



**HAL**  
open science

# Algorithmes d'ordonnement en deux étapes de graphes de tâches parallèles

Tchimou N'Takpé

► **To cite this version:**

Tchimou N'Takpé. Algorithmes d'ordonnement en deux étapes de graphes de tâches parallèles. 2006. inria-00125269v1

**HAL Id: inria-00125269**

**<https://inria.hal.science/inria-00125269v1>**

Submitted on 18 Jan 2007 (v1), last revised 7 Nov 2008 (v6)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# Algorithmes d'ordonnancement en deux étapes de graphes de tâches parallèles

**Tchimou N'takpé**<sup>1</sup>

Nancy Université / LORIA<sup>2</sup>  
Campus Scientifique - BP 239  
F-54506 Vandoeuvre-lès-Nancy Cedex  
Tchimou.Ntakpe@loria.fr

---

*RÉSUMÉ.* L'ordonnancement d'applications parallèles représentées par des graphes de tâches consiste à trouver l'ensemble de processeurs sur lesquels chaque tâche doit être exécutée afin de minimiser le temps d'exécution de ces applications tout en exploitant rationnellement les ressources. Alors que la plupart des algorithmes d'ordonnancement de graphes de tâches parallèles visent des grappes homogènes, cet article montre la nécessité d'avoir de tels algorithmes pour des agrégations de grappes de calcul qui sont de plus en plus répandues et qui peuvent permettre de déployer des applications parallèles à échelles sans précédents. Nous proposons des améliorations d'une heuristique d'ordonnancement de tâches parallèles en milieu homogène. Ensuite, nous l'adaptions au cas des plates-formes hétérogènes de type grappe hétérogène de grappes homogènes.

*ABSTRACT.* While most parallel task graphs scheduling research has been done in the context of single homogeneous clusters, heterogeneous platforms have become prevalent and are extremely attractive for deploying applications at unprecedented scales. In this paper we address the need for scheduling techniques for parallel task applications for heterogeneous clusters of clusters by proposing a method to adapt existing parallel task graphs scheduling heuristics that have proved to be efficient on homogeneous environments. Before adapting that heuristic to heterogeneous platforms, we propose some improvements for homogeneous platforms

*MOTS-CLÉS :* Heuristiques d'ordonnancement, tâches parallèles, DAGs, grappe de grappes

*KEYWORDS:* Heterogeneous scheduling, parallel tasks, DAGs, cluster of clusters

---

---

1. Cette étude a été en partie soutenue par l'ARC INRIA OTaPHe, la Région Lorraine et le Gouvernement Ivoirien  
2. UMR 7503 CNRS - INPL - INRIA - Nancy 2 - UHP, Nancy 1

## 1. Introduction

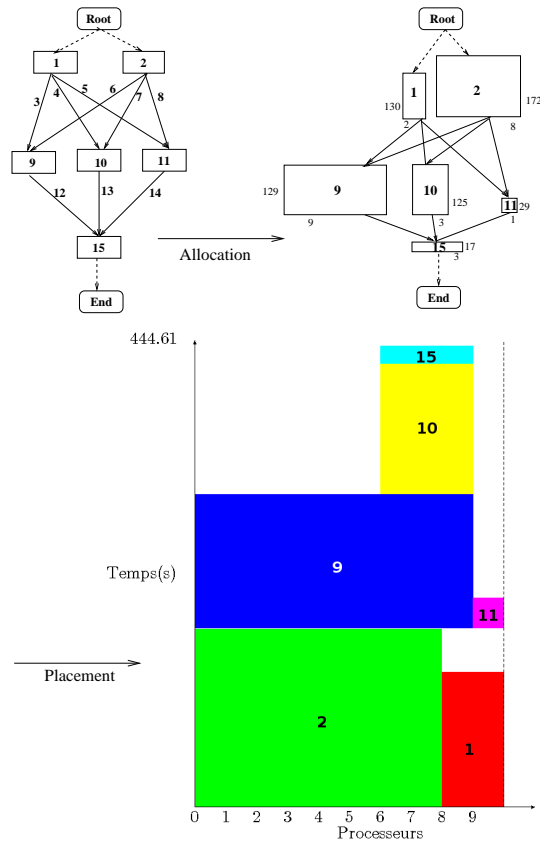
Une approche récente permettant de pallier la demande croissante en mémoire et en ressources de calcul des applications parallèles consiste à agréger des grappes de calcul existantes soit au sein d'une seule institution, soit réparties entre plusieurs institutions. Il s'agit souvent de grappes de tailles variables dont les capacités peuvent être différentes en fonction des technologies présentes au moment de leur installation. Ces plates-formes sont à la fois attractives parce qu'elles offrent une importante puissance de calcul et un déficit pour les chercheurs du fait de leur hétérogénéité.

Une des méthodes qui permettent d'exploiter la puissance de calcul ainsi disponible est de combiner les parallélismes de tâches et de données présents dans les applications scientifiques. Ces applications peuvent alors être modélisées par des graphes de tâches parallèles. De manière informelle, une tâche parallèle est une tâche qui contient des opérations élémentaires, typiquement une routine numérique ou des boucles imbriquées, qui contiennent suffisamment de parallélisme pour être exécutées par plus d'un processeur. Dans cet article, nous considérons un certain type de tâches parallèles : les *tâches modelables* (Turek *et al.*, 1992). Ce sont des tâches parallèles pouvant s'exécuter sur un nombre quelconque de processeurs. Ce nombre n'est pas fixé a priori mais est déterminé avant l'exécution et ne varie pas ultérieurement.

La figure 1 illustre le problème de l'ordonnancement de tâches parallèles qui consiste dans un premier temps à trouver le bon nombre de processeurs à allouer à chaque tâche : c'est la *phase d'allocation*. Ensuite il faut déterminer les processeurs sur lesquels placer les différentes tâches : c'est la *phase de placement*. Les algorithmes d'ordonnancement de tâches parallèles peuvent être repartis en deux catégories. La première catégorie est constituée d'algorithmes où la phase d'allocation et la phase de placement sont indissociables : ce sont les algorithmes d'*ordonnancement en une étape* (Boudet *et al.*, 2003). La seconde catégorie regroupe les algorithmes d'*ordonnancement en deux étapes* (Radulescu *et al.*, 2001a; Radulescu *et al.*, 2001b; Ramaswamy *et al.*, 1997; Rauber *et al.*, 1998) où l'allocation et le placement sont séparés en deux procédures distinctes. Notons dans cet exemple qu'à l'issue de la phase d'allocation, l'on obtient une estimation des temps d'exécution de chaque tâche (nombre sur le côté vertical des boîtes représentant les tâches) à partir son allocation (nombre sur le côté horizontal).

De nombreuses études ont été réalisées pour l'ordonnancement de graphes de tâches parallèles dans le cas de plates-formes homogènes (Lepère *et al.*, 2002; Radulescu *et al.*, 2001a; Radulescu *et al.*, 2001b; Ramaswamy *et al.*, 1997; Rauber *et al.*, 1998). Or les plates-formes hétérogènes sont de plus en plus répandues et très attractives car elles peuvent permettre de déployer des applications parallèles à échelles sans précédents. Il est donc nécessaire de développer des heuristiques d'ordonnancement pour ces systèmes hétérogènes. Une première approche a consisté à adapter une heuristique d'ordonnancement de tâches séquentielles sur plates-formes hétérogènes au cas des tâches parallèles (Casanova *et al.*, 2004). Ici, nous utilisons une approche complémentaire qui consiste à modifier un algorithme d'ordonnancement de tâches

2



**Figure 1.** Exemple d'ordonnancement en deux étapes d'un graphe de tâche synthétique sur une grappe homogène de 10 processeurs.

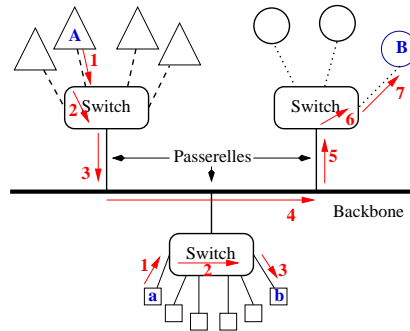
parallèles en milieu homogène (Radulescu *et al.*, 2001b) et à prendre en compte l'hétérogénéité des ressources.

Les contributions de cet article sont : (i) apporter des améliorations à l'heuristique choisie en milieu homogène, aussi bien dans la procédure d'allocation que dans la procédure de placement ; (ii) utiliser un nouveau concept de virtualisation des plateformes qui permet de gérer plus facilement les allocations de ressources en milieu hétérogène ; (iii) introduire une nouvelle méthode de placement fondée sur l'idée de l'heuristique *sufferage* (Casanova *et al.*, 2000) ; (iv) et, définir une méthode d'évaluation par simulation de nombreux scénarios.

