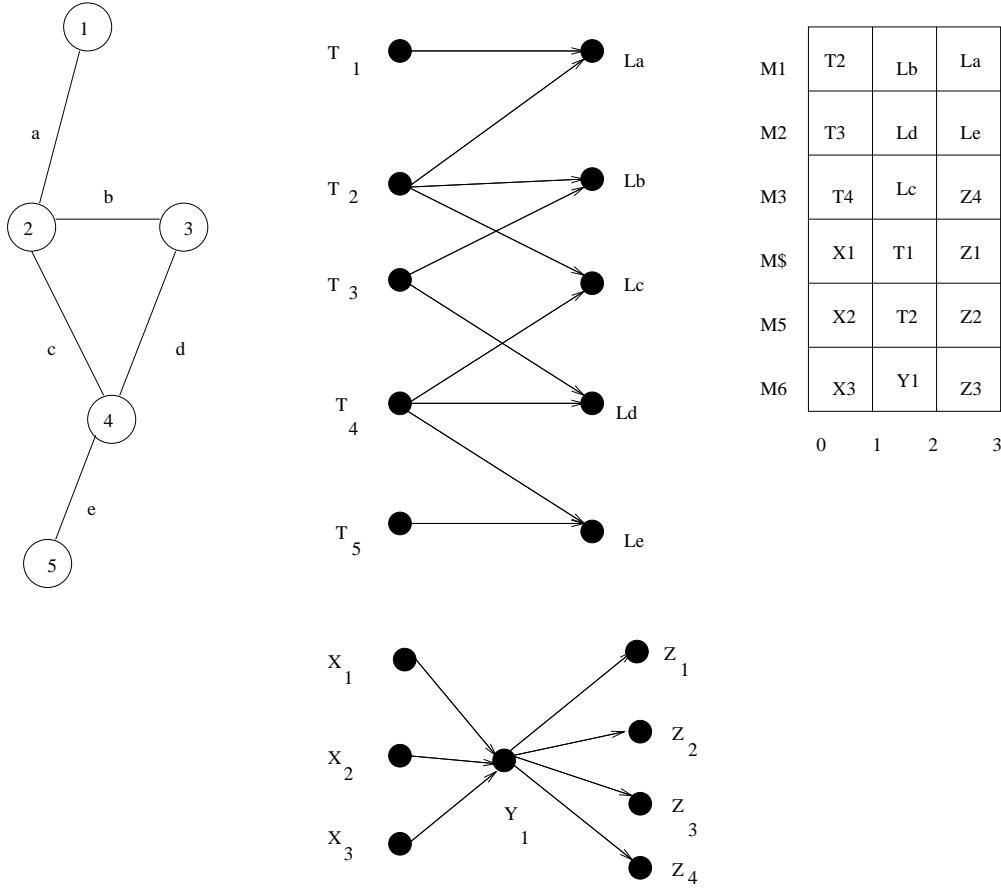


FIG. 1 – Le graphe de précédence des tâches, l'instance correspondante du problème d'ordonnancement où $m = 6$, et un ordonnancement en trois unités de temps.

- Supposons que G contient une clique de taille k , alors il existe un ordonnancement valide de longueur 3 (voir figure 1).
- Supposons qu'on a un ordonnancement valide de longueur 3, alors G contient une clique de taille k .

Il est évident que les tâches qui appartiennent à des chemins de longueur 3 doivent être exécutées au plus tôt. Ainsi X_x , Y_y et Z_z sont exécutés aux instants 1, 2 et 3 respectivement.

A l'instant 1 on peut exécuter au plus k tâches de type T . Supposons que seulement $k_1 < k$ tâches forment une clique. Ceci dit qu'il y aura certains processeurs qui vont rester inactifs pendant la deuxième unité de temps puisque les tâches de type L (c'est-à-dire les arêtes adjacentes aux k sommets exécutés à l'instant 1) qui sont libérées sont moins de l . Par conséquent, l'ordonnancement ne peut pas terminer à $t = 3$. \square

3.5 Une borne supérieure pour $P|prec, p_j = 1|C_{max}$: l'algorithme de Coffman-Graham

La solution du problème utilise un algorithme d'étiquetage qui associe à chaque tâche du graphe un entier. L'étiquetage se fait de la façon suivante [21]. R désigne l'ensemble des tâches non étiquetées dont tous les successeurs immédiats sont déjà étiquetés.

- **Étape 1 :** On distribue de façon arbitraire des étiquettes aux tâches qui n'ont pas de successeurs.
- **Étape 2 :** A chacune des tâches de R on fait correspondre la séquence décroissante des étiquettes de ses successeurs immédiats S_i . La tâche qui a la plus petite (lexicographiquement) séquence S_i prend l'étiquette suivante.
- **Étape 3 :** On itère l'étape 2 jusqu'à ce que toutes les tâches de G soient étiquetées.
- **Étape 4 :** Construction d'une liste de priorité en ordre décroissant des étiquettes.

Théorème 5 [24] *L'algorithme de CG a un ratio de performance borné par*

$$2 - \frac{2}{m}$$

La preuve assez technique de ce résultat est omise. Elle permet d'obtenir le corollaire suivant.

Corollaire 2 *L'algorithme de CG est optimal pour $P2|prec, p_j = 1|C_{max}$.*

Etant donné que la complexité en temps de l'algorithme de CG est en $O(n^2)$, il est évident qu'on a intérêt de l'utiliser seulement pour des petites valeurs de m . Pour des grandes valeurs de m les gains par rapport à un algorithme de liste quelconque deviennent minimes.

