

Évaluation et critères de performances d'un calcul parallèle

Mesurer les performances
et optimiser un
algorithme parallèle

Introduction

- Programme séquentiel : évaluation de la performance par la mesure du temps d'exécution;

- Programme parallèle : évaluation plus complexe
Les performances d'un code parallèle vont dépendre :
 - du nombre de processeurs;
 - de la machine sur laquelle il va s'exécuter;
 - et bien d'autres critères (cf plus loin)...

- En outre,
 - Un programme séquentiel rapide peut s'avérer, une fois parallélisé, très peu efficace ;
 - Inversement, un programme séquentiel "moins rapide" peut s'avérer mieux parallélisable;

Mesures des performances

Charge de travail d'un algorithme parallèle

➤ La charge de travail d'un algorithme séquentiel $W_{\text{séquentiel}}$ se décompose en deux parties :

$$W_{\text{séquentiel}} = \underbrace{W_{\text{fixe}}}_{\text{Travail non parallélisable (coût fixe)}} + \underbrace{W_{\text{parallélisable}}}_{\text{Travail pouvant être parallélisé}}$$

➤ La charge de travail d'un algorithme parallèle $W_{\text{parallèle}}$ avec p processeurs est donc :

$$W_{\text{parallèle}}(p) = W_{\text{fixe}} + W_{\text{parallélisable}}/p$$

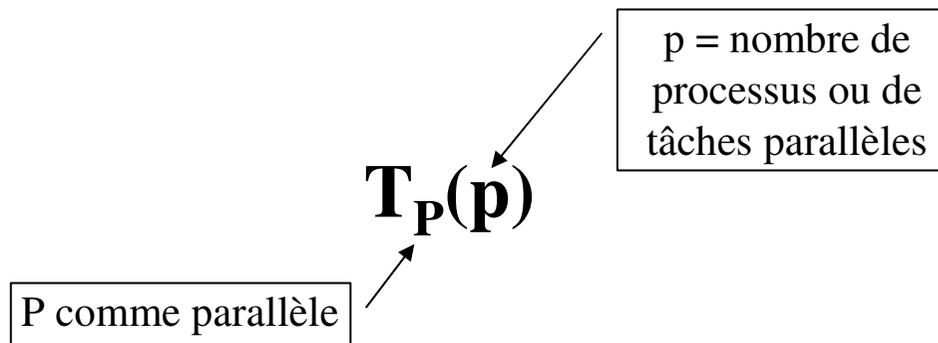
➤ Dans un "vrai" programme parallèle, il faut prendre en compte une surcharge de travail W_{com} due aux synchronisations des tâches :

$$W_{\text{parallèle}}(p) = W_{\text{fixe}} + W_{\text{parallélisable}}/p + W_{\text{com}}$$

Temps d'exécution (*elapsed time*)

➤ Définition : temps écoulé entre le début du calcul parallèle et la terminaison de la dernière tâche;

➤ Notation :



Coût total (*cost*)

- Définition : somme des coûts de tous les processeurs;
- Mesure un temps "équivalent séquentiel" ;
- Notation :

$$C_T(\mathbf{p}) = \mathbf{p} \times T_P(\mathbf{p})$$

Surcoût parallèle (*Parallel Overhead*)

$$T_O(p) = \underbrace{pT_P(p)}_{\text{Temps consommé par tous les processeurs}} - \underbrace{T_S}_{\text{Temps séquentiel}}$$

- Mesure le surcoût en temps total consommé du calcul parallèle par rapport au calcul séquentiel;
- Normalement, un calcul parallèle coûte plus cher qu'un calcul séquentiel :