Dans la section suivante, nous présentons les modèles de plates-formes et d'applications utilisés. La section 3 décrit notre méthodologie d'évaluation des algorithmes conçus. La section 4 traite des travaux précédents et formalise le problème. Après avoir apporté quelques améliorations à CPA, un algorithme d'ordonnancement (en milieu homogène) efficace issu de la littérature dans la section 5, nous adapterons l'algorithme obtenu aux plates-formes hétérogènes dans la section 6. Enfin, la section 7 nous permettra de conclure.

## 2. Modèles de plates-formes et d'applications

Dans cet article, nous considérons des agrégations hétérogènes de grappes homogènes. Ces plates-formes sont représentatives de certaines infrastructures de grilles de calcul réparties entre des institutions, qui disposent généralement de grappes homogènes. Nous avons donc  $C$  grappes et chaque grappe  $C_k$  contient  $P_k$  processeurs identiques pour un total de  $P$  processeurs sur la plate-forme. Les vitesses des processeurs et les caractéristiques des réseaux locaux ne sont pas nécessairement les mêmes entre les différentes grappes. La figure 2 montre la structure de nos plates-formes. Les processeurs d'une même grappe sont reliés entre eux à travers un commutateur (Switch). Les grappes sont reliées par le biais d'une passerelle à un lien réseau très haut débit (Backbone) commun.



**Figure 2.** *Modèle de plate-forme avec les routes entre deux processeurs A et B localisés sur deux grappes différentes et entre deux processeurs a et b de la même grappe.*

Une application parallèle peut être modélisée par un graphe acyclique orienté (DAG)  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$  où  $\mathcal{N} = \{t_i / i = 1, \dots, N\}$  est un ensemble de  $N$  nœuds (ou tâches) et  $\mathcal{E} = \{e_{i,j} / (i, j) \subset \{1, \dots, N\} \times \{1, \dots, N\}\}$  est un ensemble de  $E$  arêtes. Les nœuds représentent les tâches parallèles et les arêtes définissent les relations de précedence (dépendances de flux ou de données) entre les tâches. On associe à chaque arête  $e_{i,j}$ , la quantité de données que la tâche  $t_i$  doit transférer à la tâche  $t_j$ . Par définition, une tâche d'entrée du graphe n'a aucun prédécesseur et une tâche de sortie est sans successeur. Dans cette étude nous utilisons les nœuds **Root** et **End** pour représenter respectivement les tâches d'entrée et de sortie des DAGs. Ce sont

deux nœuds sans coût (de calcul ou de communication) qui facilitent la manipulation des graphes de tâches. Une tâche est *prête* lorsque tous ses prédécesseurs ont terminé leur exécution.

Dans une grappe homogène, le temps d'exécution d'une tâche parallèle  $t \in \mathcal{N}$  peut être modélisé par un modèle d'accélération classique. Dans cet article, nous utilisons la loi d'Amdahl (Amdahl, 1967) où le temps d'exécution parallèle est donné par :

$$T(t, p(t)) = \left( \alpha + \frac{1 - \alpha}{p(t)} \right) \cdot T(t, 1), \quad [1]$$

$p(t)$  étant le nombre de processeurs alloués à  $t$ ,  $T(t, 1)$  son temps d'exécution sur un seul processeur et  $\alpha$  sa portion non parallélisable. Dans cette étude, nous restreignons l'exécution d'une tâche parallèle à l'intérieur d'une grappe, afin de garantir la validité de ce modèle. Ce choix est également motivé par le fait qu'en pratique les communications inter-grappes peuvent être très coûteuses. On définit enfin  $T_b(t)$ , le *bottom level*, comme étant le chemin (somme des temps d'exécution et des temps de transfert de données) le plus long depuis la tâche  $t$ , incluant son propre temps d'exécution, jusqu'à une *tâche de sortie* quelconque.

### 3. Méthodologie d'évaluation

Pour l'évaluation de nos algorithmes, nous aurons recours à des simulations afin d'explorer une large variété de plates-formes et d'applications. Pour cela nous utilisons SIMGRID<sup>1</sup> (Legrand *et al.*, 2003), une boîte à outils conçue pour la simulation de grilles de calcul et/ou la mise en œuvre d'applications distribuées. Les simulations seront effectuées sur des plates-formes et des DAGs générées aléatoirement en faisant varier quelques paramètres qui permettent de fixer leurs propriétés.

Les plates-formes générées sont constituées de 1, 2, 4 ou 8 grappes. Le nombre de processeurs de chaque grappe est tiré aléatoirement entre 16 et 128 pour les plates-formes dont le nombre de grappes est supérieur ou égal à 2. Pour les plates-formes à une seule grappe, le nombre de processeurs est tiré entre 16 et 512. Les liens intra-grappes sont du Fast Ethernet (débit =  $100Mb/s$  et latence =  $100\mu s$ ) ou du Giga Ethernet (débit =  $1Gb/s$  et latence =  $100\mu s$ ). Ces liens peuvent subir des contentions. Le backbone a un débit de  $2,5Gb/s$  et une latence de  $50ms$ . Chaque passerelle a une capacité de  $1Gb/s$  avec une latence de  $100\mu s$ . La borne inférieure des vitesses des processeurs (en *GFlops*) est de 0,25, 0,5, 0,75 ou 1. En ce qui concerne les plates-formes ayant plusieurs grappes, le rapport entre cette borne inférieure et la borne supérieure des vitesses des processeurs de la plate-forme, qui constitue le degré d'hétérogénéité, est pris parmi les valeurs 1, 2 ou 5. Nous choisissons ainsi la vitesse des processeurs des différentes grappes de la plate-forme en tirant aléatoirement des vitesses entre les bornes inférieure et supérieure. Par exemple si la borne inférieure des vitesses de

1. <http://simgrid.gforge.inria.fr>

processeurs est fixée à  $0,5GFlops$  et le degré d'hétérogénéité à 5, alors les vitesses des processeurs sont tirées entre  $0,5GFlops$  et  $2,5GFlops$ . Pour chaque combinaison de ces paramètres nous générons 10 échantillons dans le cas des plates-formes à une grappe et 5 échantillons pour les autres. Au total, nous avons généré 220 plates-formes dont 40 à une grappe et 180 ( $60 \times 3$ ) pour les plates-formes ayant plusieurs grappes.

Les graphes de tâches sont générés en tirant la taille des données  $m$ , un multiple de 1024 (1 ko), entre les valeurs limites 2048 et 11268, ce qui correspond au plus à 1 Go d'occupation mémoire par tâche. Ensuite, la complexité de toutes les tâches d'un même DAG est fixée à  $a \cdot M$ ,  $a \cdot M \log M$ , ou  $a \cdot M^{3/2}$ , ou tirée au sort parmi les trois valeurs précédentes.  $M = m^2$  et  $a$  est un nombre tiré aléatoirement entre  $2^6$  et  $2^9$ . La portion non parallélisable de chaque tâche (le  $\alpha$  de la loi d'Amdahl) est un nombre aléatoire pris entre 0 et 0,2. Le coût des transferts est égal à  $M$ , où  $M$  est relatif à la tâche qui génère le transfert.

Quatre paramètres permettent de régler la forme des graphes : (i) la largeur (0,1, 0,2 ou 0,8). Une largeur de 0,8 correspond à un graphe compact avec beaucoup de parallélisme de tâches ; (ii) la régularité entre niveaux (0,2 ou 0,8). Dans un graphe peu régulier, la différence entre les nombres de tâches sur deux niveaux consécutifs peut être très importante ; (iii) la densité qui caractérise le fait que l'on a plus ou moins de relations de précédence entre les tâches de l'application (0,2 ou 0,8) ; et (iv) la longueur maximale des sauts entre niveaux (1, 2 ou 4). Ce dernier paramètre sert à générer des graphes contenant des chemins de longueur différentes (en nombre de nœuds) entre les tâches d'entrée et de sortie. Pour chaque combinaison de ces paramètres, nous générons 3 échantillons différents, soit 1296 DAGs.

Afin de comparer les algorithmes mis en jeu, nous utilisons les trois métriques suivantes :

**Le temps de complétion** des applications (*makespan*), *i.e.*, la différence entre leur date de début et leur date de fin d'exécution. Autrement dit, c'est le temps d'exécution total de l'application.

**L'énergie** consommée par chaque application (DAG) durant son exécution (en  $GFlop \cdot s$ ). Il s'agit d'une nouvelle métrique que nous introduisons afin de quantifier la consommation en ressources de calcul par l'application. Nous mesurons cette énergie totale en faisant la somme des énergies consommées par chaque tâche. L'énergie consommée par une tâche est le produit de son temps d'exécution (en  $s$ ) par la puissance de calcul utilisée (en  $GFlops$ ).