Remarque. La complexité du problème $P_m|prec, p_j = 1|C_{max}$ n'est pas connue quand $m \geq 3$.

4 Ordonnements dynamiques avec communications

Nous nous intéressons à la mise en œuvre d'un algorithme de liste sur une architecture distribuée qui fournit un ordonnancement de longueur presque optimale sur un nombre arbitraire de processeurs identiques (cf corollaire 1). L'intérêt d'un tel algorithme est d'être compatible avec une création dynamique du graphe de dépendance : une tâche t peut créer

lors de son exécution de nouvelles tâches t_k , en nombre quelconque, en ajoutant de nouvelles dépendances. Il est donc particulièrement adapté à des applications irrégulières où la génération de parallélisme est souvent réalisée en cours d'exécution.

Deux problèmes sont alors à considérer :

- l'implémentation de l'algorithme de liste lui-même : gestion du graphe et contrôle de l'ordonnancement par affectation de tâches prêtes aux processeurs inactifs.
- la prise en compte du coût des communications.

Les résultats présentés sont appliqués à l'ordonnancement de programmes dont le nombre d'opérations est grand devant le nombre de tâches et de dépendances (*granularité globale*).

Remarque : Hypothèse sur la construction à la volée du graphe de tâches. Pour faciliter la gestion du graphe de dépendance (qui induit une relation de précédence entre les tâches), nous supposons que :

- les tâches t_k créées ont une relation de précédence avec la tâche t qui les crée, c'est à dire : $\forall k$, si t crée t_k alors $t \prec t_k$;
- les nouvelles dépendances de données ajoutées au graphe ne peuvent concerner que les tâches ayant préalablement une relation de précédence avec t , c'est à dire : soit t_1 et t_2 deux tâches et $t_1 \prec t_2$ une dépendance ajoutée par t , alors, avant l'ajout de cette dépendance, on avait $t \prec t_1$ et $t \prec t_2$.

Ces contraintes, bien que non restrictives en ce qui concerne la description de l'exécution d'un algorithme sous forme de graphes de tâches (cf les chapitres sur la construction de graphes de tâches), ne sont pas fondamentales : leur introduction est motivée par l'obtention d'une implémentation concise en ce qui concerne l'algorithme de listes.

4.1 Coût des communications : simulation d'une mémoire partagée

Pour prendre en compte le coût des communications (contention et latence), nous simulons une mémoire partagée. La simulation utilise des fonctions de hachage (choisies aléatoirement dans une classe de fonctions de hachage universelles) pour distribuer les objets accessibles par tous les processeurs (i.e. les transitions dans le graphe de dépendance) sur les différents modules mémoire [31, 23].

Le délai d'une simulation est le temps requis pour l'accès à une donnée de taille unité. Il est lié à l'évaluation de la fonction de hachage, à la contention mémoire (quand plusieurs accès sont effectués vers un même module) et à la latence (temps de routage). Dans [31], une simulation de délai $\Theta(\log p)$ sur un réseau "butterfly" est donnée pour des accès exclusifs (de type EREW). Dans [22], une simulation de délai $O(\log \log p \log^* p)$ est donnée pour des accès aléatoires (de type EREW ou CRCW) sur une architecture disposant d'un réseau d'interconnexion complet (la latence n'est pas prise en compte); pour des accès concurrents (CRCW), la simulation est à un facteur $\log^* p$ de l'optimal.

Pour obtenir une simulation optimale, ces délais doivent être masqués par des calculs. La technique classique (virtualisation ou "parallel slackness" [23, 36, 22]) consiste à simuler q ($q > p$) processeurs virtuels sur les p processeurs de l'architecture physique, en utilisant généralement des processus légers. La simulation est alors dite optimale si le délai pour un accès est proportionnel à q/p .

Dans la suite de cette section, nous supposons disposer d'une simulation optimale de délai h ; le nombre de processeurs m requis pour une telle simulation peut alors être supérieur au nombre de processeurs p de l'architecture. L'ordonnancement sera donc effectué sur m processeurs et non p , mais comparé à l'ordonnancement optimal sur p processeurs. Un processeur (parmi les m) peut lire ou écrire un objet dans n'importe quel module mémoire distant. Le coût de l'accès à une donnée de taille n est supposé borné par $h.n$ où h est une constante de la machine.

Synchronisation. Les accès en mémoire distante ne permettent pas directement de synchroniser les processeurs (à la différence des communications). Sur une architecture distribuée asynchrone, des outils spécifiques de synchronisation (verrous, sémaphores) doivent alors être utilisés. Nous supposons aussi que le coût d'accès à un verrou global, lorsque celui-ci est libre, est lui aussi borné par h (il requiert un aller-retour).

4.2 Problème $P_m | prec?, p_i = ? | C_{max}$: algorithme de liste

Lorsque le graphe de dépendance est construit et déterminé en cours d'exécution, l'implémentation d'un algorithme de liste nécessite de gérer une structure permettant de déterminer les tâches prêtes (i.e. dont les prédécesseurs sont terminés) et d'affecter des tâches distinctes à chacun des processeurs devenus inactifs.