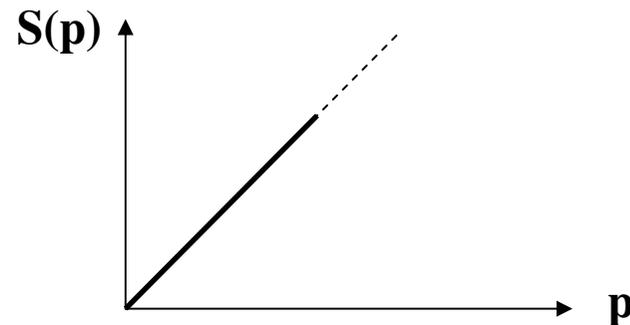
Démonstration : hypothèse : temps T proportionnel au coût W

$$\begin{aligned} T_P(p) &= T_{\text{fixe}} + T_{\text{com}} + T_{\text{parallélisable}}/p \\ \Rightarrow pT_P(p) &= pT_{\text{fixe}} + pT_{\text{com}} + T_{\text{parallélisable}} \\ &= (p-1)T_{\text{fixe}} + pT_{\text{com}} + T_{\text{fixe}} + T_{\text{parallélisable}} \\ &> 0 + T_S \\ \Rightarrow T_O(p) &> 0 \end{aligned}$$

Accélération Relative (*speed-up*)

$$S(\mathbf{p}) = T_P(\mathbf{1}) / T_P(\mathbf{p})$$

- Mesure l'accélération, le bénéfice d'un calcul parallèle par rapport au même algorithme sur un processeur;
- On dit qu'un programme est extensible (*scalable*) si et seulement si $S(\mathbf{p}) \sim \mathbf{p}$ (ou encore mieux si $S(\mathbf{p}) = \mathbf{p}$) :
 - Difficile à obtenir ;
 - Souvent extensible que sur un intervalle de \mathbf{p} ;



À suivre...

...Suite

Accélération Relative (*speed-up*)

- En théorie,
 - $S(p) \leq p$
 - Autrement dit, un calcul parallèle sur p processeurs ne peut pas s'exécuter p fois plus vite que le meilleur algorithme séquentiel ;

- Démonstration :

$$\begin{aligned} S(p) &= \frac{T_{\text{fix}} + T_{\text{par}}}{T_{\text{fix}} + T_{\text{par}}/p + T_{\text{com}}} \\ &= \frac{p(T_{\text{fix}} + T_{\text{par}})}{T_{\text{fix}} + T_{\text{par}} + (p-1)T_{\text{fix}} + pT_{\text{com}}} \\ &= \frac{p}{1 + ((p-1)T_{\text{fix}} + pT_{\text{com}})/(T_{\text{fix}} + T_{\text{par}})} \leq \frac{p}{1+0} = p \quad (*) \end{aligned}$$

À suivre...

...Suite

Accélération Relative (*speed-up*)

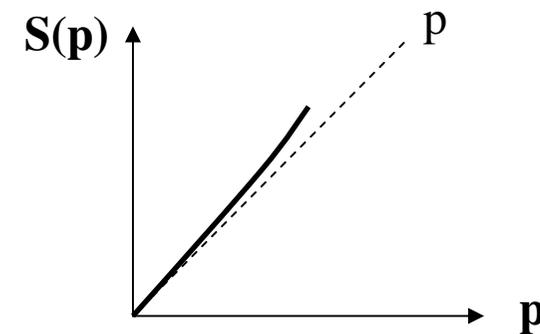
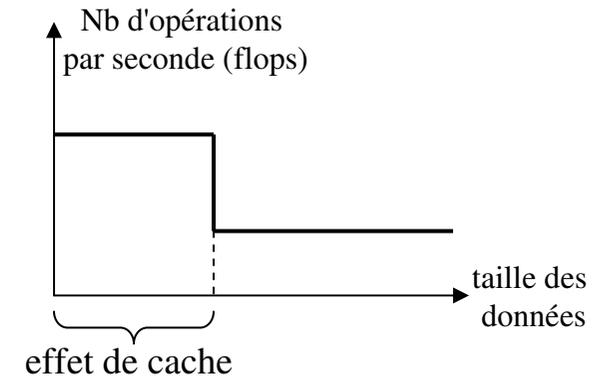
- En pratique, il existe des cas où $S(p) > p$!
 - ⇒ On dit alors que le programme est **superscalaire**;
- Pourquoi un programme peut être superscalaire ?
 - ⇒ Grâce aux effets de cache

Effets de cache en séquentiel :

la vitesse de traitement d'un vecteur "long"
est moins élevée que celle d'un vecteur "court"
Conséquence : l'hypothèse T proportionnel à
 W est fautive !