**L'efficacité** des algorithmes par rapport à un algorithme d'ordonnancement séquentiel (SEQ) qui ordonnance séquentiellement les tâches sur le processeur le plus rapide de la plate-forme considérée. Sur une grappe homogène, l'efficacité se définit comme étant le rapport de l'accélération par rapport à SEQ sur le nombre moyen de processeurs utilisés par l'application. Nous généralisons cette définition dans l'équation 2 : pour un DAG et une plate-forme donnés, l'efficacité peut se calculer en faisant le rapport de l'énergie consommée en utilisant SEQ par l'énergie consommée en utilisant l'algorithme à évaluer. Une efficacité élevée

traduit le fait que les ressources sont rationnellement utilisées dans l'exécution parallèle de l'application.

$$Eff_{ALGO} = \frac{Energie_{SEQ}}{Energie_{ALGO}} \quad [2]$$

#### 4. Travaux précédents

Le problème de l'ordonnancement de tâches en prenant en compte le coût des communications, est NP-complet au sens fort, même dans le cas où il existe un nombre infini de processeurs (Chretienne, 1988). De nombreuses heuristiques ont donc été conçues pour ordonnancer des tâches parallèles. Dans (Dutot, 2005), un algorithme d'ordonnancement de tâches modelables interdépendantes sur une hiérarchie de machines multiprocesseurs est présenté. Dans (Casanova *et al.*, 2004), les auteurs proposent l'une des rares heuristiques d'ordonnancement de tâches parallèles destinées aux plates-formes hétérogènes. Ils adaptent un algorithme d'ordonnancement de tâches séquentielles sur plates-formes hétérogènes au cas de DAGs de tâches parallèles. Notre approche complète cette dernière puisque nous partons d'un algorithme d'ordonnancement de graphes de tâches parallèles sur plates-formes homogènes pour l'adapter au cas de plates-formes hétérogènes.

Si certaines études théoriques existent (Lepère *et al.*, 2002), la plupart des algorithmes d'ordonnancement de tâches parallèles (Radulescu *et al.*, 2001a; Radulescu *et al.*, 2001b; Ramaswamy *et al.*, 1997; Rauber *et al.*, 1998) sont des algorithmes en deux étapes et ont été conçus pour des milieux homogènes. La première étape vise à déterminer un nombre de processeurs optimal pour chaque tâche. Dans la seconde étape, les auteurs utilisent des heuristiques de liste pour ordonnancer les tâches.

Dans (Ramaswamy *et al.*, 1997), les auteurs extraient des DAGs à partir de codes séquentiels puis appliquent leur algorithme TSAS (*Two Step Allocation and Scheduling*). TSAS utilise la programmation convexe, rendue possible grâce à la propriété de *posynomialité* des modèles de coût choisis, ainsi que certaines propriétés de leur structure de DAGs. Les auteurs de (Rauber *et al.*, 1998) limitent quant à eux leur étude à des graphes construits par compositions séries et/ou parallèles. Dans le premier cas, une séquence d'opérations présentant des dépendances de données est placée sur l'ensemble des processeurs. Les tâches de cette séquence sont alors exécutées séquentiellement. Dans le second cas, l'ensemble des processeurs est divisé en un nombre optimal de sous-ensembles, déterminé par un algorithme glouton. Le critère d'optimisation de cet algorithme est la minimisation du temps de complétion de l'ensemble des tâches.

Dans cet article, nous nous intéressons à l'algorithme CPA (Radulescu *et al.*, 2001b) qui est le plus performant des algorithmes en deux étapes que nous venons de citer. CPA (*Critical Path and Area-based Scheduling*) vise à obtenir le meilleur compromis entre la longueur du chemin critique, *i.e.*, le plus long chemin du graphe

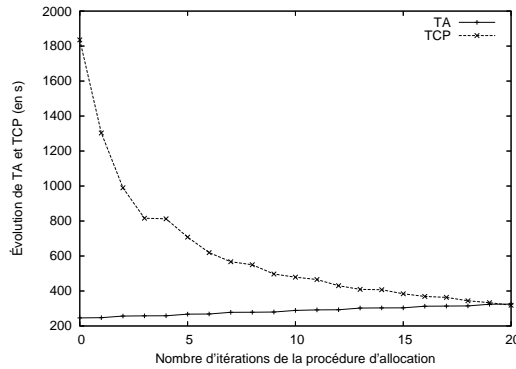


en tenant uniquement compte des temps d'exécution des tâches, et l'aire moyenne du diagramme de Gantt qui représente le temps moyen d'occupation des processeurs. Rădulescu *et al.* remarquent que le temps d'exécution d'une application parallèle peut être approché par sa borne inférieure  $T_p^e = \max\{T_{CP}, T_A\}$ , où  $T_{CP}$  est la longueur du chemin critique et  $T_A$ , l'aire moyenne :

$$T_{CP} = \max_{t \in \mathcal{N}} T_b(t), \text{ et} \quad [3]$$

$$T_A = \frac{1}{P} \sum_{i=0}^N (T(t_i, p(t_i)) \times p(t_i)). \quad [4]$$

Le but de CPA est de minimiser  $T_p^e$  au terme de la phase d'allocation. Sachant que  $T_{CP}$  diminue tandis que  $T_A$  croît lorsque le nombre de processeurs alloués aux tâches augmente, on initialise les allocations en partant du cas où  $T_{CP}$  est maximal. On alloue donc un processeur à chaque tâche à l'instant initial. Ensuite à chaque itération, on alloue un processeur de plus à la tâche la plus prioritaire jusqu'à ce que l'on obtienne  $T_{CP} \leq T_A$ . Cette tâche la plus prioritaire est celle appartenant au chemin critique et dont le rapport  $T(t, p(t))/p(t)$  diminue le plus significativement si un processeur supplémentaire lui est attribué. Dès qu'elle est vérifiée, la condition d'arrêt de cette procédure d'allocation ( $T_{CP} \leq T_A$ ) traduit le fait que  $T_p^e$  est proche de sa valeur minimale ( $T_p^e \approx T_{CP} \approx T_A$ ). La figure 3 présente l'évolution de  $T_{CP}$  et de  $T_A$  en fonction du nombre d'itérations de la procédure d'allocation de CPA appliquée à l'exemple de la figure 1.



**Figure 3.** Exemple d'évolution de  $T_{CP}$  et  $T_A$  lors de la procédure d'allocation.

À chaque itération de la procédure de placement, la tâche prête ayant le plus grand *bottom level* est choisie pour être placée. Cette phase tient compte du coût des redistributions des données pour déterminer la date de début d'exécution,  $T_s(t)$ , et la date de fin d'exécution,  $T_f(t)$ , de chaque tâche  $t$ .

Dans (Radulescu *et al.*, 2001b) il a été montré que la complexité de CPA est de l'ordre de

$$O(N(N + E)P). \quad [5]$$

## 5. Améliorations de CPA en milieu homogène

Dans cette section, nous nous plaçons dans le cas où les plates-formes sont constituées d'une seule grappe homogène et nous proposons deux améliorations de l'algorithme original : la première porte sur la phase d'allocation et la seconde sur la phase de placement.

### 5.1. Nouveau critère d'arrêt dans la procédure d'allocation

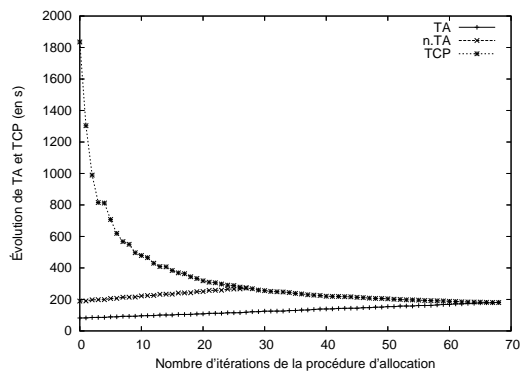
Par expérience, nous avons constaté que le calcul de l'aire moyenne de CPA n'était plus pertinent lorsque le nombre de processeurs ( $P$ ) de la plate-forme est beaucoup plus grand que le nombre de tâches ( $N$ ).  $T_A$  converge très lentement vers  $T_{CP}$  lorsque le nombre de processeurs de la plate-forme est très grand. Ce qui conduit à des allocations sur un très grand nombre de processeurs. Or, plus l'on alloue de processeurs aux tâches, plus le risque de ne plus pouvoir exécuter en parallèle certaines tâches concurrentes augmente. Il peut donc s'avérer préférable d'arrêter les allocations plus tôt afin de profiter au mieux du parallélisme de tâches. Nous proposons donc le compromis suivant qui permet d'arrêter plus vite la procédure d'allocation dans le cas où le nombre de ressources est très élevé en prenant  $\min(P, \sqrt{P \times N})$  au lieu de  $P$ . Cela nous conduit à redéfinir la notion d'aire moyenne de la façon suivante :

$$T_A = \frac{1}{\min(P, \sqrt{P \times N})} \sum_{i=0}^N (T(t_i, p(t_i)) \times p(t_i)). \quad [6]$$