Directe en séquentiel, cette affectation est cependant complexe en distribué. Dans [9], un algorithme basé sur des graphes d'expansion est utilisé pour ordonnancer n tâches indépendantes effectuant au total $T_1 = O(n)$ opérations en temps $O(\log n)$ sur $\frac{n}{\log n}$ processeurs d'une EREW PRAM; outre le caractère théorique d'un tel ordonnancement (qui suppose une synchronisation globale à chaque instruction), la constante masquée dans le O est loin d'être négligeable (> 237) devant le ratio de performance (proche de 2)¹.

Dans ce qui suit, nous considérons dans un but simplificateur une gestion séquentielle de la structure permettant de calculer les tâches prêtes (i.e. la liste) et de les affecter aux processeurs inactifs. À défaut d'être la plus performante sur un grand nombre de processeurs, elle permet d'estimer précisément le surcoût lié à la gestion de l'ordonnancement. Pour éviter de dédier un processeur à la gestion de cette structure, elle est stockée en mémoire (virtuellement) partagée et accédée par les processeurs en exclusion mutuelle; l'exclusion est réalisée par un verrou global.

Théorème 6 *Soit un programme parallèle dont l'exécution génère n tâches et n_d dépendances (transitions). Soit T_1 le temps d'une exécution séquentielle (sur un processeur) et T_∞ le temps minimal sur un nombre infini de processeurs. Alors, si les communications ne sont pas prises en compte, G peut être ordonnancé sur m processeurs identiques en temps T_m borné par :*

$$T_m \leq \frac{T_1}{m} + T_\infty + O(n + n_d + m).$$

1. Une constante beaucoup plus faible peut être obtenue si l'ordonnancement est réalisé sur moins de processeurs. Par exemple, l'utilisation d'un calcul de préfixe permet de calculer un ordonnancement à un facteur 2 de l'optimal sur $\frac{n}{\log^2 n}$ processeurs [8].

Comme pour le théorème 2, la preuve est basée sur la construction d'un ensemble de tâches situées sur un chemin critique dans G . Pour contrôler l'ordonnancement, nous utilisons une structure globale, visible par tous les processeurs et accédée en exclusion mutuelle grâce à un verrou global **verrou-graphe**; outre le graphe G , cette structure contient :

- la liste L_t des tâches prêtes (initialisée avec les tâches de G de degré entrant nul);
- la liste L_i des processeurs inactifs; à l'initialisation, seul un processeur (disons P_0) est supposé artificiellement actif – il simule la terminaison d'une tâche fictive –, les autres $m - 1$ processeurs étant mis dans L_i .

L'ordonnancement des tâches est géré de la façon suivante : lorsqu'un processeur p termine une tâche t (par exemple P_0 au top 0), il exécute les instructions de *post-traitement* liées à t suivantes :

1. Prendre le verrou **verrou-graphe**.
2. Mettre à jour G en ajoutant les nouvelles tâches créées et les dépendances associées.
3. Mettre à jour la liste L_t des tâches prêtes en ajoutant les tâches successeurs immédiats de t dont tous les prédécesseurs sont terminés.
4. Tant que ($L_i \neq$ vide) et ($L_t \neq$ vide) faire
 - Enlever une tâche t' de L_t et un processeur p' de L_i ;
 - Lancer l'exécution de t' sur p' .
5. Si ($L_t \neq$ vide) alors
 - Enlever une tâche t'
 - Relâcher le verrou **verrou-graphe**;
 - Exécuter t' (sur le processeur p)
6. sinon
 - Ajouter p à la liste des processeurs inactifs
 - Relâcher le verrou **verrou-graphe**;

De manière évidente, l'ordonnancement ainsi construit est admissible, i.e. ne viole aucune contrainte de précédence dans G . Soit T_m le temps d'exécution avec cet ordonnancement sur m machines. En reprenant le schéma de la preuve précédente (2), il s'agit de borner les tops d'inactivité *potentielle*. Les tops $\{1, \dots, T_m\}$ sont partitionnés en trois sous-ensembles, A (activité), O (ordonnancement) et I (inactivité) de la façon suivante :

- à tout top dans A , tous les processeurs exécutent une tâche, i.e. une instruction de l'application.
- à tout top dans O , une instruction de post-traitement (i.e. gestion de l'ordonnancement) est en cours.
- à tout top dans I , la liste L_i des processeurs inactifs est non vide mais aucune instruction de post-traitement n'est en cours d'exécution.

En remarquant que, lorsqu'un processeur est en attente du verrou, une instruction de post-traitement est nécessairement en cours sur un autre processeur, il n'y a pas dans $A \cup C$ de tops où un processeur est en attente du verrou.

$\#A$ peut être trivialement borné par $\frac{T_1}{m}$. L'exécution en exclusion mutuelle des instructions de post-traitement permet de borner facilement $\#O$. En effet, une fois le verrou pris, le post-traitement ne consiste qu'en la mise à jour des listes (les opérations de base sont en coût

constant) et au lancement de nouvelles exécutions éventuelles. Lorsque le verrou est libre, le temps nécessaire à sa prise est borné par h . Lorsqu'il n'est pas libre, il existe nécessairement un autre processeur en train d'exécuter une instruction de post-traitement (éventuellement en train de prendre ou de relâcher le verrou). Par suite, le nombre de tops dans $\#O$ est borné par $O(h(n + n_d + m))$ (i.e. la taille du graphe).

$\#I$ est borné de manière analogue à la preuve de Graham, en remarquant que lorsqu'un processeur est inactif et aucun en cours de post-traitement, il existe nécessairement une tâche sur un chemin critique qui est en cours d'exécution. Par suite, $\#I < T_\infty$.