Effets de cache en parallèle :

- Pour un problème de taille fixée,
lorsque le nombre de processeurs augmente,
la taille des données par processeur diminue ;
- le calcul séquentiel manipule des données
volumineuses;
- à partir d'un certain nombre de processeurs,
le calcul parallèle manipule des données qui tiennent
dans les caches des cpus, d'où une suraccélération;



À suivre...

...Suite

Accélération Relative (*speed-up*)

➤ En outre,

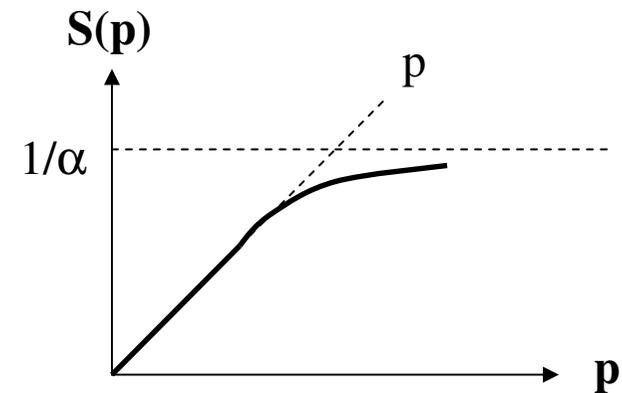
$$\forall p, \quad S(p) \leq 1/\alpha \quad (\text{Loi d'Amdahl}),$$

où $\alpha = T_{\text{fix}}/T_S$ est la part du travail non parallélisable sur le travail total

Autrement dit, l'accélération d'un algorithme parallèle est bornée par une constante, et ceci quelque soit le nombre de processeurs

➤ Démonstration : d'après (*),

$$\begin{aligned} S(p) &= \frac{p}{1 + ((p-1)T_{\text{fix}} + pT_{\text{com}})/(T_{\text{fix}} + T_{\text{par}})} \\ &= \frac{p}{1 + ((p-1)T_{\text{fix}} + pT_{\text{com}})/T_S} \quad (**) \\ &\leq \frac{p}{1 + (p-1)\alpha} = \frac{1}{1/p + (1-1/p)\alpha} \\ &\rightarrow \frac{1}{\alpha} \quad \text{pour } p \rightarrow \infty \end{aligned}$$



À suivre...

...Suite

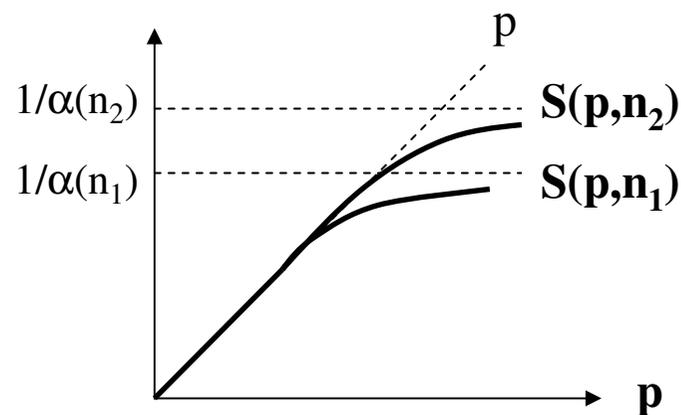
Accélération Relative (*speed-up*)

➤ Pour être exact, l'accélération dépend du nombre de processeurs mais également de la taille des données \mathbf{n} à traiter

$$\mathbf{W}_{\text{séquentiel}}(\mathbf{n}) = \mathbf{W}_{\text{fixe}} + \mathbf{W}_{\text{parallélisable}}(\mathbf{n})$$

En règle générale, $\mathbf{W}_{\text{parallélisable}}(\mathbf{n})$ augmente quand \mathbf{n} augmente donc, $\alpha(\mathbf{n}) = \mathbf{W}_{\text{fixe}} / \mathbf{W}_{\text{séquentiel}}(\mathbf{n}) \rightarrow 0$ quand $\mathbf{n} \rightarrow \infty$

Pour un nombre de processeurs fixé, l'accélération sera d'autant meilleure que la taille des données à traiter sera importante :



À suivre...