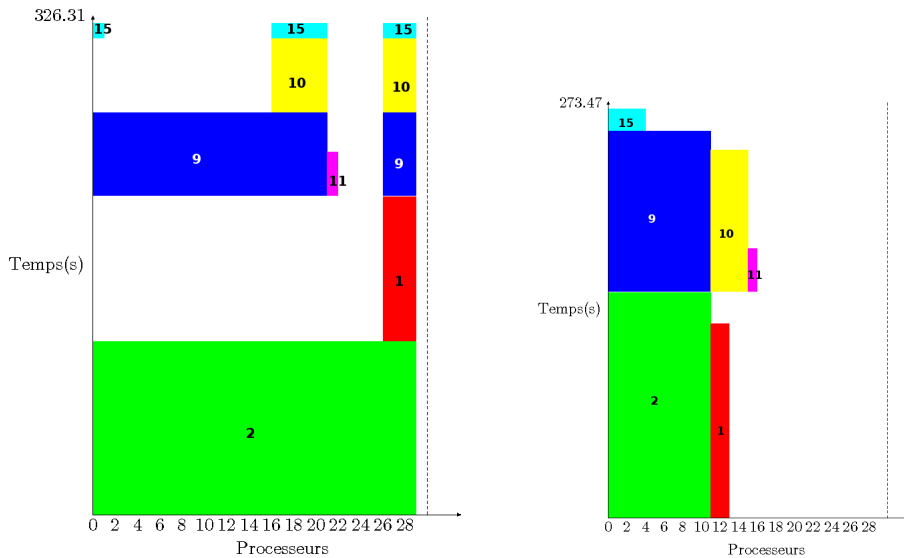
Pour  $P > N$ , cette nouvelle définition augmente la pente de croissance de  $T_A$ . La relation  $T_p^e \approx T_{CP}$  reste toujours valable à la fin de la procédure d'allocation.

La figure 4 montre l'évolution de  $T_{CP}$  et  $T_A$  en utilisant CPA (TCP - TA) ou la nouvelle procédure d'allocation (TCP - n.TA) sur le DAG de la figure 1. Dans cet exemple, le nombre de processeurs ( $P = 30$ ) est supérieur au nombre de tâches ( $N = 6$ ). D'où une convergence plus rapide de  $T_{CP}$  et  $T_A$  dans le cas de la nouvelle définition de  $T_A$  (n.TA). Dans la figure 5 qui présente le résultat de l'ordonnancement, nous pouvons observer que l'exécution en parallèle des tâches 1 et 2, puis des tâches 9, 10 et 11 permet d'obtenir un meilleur temps de complétion par rapport à l'ordonnancement de CPA en plus du fait que l'application consomme moins de ressources. L'on peut entrevoir qu'il est possible de continuer les allocations et de réduire davantage le temps de complétion. Nous sommes conscients que le compromis que nous venons de

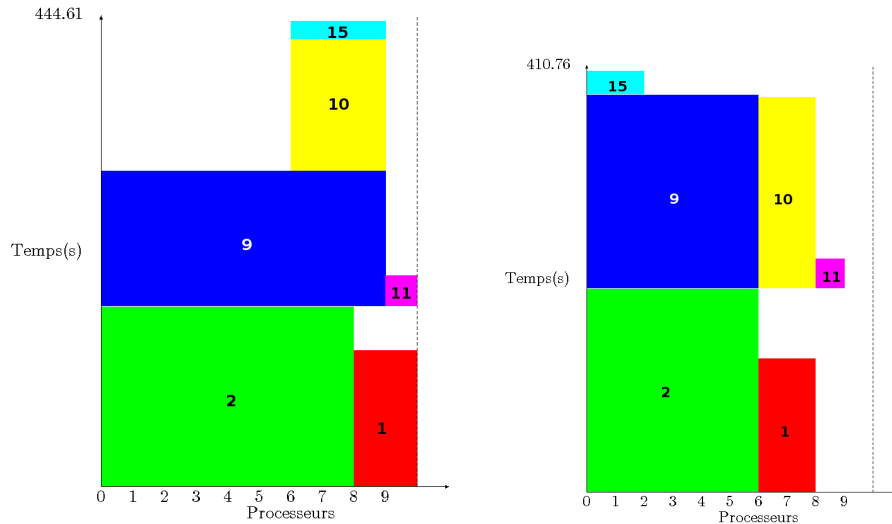
définir  $n$  n'est pas toujours optimale au niveau du temps de complétion des applications mais il permet surtout de mieux gérer l'utilisation des ressources. En effet, hormis les tâches « entièrement parallélisables », plus l'on alloue de processeurs aux tâches, plus l'efficacité diminue. Nous pouvons également remarquer dans la figure 6 qu'appliquer cette modification à l'exemple de la figure 1 (où  $N = 6$  et  $P = 10$ ) permet également de réduire le temps de complétion par rapport à CPA.



**Figure 4.** Exemple d'évolution de  $T_{CP}$  et de  $T_A$  dans la procédure d'allocation de CPA (TA) et la nouvelle procédure d'allocation (n.TA).



**Figure 5.** Ordonnancement du DAG de la figure 1 avec CPA (à gauche) et en utilisant la nouvelle allocation (à droite) sur une grappe homogène de 30 processeurs.



**Figure 6.** Ordonnancement du DAG de la figure 1 avec CPA (à gauche) et en utilisant la nouvelle allocation (à droite) sur une grappe homogène de 10 processeurs.

## 5.2. Tassage lors du placement

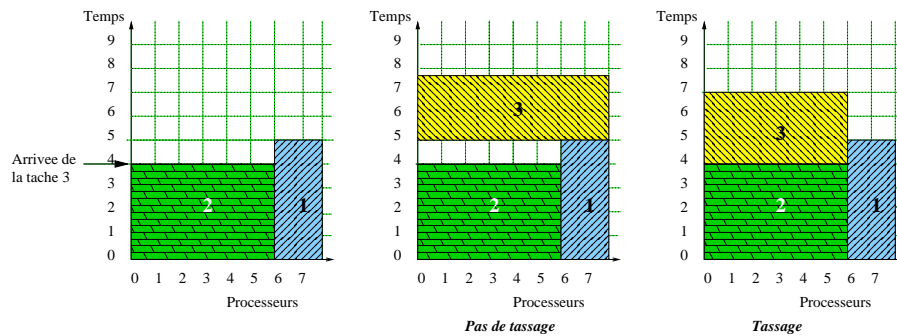
La phase d'allocation et la phase de placement des tâches étant découplées, il peut arriver qu'une tâche prête se mette à attendre qu'une partie des processeurs qui lui sont alloués soient disponibles alors que la majorité des processeurs dont elle aurait besoin le sont déjà. Cette tâche pourrait donc avoir une meilleure date de fin d'exécution si l'on réduisait son allocation de sorte qu'elle puisse démarrer à la date où elle est prête.

Le *tassage* permet d'éviter ce genre de situation. Il permet dans certains cas d'améliorer la date de fin d'exécution d'une tâche prête tout en réduisant le nombre de processeurs qui lui sont alloués. La procédure de placement avec tassage se comporte de la manière suivante :

1) À l'instant où la tâche prête la plus prioritaire est choisie, l'on regarde si celle-ci peut débuter son exécution immédiatement avec son allocation initiale, *i.e.*, le nombre de processeurs disponibles est supérieur ou égal à cette allocation. Dans ce cas, la tâche est placée sur les processeurs disponibles.

2) Si le nombre de processeurs disponibles est inférieur à l'allocation, il faut vérifier si la tâche terminera son exécution plus tôt en utilisant seulement les processeurs déjà disponibles plutôt que si elle attendait que tous les processeurs qui lui sont alloués soient libres. Si c'est le cas, alors une nouvelle allocation lui est attribuée et la tâche est placée sur les processeurs disponibles. Sinon elle attend qu'il y ait suffisamment de processeurs libres en vue d'être placée selon son allocation initiale.

La figure 7 illustre un exemple de tassage pour une tâche prête. Le fait de réduire l'allocation de la tâche 3 de 8 processeurs à 6 processeurs permet d'avoir une meilleure date de fin d'exécution pour cette tâche.



**Figure 7.** Illustration du placement avec tassage.

Dans la section suivante, nous allons comparer les performances des versions modifiées de CPA avec l'algorithme original en milieu homogène. Nous regarderons également comment l'algorithme se comporte lorsque nous combinons les deux modifications que nous venons de décrire.