Pour conclure, $mT_m = m\#O + m\#A + (m - 1)\#I$. D'où² $T_m \leq \frac{T_1}{m} + T_\infty + O(n + n_d + m)$.
□

Conclusion sur le surcoût d'ordonnancement. Borner le surcoût requiert d'analyser l'algorithme qui contrôle l'ordonnancement. Dans tous les cas, ce surcoût fait apparaître un terme lié à la taille du graphe de dépendance. Pour le diminuer, on peut s'intéresser à structurer ce graphe; mais les remarques énoncées à la section précédente montre que cela ne peut qu'entraîner une perte en terme de ratio.

Implémentation d'un algorithme de liste. En pratique, la gestion centralisée des listes et l'accès en exclusion mutuelle est évité. D'une part, le graphe est stocké de manière distribuée; d'autre part, lorsqu'un processeur est inactif, il scrute les autres processeurs pour trouver une nouvelle tâche à exécuter (on parle de "vol" de tâche). Une telle stratégie est analysée en moyenne d'un point de vue théorique dans [6, 20] et permet d'obtenir un résultat similaire à 6.

4.3 Problème $P_m | prec?, p_i =?, com? | C_{max}$: communications

L'algorithme précédent ne prend pas en compte le délai entre l'affectation d'une tâche prête à un processeur et le moment où elle démarre son exécution, les données dont elle dépend étant disponibles localement sur ce processeur. De part la technique de simulation utilisée pour la mémoire partagée (§4.1), ce délai est borné par h fois le volume des données accédées. Le résultat s'étend donc directement.

Théorème 7 *Soit une machine simulant m processeurs identiques accédant avec un délai h à une mémoire virtuellement partagée.*

Soit un programme parallèle dont l'exécution génère n tâches et n_d dépendances (transitions). Soit T_1 le temps de l'exécution séquentielle (sur un processeur), W_c le volume total maximal d'accès distants et T_∞ le temps minimal sur un nombre infini de processeurs (avec communications). Alors, G peut être ordonnancé sur m processeurs identiques en temps T_m borné par :

$$T_m \leq \frac{T_1}{m} + h \frac{W_c}{m} + T_\infty + hO(n + n_d + m).$$

Le résultat s'obtient à partir de la preuve précédente (théorème 6) en bornant $m\#A$ par $T_1 + W_c$ et $\#I$ par T_∞ (qui prend en compte les communications sur le chemin critique). □

2. Le ratio $(2 - 1/m)$ de Graham se retrouve directement dans le cas où $\#O = 0$.

4.4 Granularité globale : $P_p | prec?, p_i =?, com? < \alpha \sum p_i | C_{max}$

Le théorème précédent ne permet pas d'évaluer le ratio de performance par rapport à l'ordonnancement optimal. D'une part, ce dernier peut éventuellement supprimer toutes les communications et d'autre part, la simulation de délai optimal utilisé requiert un nombre de processeurs m supérieur au nombre de processeurs p physiques. Le temps effectif sur p processeurs est alors $\frac{m}{p}T_m$; il ne peut être garanti à un facteur constant de l'optimal que si hW_c négligeable devant T_1 .

Lorsque le volume maximal de communications est borné par αT_1 (problème $P_p | prec?, p_i =?, com? < \alpha \sum p_i | C_{max}$), l'algorithme de liste précédent (théorème 7) permet d'obtenir un ordonnancement de longueur bornée par :

$$\frac{m}{p} \left((1 + \alpha h) \frac{T_1}{m} + T_\infty + hO(n + n_d + m) \right).$$

Si pour une entrée de taille suffisante, les valeurs de α (borne supérieure du rapport entre volume de communications et volume de calcul), T_∞ (temps parallèle minimal) et $(n + n_d)$ (nombre de tâches et de dépendances) sont asymptotiquement négligeables devant le temps séquentiel d'exécution, le ratio de performance est asymptotiquement optimal (i.e. 1). D'où le théorème suivant.

Théorème 8 *Soit un réseau de processeurs constitué de p processeurs identiques communiquants via un réseau d'interconnexion.*

Soit π un programme parallèle dont l'exécution construit un graphe de dépendance avec un volume de dépendances de données asymptotiquement négligeable devant le temps séquentiel T_1 et un nombre d'arêtes de dépendances asymptotiquement négligeable devant le temps parallèle sur un nombre infini de processeurs. Alors π peut être exécuté en temps T_p asymptotiquement optimal par un algorithme de liste :

$$T_p \leq (1 + \epsilon) \frac{T_1}{p}.$$

On obtient donc un ordonnancement asymptotiquement optimal pour des problèmes très parallèles et requérant peu de communications. Ce cadre est plus général que celui de la granularité grossière (certaines tâches pouvant requérir un temps de communication grand devant le temps de calcul [13]). Il est à noter que la création dynamique de tâches et la gestion du graphe n'implique pas de surcoût pour autant que la taille du graphe des tâches reste négligeable devant le temps parallèle optimal.

5 Ordonnements statiques avec communications

En reprenant le même modèle de communication qu'au paragraphe 4.1, nous considérons que le temps de communication entre deux tâches ayant une dépendance directe est nul si les deux tâches sont exécutées sur le même processeur et est de durée constante et connue (proportionnelle à la taille de la donnée communiquée) sinon. Soit d_i la date de début de la

tâche t_i sur le processeur P_i ; si (i, j) est un arc du graphe de précédence et P_i le processeur qui exécute la tâche i , alors

$$\begin{cases} d_j - d_i \geq p_i, & \text{si } P_i = P_j \\ d_j - d_i \geq p_i + c_{ij}, & \text{si } P_i \neq P_j \end{cases}$$

où c_{ij} est le coût de la communication interprocesseur entre les tâches i et j .