...Suite

Accélération Relative (*speed-up*)

- Il ne faut pas oublier l'influence des synchronisations des tâches :

Plus il y'a de processeurs, plus le temps pris par les communications est important :

$$\mathbf{T}_{\text{com}}(\mathbf{p}) \rightarrow \infty \text{ quand } \mathbf{p} \rightarrow \infty$$

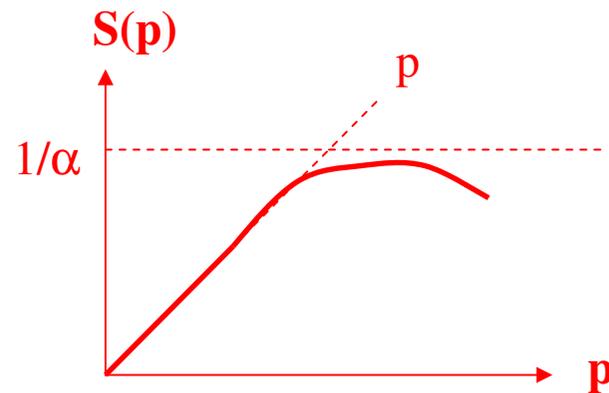
Conséquence : $\mathbf{S}(\mathbf{p}) \rightarrow 0$ quand $\mathbf{p} \rightarrow \infty$

Démonstration : d'après (**),

$$\mathbf{S}(\mathbf{p}) = \frac{1}{1/\mathbf{p} + (1 - 1/\mathbf{p})\alpha + \mathbf{T}_{\text{com}}(\mathbf{p})/\mathbf{T}_S} \underset{\mathbf{p} \rightarrow \infty}{\sim} \frac{1}{\alpha + \mathbf{T}_{\text{com}}(\mathbf{p})/\mathbf{T}_S}$$

Or, $\mathbf{T}_{\text{com}}(\mathbf{p}) \rightarrow \infty$ quand $\mathbf{p} \rightarrow \infty$, donc

$$\mathbf{S}(\mathbf{p}) \rightarrow 0 \text{ quand } \mathbf{p} \rightarrow \infty$$



Accélération Absolue (*speed-up*)

$$S^*(p) = T_S / T_P(p) ,$$

où T_S est le temps du meilleur algorithme séquentiel

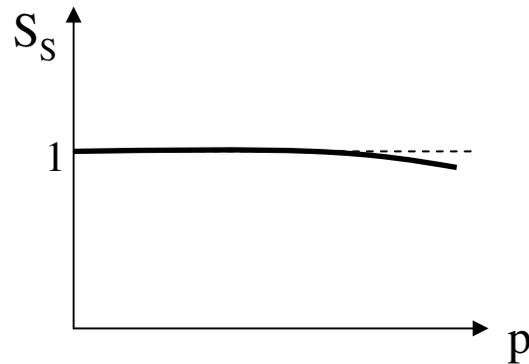
- Mesure l'accélération, le bénéfice d'un calcul parallèle par rapport au meilleur algorithme séquentiel ;
- Rappel : en règle générale, $S(p) \neq S^*(p)$ car $T_P(1)$ n'est pas forcément le temps du meilleur algorithme séquentiel;
- Il est impossible que $S^*(p) > p$ (superscalaire) sinon T_S ne serait pas le temps du meilleur algorithme séquentiel (manipulation des données par bloc pour bénéficier des effets de caches);

Scaled speed-up

$$S_s(p, n) = T_p(1, n) / T_p(p, p \times n) ,$$

où n est la taille des données à traiter

- Interprétation : si je multiplie par p le nombre de processeurs et la taille des données, est-ce que cela me prend le même temps ?
- Le programme est extensible si $S_s(p) \sim 1$;