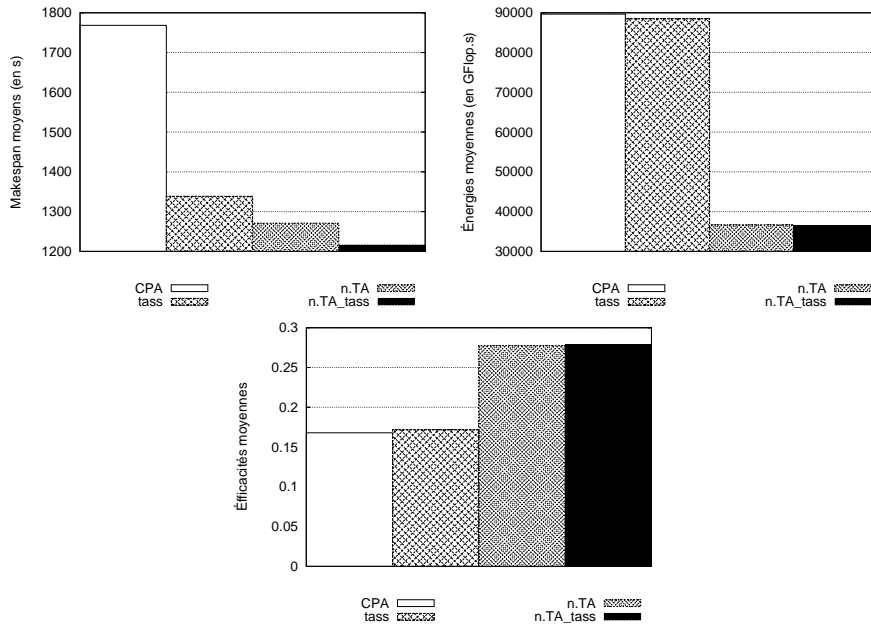
### 5.3. Évaluation des algorithmes

La figure 8 compare les moyennes du temps de complétion, de l'énergie consommée et de l'efficacité des algorithmes dans les cas suivants :

- CPA : l'algorithme original.
- tass : tassage pendant le placement.
- n.TA : utilisation de la nouvelle définition de  $T_A$ .
- n.TA\_tass : combinaison de la nouvelle définition de  $T_A$  et du tassage.

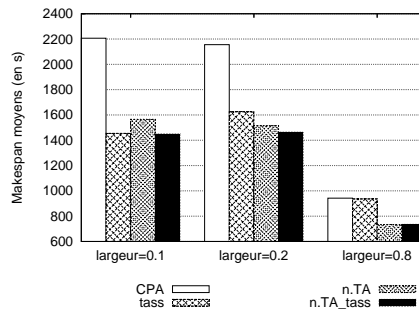
Nous observons qu'en moyenne, le tassage permet de réduire le temps de complétion des applications de 24% et de gagner légèrement en consommation de ressources par rapport à CPA. En effet, cette amélioration de l'heuristique de placement permet de réduire les trous dans l'ordonnancement en diminuant légèrement l'allocation de certaines tâches. En ce qui concerne l'utilisation de la nouvelle définition de  $T_A$ , non seulement elle permet de gagner de manière significative sur le temps de complétion des applications (réduction de 28%), mais aussi elle consomme considérablement moins de ressources par rapport à CPA (–60%). Ceci confirme le fait qu'il est important d'arrêter les allocations suffisamment tôt afin de pouvoir exécuter les tâches prêtes en parallèle. La combinaison des deux améliorations fournit à la fois un algorithme plus efficace et un gain plus important sur la moyenne des temps de complétion.

Lorsque nous regardons l'évolution du temps de complétion des algorithmes en fonction de la largeur des DAGs (figure 9), il apparaît que pour les DAGs peu larges



**Figure 8.** Performances globales sur une grappe.

(*largeur* = 0,1) le tassage fournit un meilleur temps de complétion par rapport à l'utilisation du nouveau critère d'arrêt. Cela s'explique par le fait que pour les graphes quasi filiformes, il existe peu de tâches pouvant s'exécuter en parallèle. Par contre dès qu'il existe suffisamment de tâches concurrentes dans les applications (*largeur* > 0,2), le nouveau critère d'arrêt de la procédure d'allocation permet de gagner en temps de complétion par rapport au tassage qui améliore de moins en moins l'algorithme original.



**Figure 9.** Makespan moyens en fonction de la largeur des graphes.

Dans la section suivante, nous adaptons l'algorithme obtenu en combinant les deux améliorations au cas des plates-formes hétérogènes. Nous mettrons en place deux algorithmes qui diffèrent par leur phase de placement.

## 6. Adaptation aux plates-formes hétérogènes

La plate-forme cible étant maintenant constituée de  $C$  grappes, notre idée consiste à attribuer, lors de la première phase, une *allocation de référence* à chaque tâche qui représente ses  $C$  allocations potentielles. Nous définirons comment déterminer le nombre de processeurs à allouer à une tâche sur une grappe en fonction de son allocation de référence dans la section 6.1. Lors de la phase de placement nous retiendrons parmi ces allocations potentielles celle qui minimise la date de fin d'exécution de la tâche.

Étant donné qu'il existe désormais plusieurs allocations possibles pour une même tâche, il convient de redéfinir formellement les notions de chemin critique et d'aire moyenne qui déterminent la condition d'arrêt de la phase d'allocation. Pour cela nous introduisons la notion de *grappe de référence* sur laquelle nous ferons évoluer l'allocation des processeurs. Cette grappe de référence est une plate-forme homogène virtuelle ayant une puissance de calcul cumulée équivalente à celle de l'ensemble de la plate-forme réelle et dont les processeurs ont la plus petite vitesse de la plate-forme initiale. Le nombre total de processeurs contenus dans la grappe de référence est donc :

$$P_{ref} = \left\lceil \sum_{k=0}^{C-1} \frac{P_k}{r_k} \right\rceil \quad [7]$$

où  $r_k$  est le rapport de la puissance d'un processeur de la grappe de référence sur celle d'un processeur de la grappe  $C_k$ . L'utilisation de cette grappe homogène virtuelle permet de conserver une faible complexité dans le nouvel algorithme. Notons  $p^{ref}(t)$ , l'allocation de référence pour la tâche  $t$ , *i.e.*, le nombre de processeurs qui lui seraient attribués sur la grappe de référence. L'allocation de référence est définie de sorte que le temps d'exécution effectif de chaque tâche soit relativement proche du temps qu'elle mettrait sur la grappe de référence :  $T^{ref}(t, p^{ref}(t))$ . La longueur du chemin critique ( $T_{CP}$ ) et  $T_A$  se définissent maintenant par rapport aux allocations de référence :

$$T_{CP} = \max_{t \in \mathcal{N}} T_b^{ref}(t), \quad [8]$$

$$T_A = \frac{1}{\min\{P_{ref}, \sqrt{P_{ref} \times N}\}} \sum_{i=0}^N (T^{ref}(t_i, p^{ref}(t_i)) \times p^{ref}(t_i)). \quad [9]$$

où  $T_b^{ref}(t)$  est le *bottom level* calculé à partir des allocations de référence des tâches.

Les deux sections qui suivent décrivent les deux étapes des algorithmes que nous avons conçus.

### 6.1. Allocation de processeurs

Comme dans CPA, à chaque itération de la procédure d'allocation, la tâche du chemin critique qui en bénéficie le plus se voit attribuer un processeur supplémentaire. Dans nos algorithmes ce processeur est ajouté à son allocation de référence. Pour déterminer le nombre de processeurs à allouer à une tâche sur une grappe donnée, d'après son allocation de référence, nous utilisons la loi d'Amdahl. L'égalité :

$$T^k(t, p^k(t)) = T^{ref}(t, p^{ref}(t))$$

conduit à :

$$f(p^{ref}(t), t, k) = \frac{(1 - \alpha) \cdot T^k(t, 1)}{T^{ref}(t, p^{ref}(t)) - \alpha \cdot T^k(t, 1)},$$

$f$  étant la fonction qui permet de déduire  $p^k(t)$ . Puisque l'on ne peut allouer plus de processeurs qu'il n'en existe sur une grappe et que le nombre de processeurs alloués doit être un nombre entier, on en déduit :

$$p^k(t) = \min\{P_k, \lceil f(p^{ref}(t), t, k) \rceil\} \quad [10]$$

Pour éviter une boucle infinie dans cette procédure, nous définissons une condition d'arrêt supplémentaire qui est la notion de *chemin critique saturé*. Le chemin critique est dit saturé si les allocations de référence sont telles qu'il est impossible de rajouter des processeurs aux tâches qui le composent : le nombre de processeurs à leur allouer sur toute grappe  $C_k$  est le nombre total de processeurs de cette grappe ( $P_k$ ). Les temps d'exécution des tâches du chemin critique ne peuvent donc plus être réduits en augmentant leurs allocations de référence.  $T_{CP}$  est alors minimal et nous arrêtons la procédure dès que cet état est atteint.