5.1 Granularité grossière : $\bar{P}|prec, c = 1, p_j = 1|C_{max}$

Théorème 9 *Le problème de décider si une instance de $\bar{P}|prec, c = 1, p_j = 1|C_{max}$ a un ordonnancement de longueur au plus 6 est NP-complet.*

La preuve est basée sur une transformation par le problème 3-SAT (voir [19]).

Corollaire 3 *Il n'existe pas d'heuristique ρ -approchée telle que $\rho \leq \frac{7}{6}$.*

Ainsi nous savons que nous sommes dans l'impossibilité de trouver une heuristique qui permette de garantir que quelque soit le graphe en donnée la solution obtenue par l'heuristique est à moins de 16% de la solution optimale.

Ce problème avec $c = 0$ était un problème facile. Une heuristique triviale pour $c = 1$ est d'intercaler entre 2 étapes de calcul une phase de communication. Cette heuristique a trivialement un ratio égal à deux. Récemment dans [29], les auteurs donnent une heuristique de ratio $\frac{4}{3}$. Cette heuristique est basée sur la résolution d'un problème de programmation linéaire.

En effet, soit $G=(V,E)$ le graphe représentant une instance du problème $\bar{P}|prec, c = 1, p_j = 1|C_{max}$. Le problème peut-être formulé par le problème linéaire en nombres entiers suivant :

$$\begin{aligned} & \text{Min } w \\ & \forall i \in V, \quad t_i \geq 0 \\ & \forall (i, j) \in E, \quad x_{(i,j)} \in \{0, 1\} \quad (*) \\ & \forall (i, j) \in E, \quad t_i + 1 + x_{(i,j)} \geq t_j \\ & \forall i \in V, \quad \sum_{j \in \Gamma^+(i)} x_{(i,j)} \geq |\Gamma^+(i)| - 1 \\ & \forall i \in V, \quad \sum_{j \in \Gamma^-(i)} x_{(j,i)} \geq |\Gamma^-(i)| - 1 \\ & \forall i \in V, \quad t_i + 1 \leq w \end{aligned}$$

Avec cette formulation $x_{(i,j)} = 0$ implique que les tâches i et j sont exécutées sur le même processeur. Lorsqu'on relâche les contraintes d'intégrité (c'est-à-dire on remplace $(*)$ par $\forall (i, j) \in E, \quad x_{(i,j)} \in [0, 1]$) on obtient un problème de programmation linéaire qui peut être résolu en temps polynômial.

Intuitivement, dans la solution réelle obtenue plus $x_{(i,j)}$ est proche de 0 plus on a intérêt à supprimer la communication entre la tâche i et la tâche j . Dans [29], les auteurs montrent

qu'effectivement en mettant à 0 tous les $x_{(i,j)}$ strictement inférieurs à $\frac{1}{2}$ et à 1 les autres, on obtient une solution réalisable à moins de 33% de la solution réelle (qui est elle-même une borne inférieure de la solution optimale).

Théorème 10 *L'heuristique proposée est un algorithme $\frac{4}{3}$ -approché.*

Dans [10], il est montré que ce problème devient facile si on autorise la duplication des tâches.

5.2 Le problème $P|prec, c = 1, p_j = 1|C_{max}$

Théorème 11 *Le problème de décider si une instance de $P|biparti, c = 1, p_j = 1|C_{max}$ a un ordonnancement de longueur au plus 4 est NP-complet.*

La preuve est omise; elle utilise une réduction du problème de la clique [1].

Ainsi nous savons que nous sommes dans l'impossibilité de trouver une heuristique qui permettent de garantir que quelque soit le graphe en donnée la solution obtenue par l'heuristique est à moins de 25% de la solution optimale pour le problème $P|prec, c = 1, p_j = 1|C_{max}$.

Une solution classique pour obtenir un ordonnancement sur un nombre limité de processeurs est de partir d'un ordonnancement sur un nombre illimité de processeurs et de plier cet ordonnancement sur le nombre de processeurs considéré. Le théorème suivant permet d'évaluer une telle démarche.

Théorème 12 *A partir de n'importe quelle heuristique de rapport de performance ρ pour le problème $\bar{P}|prec, c = 1, p_j = 1|C_{max}$ on peut obtenir une heuristique pour le problème $P|prec, c = 1, p_j = 1|C_{max}$ de rapport de performance $1 + \rho$.*

Soit l'heuristique h suivante : on exécute les tâches dans le même ordre que pour l'heuristique, notée h^* dans la suite, utilisée pour le problème $\bar{P}|prec, c = 1, p_j = 1|C_{max}$. Plus précisément, si X_i est l'ensemble des tâches exécutées à l'instant i avec l'heuristique h^* , et β_i le cardinal de X_i , alors dans h , si $\beta_i \neq 0$, on exécute les tâches de X_i en $\lceil \frac{\beta_i}{m} \rceil$ instants où m est le nombre de processeurs. Dans le cas où $\beta_i = 0$, on conserve l'instant d'inactivité.