Efficacité parallèle (*efficiency*)

$$E(p) = T_S / C_T(p)$$

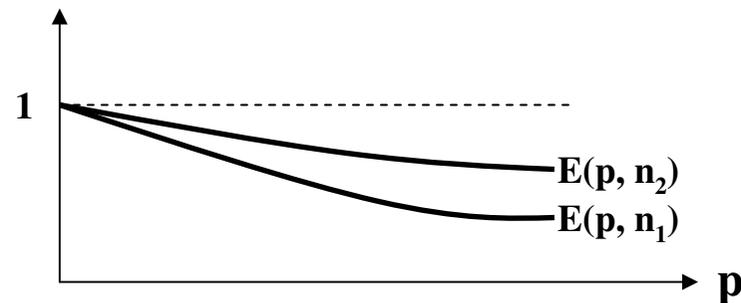
- Mesure le "rendement" du calcul parallèle;
- En théorie, l'efficacité est un nombre compris entre 0 et 1 :

$$E(p) = T_S / C_T(p) = T_S / (pT_p(p)) = S(p) / p \leq 1$$

Mais, en pratique, l'efficacité peut être supérieure à 1 dans le cas des calculs superscalaires;

- Programme extensible si $E = 1$: 100 % des ressources parallèles ont été utilisées pour du calcul "utile" ;
- $S(p) \leq 1/\alpha \Rightarrow E(p) \leq 1/(\alpha p)$

(rappel : α = part de la charge non parallélisable sur la charge totale);



Équilibrage de charge

- Mesure la "qualité" de la répartition de la charge de travail sur les processeurs;

Soit W_i la charge de travail du processeur i

Soit W_{\max} la charge maximale : $W_{\max} = \max_{i \in [1,p]} W_i$

Soit W_{tot} la charge totale de tous les processeurs : $W_{\text{tot}} = \sum_{i \in [1,p]} W_i$

Soit W_{moy} la charge moyenne par processeur : $W_{\text{moy}} = W_{\text{tot}} / p$

On définit l'équilibrage de charge $\epsilon(p)$ comme :

$$\epsilon(p) = (W_{\max} - W_{\text{moy}}) / W_{\max}$$

- Si $\epsilon(p) = 0$ (cas où $W_{\max} = W_{\text{moy}}$), alors l'équilibrage est parfait ;

Critères de performances

Critères de performances ?

Pour résoudre le même problème, le coût d'un calcul parallèle est supérieur au coût d'un calcul séquentiel;

Qu'est-ce qui engendre ce surcoût ?

- W_{fixe} : travail non parallélisable :
 - peu de marge de manœuvre;
 - de plus, la part du travail fixe sur le travail total diminue quand le nombre de données à traiter augmente;

⇒ ne pas perdre son temps à diminuer W_{fixe} !
- W_{com} : travail dû au traitement parallèle :
 - le travail augmente quand le nombre de tâches parallèle augmente;
 - le travail augmente quand le nombre de données à traiter par tâche augmente;

⇒ il faut concentrer l'effort sur W_{com} car c'est de lui que vont dépendre en partie les bonnes performances du calcul parallèle !

Quelles sont les sources du surcoût dû au traitement parallèle ?
Comment diminuer cette surcharge de travail ?

...Suite

Les communications

Conséquences (et solutions):

- chaque message a un coût fixe important \Leftrightarrow **il faut le moins de messages possibles;**
- il faut des **messages les plus longs possibles** pour minimiser l'impact des latences;
- il faut souvent **réorganiser l'algorithme séquentiel** pour minimiser le nombre de communications lors de sa parallélisation;
- lorsque c'est possible, **recouvrir les communications par du calcul** (effectuer du calcul pendant que les communications se déroulent, notion de communications non bloquantes).

Remarque sur les communications collectives : elles coûtent encore plus chers que les communications point à point (coût de l'ordre de p fois un échange, où p est le nombre de processeurs).