L'algorithme 1 présente notre version hétérogène de la phase d'allocation.

### 6.2. Placement des tâches

Le placement consiste à attribuer à chaque tâche prête choisie  $t$ , les processeurs qui lui garantissent la plus petite échéance  $T_f(t)$ . Notons  $T_r^k(t_i, t_j)$  le coût de la redistribution des données lorsque l'on passe d'une tâche  $t_i$  qui vient de s'exécuter sur un ensemble de processeurs connus à l'exécution d'une de ses tâches filles  $t_j$  sur une grappe  $C_k$ . Ce coût dépend entre autres des caractéristiques du réseau, de la quantité des données à transférer et des nombres de processeurs alloués aux tâches  $t_i$  et  $t_j$ .



**Algorithme 1** Allocation de processeurs

---

```

1: pour tout  $t \in \mathcal{N}$  faire
2:    $p^{ref}(t) \leftarrow 1$ ;
3:    $p^k(t) \leftarrow 1, \forall k \in [0, C - 1]$ ;
4: fin pour
5: tant que  $T_{CP} > T_A$  et chemin critique non saturé faire
6:    $t \leftarrow$  tâche du chemin critique /  $(\exists C_k / [f(p^{ref}(t), t, k)] < P_k)$  et
7:    $\left( \frac{T^{ref}(t, p^{ref}(t))}{p^{ref}(t)} - \frac{T^{ref}(t, p^{ref}(t)+1)}{p^{ref}(t)+1} \right)$  est maximum;
8:    $p^{ref}(t) \leftarrow p^{ref}(t) + 1$ ;
9:   Mettre à jour les  $T_b^{ref}$ ;
10: fin tant que

```

---

Soit  $T_m^k(t)$  la date d'arrivée du dernier message d'une tâche prête  $t$  si celle-ci devait s'exécuter sur la grappe  $C_k$ . L'on a :

$$T_m^k(t) = \max_{t_j \in Pred(t)} (T_f(t_j) + T_r^k(t_j, t)), \quad [11]$$

où  $Pred(t)$  est l'ensemble des prédécesseurs de  $t$ . La date à laquelle la tâche  $t$  peut effectivement démarrer son exécution est donc :

$$T_s^k(t) = \max\{dispo(p^k(t)), T_m^k(t)\} \quad [12]$$

où  $dispo(p^k(t))$  est la date à laquelle la grappe  $C_k$  aura au moins  $p^k(t)$  processeurs libres.

Les allocations potentielles des tâches définies, nous avons étudié deux politiques différentes pour le placement des tâches.

Dans la première nous adaptons l'algorithme d'ordonnancement de liste de CPA au cas où l'on dispose de plusieurs allocations possibles pour chaque tâche. Comme dans CPA, la tâche prête la plus prioritaire est celle qui a le *bottom level* le plus élevé. Une fois que cette tâche  $t$  est déterminée, nous choisissons l'allocation qui minimise sa date de fin d'exécution :

$$T_f(t) = \min_k (T_s^k(t) + T^k(t, p^k(t))) \quad [13]$$

Nous en déduisons  $C_k$ , la grappe sur laquelle son exécution est prévue ainsi que sa date de début d'exécution :  $T_s(t) = T_s^k(t)$ . La combinaison de la procédure d'allocation et cet algorithme de placement donne lieu à l'heuristique HCPA (*Heterogeneous Critical Path and Area-based*).

Dans notre seconde heuristique d'ordonnancement, la tâche la plus prioritaire est déterminée selon le principe de l'heuristique *sufferage* (Casanova *et al.*, 2000). Ce principe consiste à choisir parmi les tâches prêtes celle qui serait la plus pénalisée s'il lui était attribué sa deuxième meilleure allocation au lieu de la première. Cette tâche est donc celle qui accuse la plus grande différence entre les deux dates de fin

d'exécution prédites. Cette heuristique ne tient donc pas compte du chemin critique. Le but de sa mise en œuvre est de voir s'il peut être parfois intéressant d'utiliser des heuristiques de placement autres que celles fondées sur la réduction du chemin critique. Nous obtenons ainsi l'algorithme S-HCPA (*Sufferage-based Heterogeneous Critical Path and Area-based*). Ici également, le placement de la tâche prête la plus prioritaire vise à minimiser sa date de fin d'exécution en lui attribuant sa meilleure allocation.

### 6.3. Complexité de HCPA et S-HCPA

Soit le nombre de processeurs maximum qu'on pourrait attribuer à une tâche sur la grappe de référence :

$$K = \max_{(k,t)} \lceil f^{-1}(P_k, t, k) \rceil \quad [14]$$

Dans le pire des cas, il faudrait  $K$  itérations pour chaque tâche dans la procédure d'allocation. D'où un total de  $K \times N$  itérations de la procédure d'allocation. la complexité du corps de la boucle (choix du chemin critique, calcul de  $T_{CP}$ ,  $T_b$ , et  $T_A$ ) est de l'ordre de  $O((N + E)C)$ . Nous en déduisons la complexité de la phase d'allocation :  $O(N(N + E)C \times K)$ .

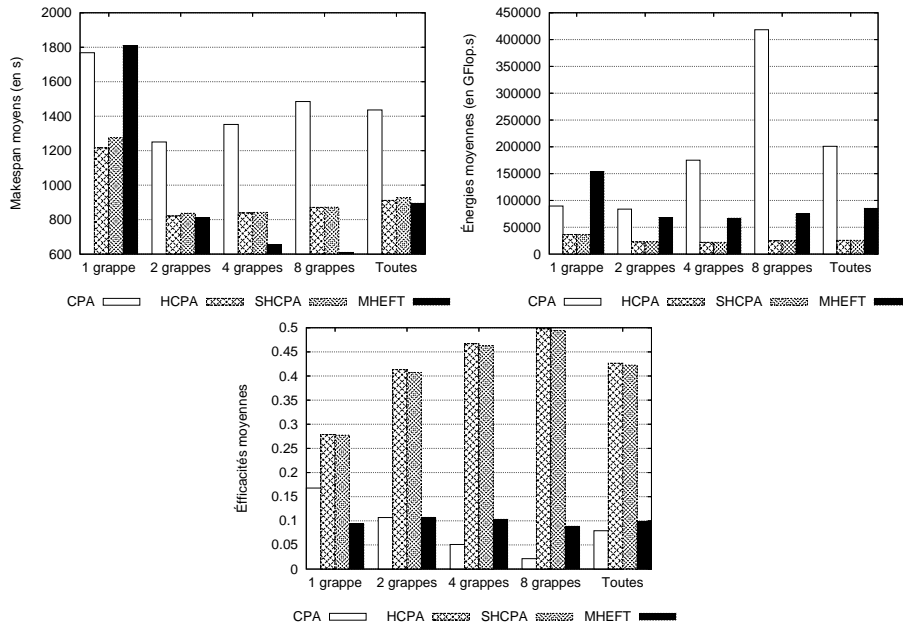
En ce qui concerne la phase de placement, dans le cas de HCPA, les tâches prêtes sont triées par ordre de priorité à l'aide d'une procédure de l'ordre de  $O(N^2)$  dans le pire des cas (cas où toutes les tâches du DAG sont prêtes en même temps du premier coup). Pour chaque tâche prête à placer, l'on teste les  $C$  allocations possibles. Au total, la procédure de placement de HCPA est de l'ordre de  $O(N(N + C))$ . Cette complexité est négligeable par rapport à celle de la phase d'allocation. Pour la phase de placement de S-HCPA, dans le pire des cas, toutes les  $N$  tâches sont prêtes en même temps et la détermination de la  $i^{eme}$  tâche la plus prioritaire a un coût de l'ordre de  $C \times (N - i)$ . En effet l'on examine les  $C$  allocations possibles pour chacune des  $N - i$  tâches prêtes pas encore placées. La complexité de la procédure de placement de S-HCPA est donc de l'ordre de  $C \times N^2$ . Cette complexité peut également être négligée par rapport à la celle de la phase d'allocation.