On remarque que le nombre d'instants où il y a au moins un processeur inactif est inférieur à la durée de l'ordonnancement avec un nombre illimité de processeurs. Les autres instants, tous les processeurs sont actifs. Ainsi l'écart entre le temps donné par l'heuristique et le temps optimal est borné par la durée de l'ordonnancement avec un nombre illimité de processeurs. Si l'on note t^* la durée de l'ordonnancement donné par h^* et t^m la durée de l'ordonnancement donné par h alors $t^m \leq t^* + t_{opt}^m$. Comme par hypothèse $t^* \leq \rho t_{opt}^*$ et comme $t_{opt}^* \leq t_{opt}^m$, la conclusion s'impose. \square

Le corollaire immédiat de ce théorème est l'existence d'une heuristique de ratio de performance $\frac{7}{3}$ pour le problème $P|prec, c = 1, p_j = 1|C_{max}$. Cette heuristique est la meilleure connue à l'heure actuelle.

Dans les sections suivantes nous allons citer quelques familles de graphes ayant un comportement spécifique vis-à-vis de ce modèle.

5.3 Le problème $P2|arbre, p_j = 1, c = 1|C_{max}$

Une famille des graphes de précédence qui a un intérêt important est la famille des arbres. Depuis 1961, on sait que le problème $P|arbre, p_j = 1|C_{max}$ est un problème polynômial. En fait, T.C. Hu a proposé un algorithme simple qui consiste à construire une liste de priorité où la tâche dont la distance de la racine (le niveau) de l'anti-arbre est la plus grande, est la tâche la plus prioritaire.

Par contre, quand on prend en compte les communications Lenstra et al [28, 38] ont montré que le problème $P|arbre, p_j = 1, c = 1|C_{max}$ devient NP-complet. Lorsqu'on fixe le nombre de processeurs, un algorithme optimal de programmation dynamique de complexité $O(n^{2(m-1)})$ a été présenté dans [37].

Restrictions sur la structure du graphe. Parmi les autres classes de graphes pour lesquels un ordonnancement optimal est facilement calculable figure les graphes ordonnés par intervalles. Un graphe de dépendance est dit ordonné par intervalles si à toute tâche t peut être associé un intervalle réel I_t tel que $t \prec t'$ si et seulement si $\forall x \in I_t, \forall y \in I_{t'} on a $x < y$. Pour de tels graphes, le calcul d'un ordonnancement optimal peut être ramené à un tri [30].$

Étant donné un graphe quelconque, le problème est alors d'ajouter un nombre minimal de dépendances qui, sans limiter de manière trop importante le parallélisme, permettent de le ramener à une famille de graphes pour laquelle le calcul de l'ordonnancement optimal est facile.

5.4 Granularité fine : $P|prec, c, p_j = 1|C_{max}$

Sur ce modèle, les résultats connus sont très décevants et de nombreux travaux théoriques et pratiques sont en cours.

Les meilleurs résultats connus pour la borne inférieure du ratio de performance est $1 + \frac{1}{c+3}$ [2]. Dès que c devient grand et donc que la granularité est relativement fine, cette borne ne représente aucun intérêt. Réciproquement les meilleures heuristiques connues donnent un ratio de performance égal à $c + 2$. Là aussi, la borne ne représente aucun intérêt puisque ce ratio est le même qu'une heuristique triviale consistant à prendre un ordonnancement négligeant les communications et à introduire entre deux étapes de calcul une inactivité des processeurs de durée c pour permettre les communications.

5.5 Tableau récapitulatif

Le tableau ci-dessous résume les différents résultats pour l'ordonnancement statique lorsque les tâches sont de même durée ($p_i = 1$) et les communications de durée c .

	m entrée		m fixé	m suffisant	
	B_{inf}	B_{sup}		B_{inf}	B_{sup}
$c = 0$	$\frac{4}{3}$	2	m = 2 (Classe P) m = 3 (Pb Ouvert)	1	1
$c = 1$	$\frac{5}{4}$	$\frac{7}{3}$	m = 2 (Pb Ouvert)	$\frac{7}{6}$	$\frac{4}{3}$
c grand	$1 + \frac{1}{c+3}$	$c + 2$	m = 2 (NP-Complet)	1	$1 + c$

6 Systèmes hétérogènes

Le rôle croissant que joue les réseaux de stations de travail dans le domaine du calcul parallèle, nous oblige à intégrer ce paramètre dans nos modèles. La plupart de ces réseaux sont des réseaux constitués de stations qui n'ont pas tous la même puissance et il est évident que cette hétérogénéité pose de nouveaux problèmes sur le plan de la modélisation et de l'analyse des problèmes d'ordonnancement associés.

Dans le cas général, même si l'on ne considère pas les coûts de communication, le temps d'une tâche ne peut être connu qu'après sa terminaison et dépend non seulement du processeur auquel elle a été affectée mais encore de l'état de ce processeur durant son exécution. Le problème apparaît donc complexe à modéliser. Le but de cette section est d'étudier cette complexité même sous des hypothèses simplificatrices.

6.1 Difficulté d'un ordonnancement optimal

Dans cette section nous nous limitons au cas de systèmes comportant une machine puissante, que l'on va noter P_1 et un nombre illimité de machines de faible puissance, que nous appellerons de type P_2 . Ce problème est noté $P_1 \cup \bar{P}_2 | prec, p_i^1, p_i^2 | C_{max}$ où p_i^1 (resp. p_i^2) représente la durée d'exécution sur P_1 (resp. sur un processeur de type P_2).