Calculs supplémentaires

- Calculs non présents dans l'algorithme séquentiel mais indispensables dans l'algorithme parallèle ;
- Différentes sources :
 - éviter des communications :
 - dans certains cas, il est préférable que toutes les tâches parallèles effectuent le même calcul, plutôt qu'une seule tâche suivie d'une phase de communication;
 - dans quels cas ? On évalue le coût de la phase de calcul par rapport au coût d'une communication. En règle générale, le coût de la phase de calcul ne dépend pas de la taille des données à traiter;
 - cas où l'algorithme parallélisé est différent de l'algorithme séquentiel:
 - conséquence : plus d'opérations en parallèle par processeur que pour l'algorithme séquentiel;
 - dans quels cas ? L'algorithme séquentiel est difficilement parallélisable (mauvaises accélérations) ou non parallélisable.

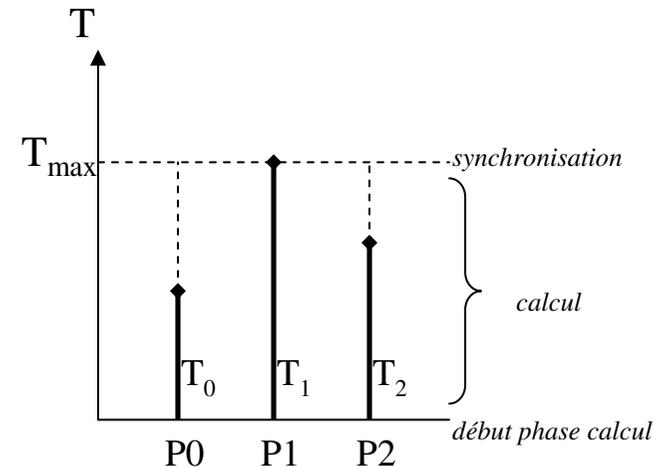
Remarque : ces calculs supplémentaires sont inévitables voire indispensables (donc, pas de solution).

Déséquilibre de charge

(*Imbalance workload*)

Que se passe t'il quand les tâches n'ont pas la même charge de travail ?

- temps de la phase parallèle : $T_{\max} = T_1$;
- coût de la phase parallèle : $C_T = 3 \times T_1$;



Qu'ont fait P0 et P2 entre le moment où ils ont terminé et le moment où P1 a terminé ?

- ils ont attendu ! (*idling*)
- temps passé dans l'attente : $\sum_i (T_{\max} - T_i) = (T_1 - T_0) + (T_1 - T_2)$;
- si on rapporte cette attente au coût total du calcul, on retrouve la mesure de l'équilibrage de charge : $\epsilon(p) = \sum_i (T_{\max} - T_i) / (p T_{\max})$;

À suivre...

...Suite

Déséquilibre de charge

(*Imbalance workload*)

- Équilibrage parfait (équidistribution du travail) \Leftrightarrow pas de perte de temps dans des attentes ;
- Un calcul est rarement parfaitement équilibré : comment répartir **n** données à traiter sur **p** tâches ?
- Un bon équilibrage sera **facile** à obtenir si :
 - **n** est grand devant **p** : on parle alors de grain fin (*fine granularity*);
 - répartition statique (i.e. ne change pas au cours du calcul);
- Un bon équilibrage sera **difficile** à obtenir si :
 - **n** est du même ordre de grandeur que **p** : on parle alors de grain grossier (*coarse granularity*);
 - la charge de travail est dynamique : la taille de la charge n'est pas prévisible car générée par l'algorithme :
 - \Leftrightarrow répartition dynamique de la charge (*dynamic load balancing*), mais coûteux (\Rightarrow calculs supplémentaires);
 - les tâches s'exécutent sur des unités de traitement de puissances différentes ou d'architectures différentes (réseaux hétérogènes) : ex : réseau de PCs avec des fréquences différentes, ou réseau constitué de machines d'architectures différentes ...

En bref ...

- Les performances d'un calcul parallèle s'évaluent à l'aide de plusieurs mesures comme l'accélération ou l'équilibrage ;
- Les performances parallèles sont essentiellement impactées par :
 - les communications : il en faut le moins possible, et leurs messages doivent être les plus grands possibles ;
 - le déséquilibre de charge induit des attentes inutiles sur les processeurs les moins chargés : veillez à une répartition équitable du travail !