Il en découle que la complexité de HCPA et S-HCPA est de l'ordre de

$$O(N(N + E)C \times K). \quad [15]$$

### 6.4. Évaluation des algorithmes

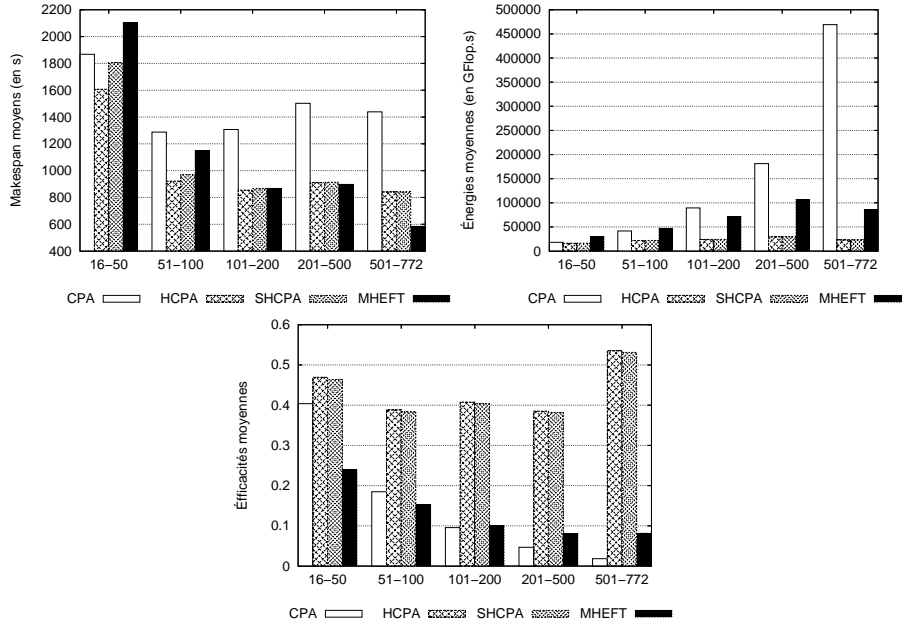
Dans cette section, nous comparons les deux algorithmes obtenus (HCPA et S-HCPA) à CPA et MHEFT (Casanova *et al.*, 2004) : l'heuristique qui adapte un algo-



**Figure 10.** *Impact du nombre de grappes.*

rythme d'ordonnement de tâches séquentielles sur plates-formes hétérogènes au cas de DAGs de tâches parallèles. En vue de comparer CPA à nos algorithmes, nous générons une plate-forme homogène (en terme de vitesse des processeurs) équivalente à la plate-forme hétérogène considérée. Cette plate-forme équivalente possède la même structure que la plate-forme initiale mais tous les processeurs ont comme vitesse, la vitesse moyenne des processeurs de la plate-forme hétérogène.

La figure 10 montre les performances moyennes des algorithmes en fonction du nombre de grappes (1, 2, 4 ou 8) et sur la totalité des plates-formes générées (toutes). En regardant les moyennes sur l'ensemble des plates-formes, nous observons que nos deux heuristiques sont en moyenne plus performantes que CPA par rapport au *makespan* (−36% environ) bien que ce dernier, qui ne restreint pas l'exécution d'une tâche à l'intérieur d'une même grappe, soit favorisé par le modèle d'Amdhal, ce modèle ne tenant pas compte des communications intra-tâches. M-HEFT est légèrement meilleur au niveau du *makespan* mais n'utilise le parallélisme de tâches que si la plate-forme comprend plusieurs grappes. En effet, le modèle utilisé fait que sur une seule grappe, toutes les tâches sont placées sur tous les processeurs l'une après l'autre. L'exécution en parallèle des tâches fait que nos heuristiques fournissent en moyennes un meilleur *makespan* par rapport à M-HEFT sur une grappe. En outre, du fait que nous allouons des processeurs supplémentaires qu'aux tâches du chemin critique et que nous arrêtons les allocations plus tôt que CPA, nos deux heuristiques consomment beaucoup



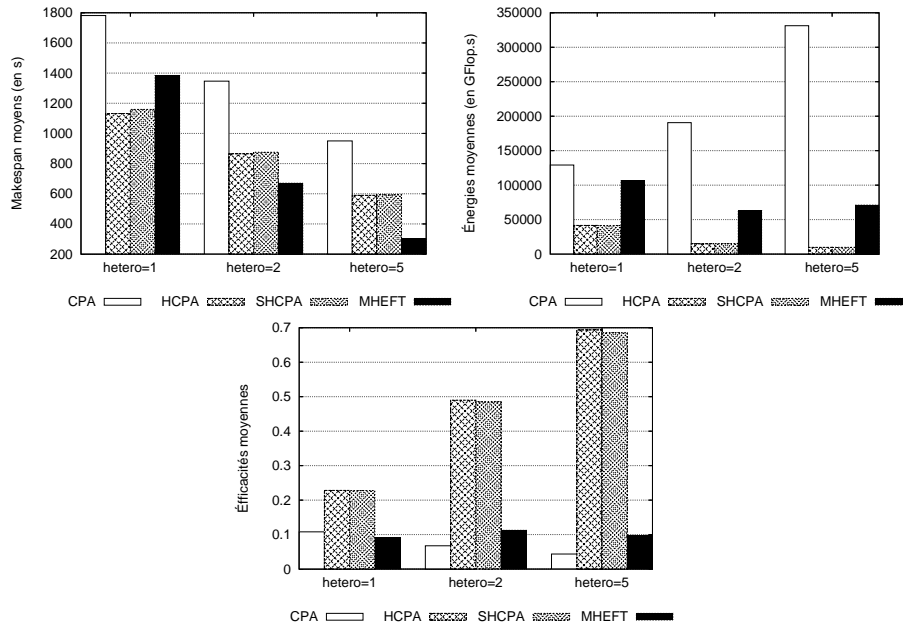
**Figure 11.** Impact du nombre de nœuds des plates-formes.

moins de ressources que CPA et M-HEFT quelque soit le nombre de processeurs de la plate-forme. Ainsi, HCPA et S-HCPA permettent de consommer en moyenne 87% moins d'énergie par rapport à CPA et 70% moins d'énergie par rapport à M-HEFT. Cela se traduit par une efficacité relativement importante pour nos deux algorithmes.

Dans la figure 11, nous avons regroupé les plates-formes en fonction du nombre de processeurs. Pour un nombre de processeurs  $P$  compris entre 16 et 50, Tous les DAGs ont leur nombre de tâches  $N$  proches de  $P$ . CPA et HCPA ont alors des allocations similaires. En revanche, dès que  $P \gg N$ , nous pouvons observer l'intérêt de notre contribution sur le calcul de  $T_A$ .

Dans la figure 12, lorsque l'hétérogénéité des plates-formes est égale à 1, nous avons encore une illustration de l'impact du changement du calcul de  $T_A$ . Ces plates-formes sont homogènes en termes de vitesse de processeurs mais peuvent être composées de plusieurs grappes et d'un grand nombre de processeurs. CPA pouvant gérer des allocations multi-sites, il obtient des performances moins bonnes. Lorsque l'hétérogénéité augmente, l'écart se creuse entre nos heuristiques et CPA aussi bien vis à vis du *makespan* que sur l'utilisation des ressources.

Le tableau 1 illustre les performances relatives des algorithmes par rapport à un ordonnancement séquentiel où nous avons calculé les moyennes sur la totalité des simulations. L'on peut constater que HCPA et S-HCPA sont largement plus efficaces



**Figure 12.** Impact de l'hétérogénéité des plates-formes.

que CPA et M-HEFT, l'efficacité étant le ratio entre le gain en temps de complétion et l'utilisation des ressources par rapport au meilleur algorithme séquentiel.

Algorithme	Accélération % SEQ	Efficacité
CPA	6,85	7,94%
HCPA	9,82	42,65%
S-HCPA	9,78	42,26%
MHEFT	11,23	9,85%

**Tableau 1.** Performances relatives par rapport à SEQ.

Enfin pour ce qui est de la comparaison HCPA et S-HCPA, il apparaît que HCPA est en moyenne légèrement meilleure au niveau des *makespan*. Mais une observation plus fine des résultats de simulations révèle que S-HCPA produit un meilleur *makespan* que HCPA dans 16,9% des cas et, un *makespan* égal à celui de HCPA dans 45,5% des cas. Cela signifie que les heuristiques d'ordonnancement de liste fondées sur la réduction du chemin critique ne minimisent pas toujours le temps de complétion des applications parallèles et donc que d'autres stratégies peuvent être explorées.