Un ordonnancement consiste ici en deux types d'information :

- l'affectation des tâches au processeur (choix du type de processeur)
- l'ordre d'exécution sur chaque processeur (ici sur P_1)

En fait, il a été montré que même si on se limite à des graphes bipartis de profondeur 1, le problème est déjà NP-difficile. Même lorsque l'affectation des tâches au processeur est imposée le problème reste NP-complet dans des versions très simples (par exemple lorsque le graphe est une réunion de chemins de longueur 3). Pour plus de détails voir [2].

6.2 Machines uniformes et non-uniformes

Lorsque la durée d'une tâche est inconnue mais ne dépend que du processeur sur lequel elle a été affectée, deux modèles de machine sont distingués (m désigne le nombre de processeurs

de la machine) :

- machine *uniforme*: les processeurs sont de vitesses proportionnelles, le rapport des vitesses étant éventuellement inconnu).
- machine *non-uniforme*: le temps d'une tâche dépend du processeur sur lequel elle s'exécute.

Dans ces deux cas, $\Omega(\log m)$ est une borne inférieure pour le ratio de performance d'un algorithme d'ordonnancement dynamique avec ou sans préemption [33, 18].

Par ailleurs, des algorithmes polynomiaux ont été proposés pour ces deux modèles qui permettent d'atteindre un ratio $O(\log n)$, n étant le nombre de tâches [27]. En regroupant les tâches affectées à un même processeur, [33] donnent un algorithme d'ordonnancement dynamique de ratio $O(\log \min(m, R))$ où R le rapport de vitesse entre le processeur le plus rapide et le plus lent. Il est à noter que tous ces ratio peuvent être obtenus par des algorithmes d'ordonnancement non-préemptifs; la migration n'apporte un gain que d'un facteur constant dans le ratio.

6.3 Simulation d'une machine à processeurs identiques

Une approche pragmatique pour construire des algorithmes ayant une justification théorique (par exemple un algorithme de liste) est de simuler une machine constituée de processeurs identiques sur une architecture hétérogène. Une telle simulation est facilitée par la manipulation de processus légers dont le nombre sur chaque processeur peut varier en fonction de l'état du processeur.

7 Conclusion

Les différents résultats énoncés dans ce chapitre montrent la difficulté d'obtenir des résultats analytiques non triviaux lorsqu'on tente d'intégrer dans les modèles d'ordonnancement certains paramètres primordiaux d'un point de vue pratique comme les grands délais de communication. Dans le cas général, il n'existe pas d'algorithmes d'ordonnancement ayant un ratio de performance constant (indépendant de l'architecture) sur une architecture hétérogène. Il est donc important de considérer des classes spécifiques de problèmes à ordonnancer.

Une approche harmonieuse semble être de se restreindre à des graphes de dépendance (éventuellement construits dynamiquement en fonction des entrées) dont le volume maximal de communication est négligeable devant le volume de calcul. Aussi bien dans le cadre dynamique (algorithmes de liste) que statique (granularité grossière, $p_i = c = 1$), des ordonnancements efficaces sont alors connus. Cependant, l'obtention d'un tel graphe par compilation à partir d'un programme de grain fin nécessite des algorithmes de partition efficaces qui risquent de réduire fortement le degré de parallélisme sur des problèmes ayant un comportement irréguliers.

Références

- [1] R.J. Anderson and G.L. Miller. Deterministic parallel list ranking. In *AWOC88*, number 319 in LNCS, pages 81–90. Springer-Verlag, 1988.
- [2] E. Bampis, A. Giannakos, and J-C. Konig. From simple-machine problems to heterogeneous scheduling. Rapport technique 25, LaMI, Université d'Evry, France, 1997.
- [3] J. Blazewicz, K. Exker, G. Schmidt, and J. Węglarz. *Scheduling in Computer and Manufacturing Systems*. Springer-Verlag, Germany, 1993.
- [4] Guy E. Blelloch, Phillip B. Gibbons, and Yossi Matias. Provably efficient scheduling for languages with fine-grained parallelism. In *Proceedings of the 7th Symposium on Parallel Algorithms and Architectures*, pages 1–12, Santa-Barbara, California, 1995. ACM Press.
- [5] Guy E. Blelloch, Phillip B. Gibbons, Yossi Matias, and Girija J. Narlikar. Space efficient scheduling of Parallelism with Synchronization Variables. In *Proceedings of the 9th Symposium on Parallel Algorithms and Architectures*. ACM Press, June 1997.
- [6] Robert D. Blumofe. *Executing Multithreaded Programs Efficiently*. PhD thesis, Massachusetts Institute of Technology, Boston, 1995.
- [7] Robert D. Blumofe and Charles E. Leiserson. Space-efficient scheduling of multithreaded computations. *SIAM Journal on Computing*, 27(1):202–229, 1998.
- [8] J. Briat, T. Gautier, and J.-L. Roch. Application irrégulière et ordonnancement en ligne. In G. Bernard, J. Chassin de Kergommeaux, Folliot B., and C. Roucairol, editors, *Placement dynamique et répartition de charge: application aux systèmes répartis et parallèles*, pages 81–106. Collection didactique INRIA, 1996.
- [9] R. Cole and U. Vishkin. Approximate Parallel Scheduling. Part I: The Basic Technique with Applications to Optimal Parallel List Ranking in Logarithmic Time. *SIAM Journal on Computing*, 17(1), 1988.
- [10] J.Y. Colin and P. Chretienne. C.P.M. Scheduling with small communication delays and task duplication. *Operations Research*, 39(3):680–684, 1991.
- [11] David E. Culler, Richard M. Karp, David Patterson, Abhijit Sahay, Eunice E. Santos, Klaus Erik Schausser, and Thorsten von Eicken Ramesh Subramonian. LogP: A Practical Model of Parallel Computation. *Communications ACM*, 39(11):78–85, 1996.
- [12] François Galilée, Jean-Louis Roch, Gerson Cavalheiro, and Matthias Doreille. Athapascan-1: On-line Building Data Flow Graph in a Parallel Language. In IEEE, editor, *International Conference on Parallel Architectures and Compilation Techniques, PACT'98*, pages 88–95, Paris, France, October 1998.
- [13] T. Gautier, J.L. Roch, and G. Villard. Regular versus irregular problems and algorithms. In *Proc. of IRREGULAR'95, Lyon, France*. Springer-Verlag, Sep. 1995.
- [14] P.B. Gibbons. A more practical PRAM model. In *Proceedings of the 1989 ACM Symposium on Parallel Algorithms and Architectures*, 1989.
- [15] R.L. Graham. Bounds for Certain Multiprocessor Anomalies. *Bell System Tech J.*, 45:1563–1581, 1966.
- [16] R.L. Graham. Bounds on multiprocessing timing anomalies. *SIAM J. Appl. Math.*, 17(2):416–426, 1969.

- [17] R.L. Graham, E.L. Lawler, J.K. Lenstra, and A.H.G. Rinnooy Kan. Optimization and approximation in deterministic scheduling: A survey. *Ann. Disc. Math.*, 5:287–326, 1979.
- [18] L. A. Hall, A. S. Schulz, D. B. Shmoys, and J. Wei n. Scheduling to Minimize Average Completion Time: Off-line and On-line Approximation Algorithms. Technical Report 516/1996, Technische Universität, Berlin, 1996.
- [19] J. A. Hoogeveen, J.K. Lenstra, and B. Veltman. Three, four, five, six or the complexity of scheduling with communication delays. *Operations Research Letters*, 16:129–137, 1994.
- [20] C.F. Joerg. *The Cilk system for parallel multithreaded computing*. PhD thesis, Massachusetts Institute of Technology, january 1996.
- [21] Coffman E.G. Jr. and R.L. Graham. Optimal scheduling for two-processor systems. *Acta Informatica*, (1):200–213, 1972.
- [22] R. M. Karp, M. Luby, and F. Meyer auf der Heide. Efficient PRAM Simulation on a Distributed Memory Machine. *Algorithmica*, 16:517–542, 1996.
- [23] Clyde P. Kruskal, Larry Rudolph, and Marc Snir. A complexity theory of efficient parallel algorithms. *Theoretical Computer Science*, 71:95–132, 1990.
- [24] S. Lam and R. Sehti. Worst case analysis of two scheduling algorithms. *SIAM J. Computing*, 6(3):518–538, 1977.
- [25] J. K. Lenstra and D. B. Shmoys. Computing near-optimal schedules. In P. Chrétienne and al., editors, *Scheduling, Theory and its Applications*, pages 1–14. J. Wiley, 1995.
- [26] J.K. Lenstra and A.H.G. Rinnooy Kan. Complexity of scheduling under precedence constraints. *Operations Research*, 26:22–35, 1978.
- [27] J.K. Lenstra, D.B. Shmoys, and É. Tardos. Approximation algorithms for scheduling unrelated parallel machines. *Math. Programming*, 46:259–271, 1990.
- [28] J.K. Lenstra, M. Veldhorst, and B. Veltman. The complexity of scheduling trees with communication delays. *Journal of Algorithms*, 20(1):157–173, 1996.
- [29] A. Munier and J.-C. König. A heuristic for a scheduling problem with communication delays. A paraître dans *Operations Research*.
- [30] C. Picouleau. *Etude des Problèmes d’Optimisation dans les Systèmes Distribués*. PhD thesis, Université Paris VI, France, 1992.
- [31] Abhiram G. Ranade. How to emulate shared memory. In *Proceedings 28th Annual Symposium on Foundations of Computer Science*, pages 185–192. IEEE, 1987.
- [32] M.C. Rinard. The design, implementation and evaluation of Jade : a portable, implicitly parallel programming language. *ACM Transactions on Programming Languages and Systems*, 20(3):483–545, May 1998.
- [33] D. B. Shmoys, J. Wein, and P. Williamson. Scheduling parallel machines on-line. *SIAM Journal on Computing*, 24(6):1313–1331, 1995.
- [34] Leslie G. Valiant. A Bridging Model For Parallel Computation. *Communications of the ACM*, 33(8):103–111, 1990.
- [35] Leslie G. Valiant. A Bridging Model For Parallel Computation. *Communications of the ACM*, 33(8):103–111, 1990.

- [36] Leslie G. Valiant. General purpose parallel architectures. In J. van Leuwen, editor, *Algorithms and Complexity*, pages 944–971. Elsevier, 1990.
- [37] T. Varvarigou, V.P. Roychowdhury, and T. Kailath. Scheduling in and out forests in the presence of communication delays. In *Proceedings International Parallel Processing Symposium*, pages 222–229, Newport Beach, CA, 1993.
- [38] B. Veltman. *Multiprocessor Scheduling with Communication Delays*. PhD thesis, CWI-Amsterdam, Holland, 1993.
- [39] B. Veltman, B.J. Lageweg, and J.K. Lenstra. Multiprocessor scheduling with communication delays. *Parallel Computing*, 16:173–182, 1990.