La figure 13 montre les résultats de la simulation des algorithmes sur une grappe de grappes tirée d'une plate-forme réelle. Cette plate-forme est un sous-ensemble de

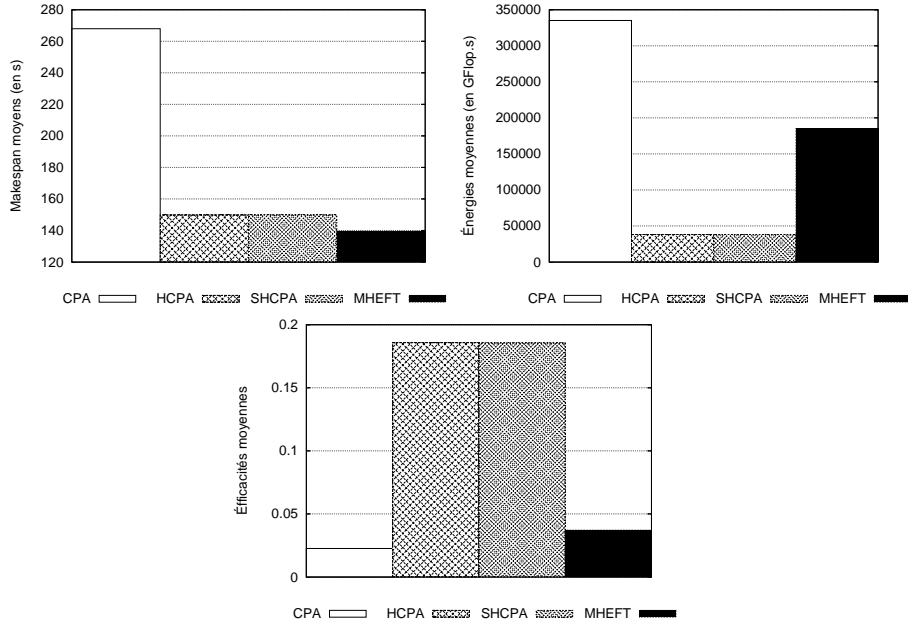


Figure 13. Comparaison sur un sous-ensemble de grid'5000.

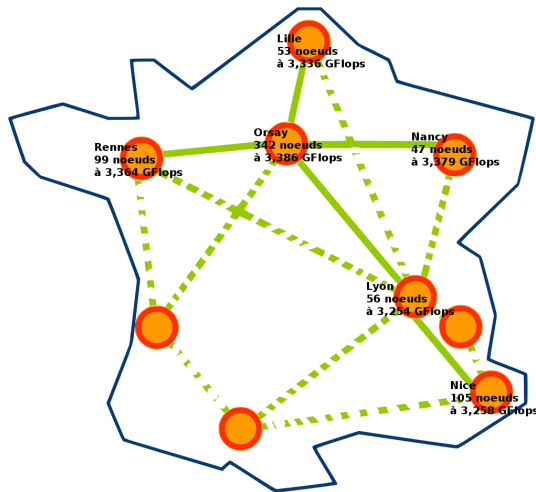


Figure 14. Exemple de grappe de grappes extraite de grid'5000.

la grille expérimentale *grid'5000*<sup>2</sup> (voir figure 14). Nous avons retenu 6 grappes sur

2. <https://www.grid5000.fr>

les sites de Lille, Lyon, Nancy, Nice, Orsay et Rennes. Elle comprend au total 702 processeurs et 5 Backbones à  $10Gb/s$ . Sur cette plate-forme de faible hétérogénéité (tous les processeurs ont une vitesse environ égale à  $3,3GFlops$ ) où le nombre de processeurs ( $P = 702$ ) est très élevé par rapport au nombre de tâches des applications ( $N \leq 50$ ), nous constatons que les efficacités de M-HEFT et CPA sont très faibles par rapport à HCPA et S-HCPA. Or HCPA et SHCPA fournissent des temps de complétion comparables à ceux de M-HEFT. Cela confirme l'utilité de notre nouveau critère d'arrêt de la phase d'allocation.

## 7. Conclusion

Dans cet article, nous avons adapté un algorithme d'ordonnancement de tâches parallèles sur plates-formes homogènes au cas des plates-formes hétérogènes. Pour cela nous avons choisi un algorithme d'ordonnancement en deux étapes de la littérature : CPA (Radulescu *et al.*, 2001b) dont les performances sont relativement bonnes par rapport aux heuristiques concurrentes.

Dans un premier temps nous avons proposé deux améliorations de CPA dans le cas des plates-formes homogènes. La première concerne la phase d'allocation et la seconde la phase de placement. Dans la phase d'allocation, nous avons redéfini la notion d'aire moyenne afin d'arrêter plus tôt la procédure d'allocation de processeurs et de profiter au mieux de l'exécution en parallèle des tâches concurrentes. Nous avons expérimentalement vérifié que cette modification permet en moyenne de réduire le temps de complétion des applications tout en utilisant peu de ressources par rapport à CPA. Dans la phase de placement, nous avons mis en place une technique de *tassage* dont le but est de réduire le temps d'attente et la date de fin d'exécution de la tâche prête à placer en diminuant son allocation si cela est nécessaire. Nous avons également vérifié que cette amélioration réduit le temps de complétion des applications par rapport à CPA. De plus la combinaison de ces deux améliorations fournit un léger gain par rapport à chacune des améliorations prise individuellement.

Ensuite, nous avons adapté l'algorithme obtenu en combinant les deux améliorations au cas de plates-formes hétérogènes de type grappe hétérogène de grappes homogènes. À cet effet, nous avons introduit la notion de *grappe de référence* qui nous a permis de mieux gérer l'hétérogénéité des plates-formes dans la procédure d'allocation. Le but de cette virtualisation des plates-formes était surtout de conserver une faible complexité pour nos heuristiques. Ainsi, nous avons mis en place deux heuristiques qui se distinguent par leur phase de placement : HCPA et S-HCPA. Dans HCPA, nous adaptons la technique de placement de CPA au cas où nous disposons de plusieurs grappes, et donc de plusieurs allocations possibles. La phase de placement de S-HCPA est quant à elle basée sur l'heuristique *sufferage* (Casanova *et al.*, 2000).

La complexité de nos deux heuristiques reste de l'ordre de celle de CPA. Leur évaluation en milieu hétérogène révèle qu'elles sont plus performantes que CPA et qu'elles consomment beaucoup moins de ressources par rapport à M-HEFT (Casanova

*et al.*, 2004) tout en menant à des temps de complétion en moyenne comparables à ceux de M-HEFT.

## 8. Bibliographie

- Amdahl G., « Validity of the Single Processor Approach to Achieving Large Scale Computing Capabilities », *AFIPS 1967 Spring Joint Computer Conference*, vol. 30, p. 483-485, April, 1967.
- Boudet V., Desprez F., Suter F., « One-Step Algorithm for Mixed Data and Task Parallel Scheduling Without Data Replication », *Proceedings of the 17th International Parallel and Distributed Processing Symposium (IPDPS'03)*, April, 2003.
- Casanova H., Desprez F., Suter F., « From Heterogeneous Task Scheduling to Heterogeneous Mixed Parallel Scheduling », *10th Int. Euro-Par Conference*, vol. 3149 of LNCS, Springer, Pisa, Italy, p. 230-237, August, 2004.
- Casanova H., Legrand A., Zagorodnov D., Berman F., « Heuristics for Scheduling Parameter Sweep Applications in Grid Environments », *9th Heterogeneous Computing Workshop (HCW'00)*, Cancun, Mexico, p. 349-363, 2000.
- Chretienne P., « Task Scheduling Over Distributed Memory Machines », *Parallel and Distributed Algorithms*, North Holland, p. 165-176, 1988.
- Dutot P.-F., « Hierarchical Scheduling for Moldable Tasks », *11th Int. Euro-Par Conference*, vol. 3648 of LNCS, Springer, Lisbon, Portugal, p. 302-311, August, 2005.
- Legrand A., Marchal L., Casanova H., « Scheduling Distributed Applications : The SimGrid Simulation Framework », *3rd IEEE Symposium on Cluster Computing and the Grid (CC-Grid'03)*, Tokyo, p. 138-145, May, 2003.
- Lepère R., Trystram D., Woeginger G., « Approximation Algorithms for Scheduling Malleable Tasks Under Precedence Constraints », *IJFCS*, vol. 13, n° 4, p. 613-627, 2002.
- Radulescu A., Nicolescu C., van Gemund A., Jonker P., « CPR : Mixed Task and Data Parallel Scheduling for Distributed Systems », *15th International Parallel and Distributed Processing Symposium (IPDPS)*, apr, 2001a. Best Paper Award.
- Radulescu A., van Gemund A., « A Low-Cost Approach towards Mixed Task and Data Parallel Scheduling », *15th International Conference on Parallel Processing (ICPP)*, Valencia, Spain, September, 2001b.
- Ramaswamy S., Sapatnekar S., Banerjee P., « A Framework for Exploiting Task and Data Parallelism on Distributed Memory Multicomputers », *IEEE Transactions on Parallel and Distributed Systems*, vol. 8, n° 11, p. 1098-1116, Nov, 1997.
- Rauber T., Rünger G., « Compiler Support for Task Scheduling in Hierarchical Execution Models », *Journal of Systems Architecture*, vol. 45, p. 483-503, 1998.
- Turek J., Wolf J., Yu P., « Approximate Algorithms for Scheduling Parallelizable Tasks », *4th ACM Symp. on Parallel Algorithms and Architectures*, p. 323-332, 1992.