

CentraleSupélec

Mineure HPC-SBD

Mesure, métriques et analyse de performances

Stéphane Vialle

Stephane.Vialle@centralesupelec.fr
http://www.metz.supelec.fr/~vialle

CentraleSupélec

Mesure, métriques et analyse de performances

- 1 – Mesure des temps d'exécution
- 2 – Métriques de performances (S, e)
- 3 – Métriques de *Size Up*
- 4 – Critères de passage à l'échelle
- 5 – Loi d'Amdahl
- 6 – Loi de Gustafson
- 7 – Liens Amdahl-Gustafson
- 8 – Roof Line Model

CentraleSupélec

Mesure, métriques et analyse de performances

1 – Mesure des temps d'exécution

Le temps d'exécution (d'un code) est souvent LA grandeur que l'on mesure, est qui entachée d'imprécisions

- Différents temps d'exécution
- Les outils de mesure
- Reproductibilité des mesures
- Stratégie de mesure

CentraleSupélec

Mesure des temps d'exécution

Méthodologie de mesures

Mesures externes :

```
>time myAppli
>/usr/bin/time myAppli
>times myAppli
>timeX myAppli
.....
```

12.002u	user
0.128s	system
12.150	total

Nom et fonctionnement variables selon le système utilisé !!

Fréquemment : total > user + system !!

Simple à utiliser
Pas de modifications des codes sources

Peu précis: $\pm 0.5s$

CentraleSupélec

Mesure des temps d'exécution

Méthodologie de mesures

Mesures internes :

```
time ()
clock ()
gettimeofday ()
omp_get_wtime ()
```

Compte les clics d'horloge

Compte le temps écoulé en s

- Toutes ces routines ne sont pas toujours disponibles !
- "gettimeofday" est en général une bonne solution.
- Parfois il existe des outils plus précis pour mesurer de petites durées.

Plus précis que les mesures externes

Besoin de modifier le code source
Pas toujours totalement portable

CentraleSupélec

Mesure des temps d'exécution

Méthodologie de mesures

Précision des outils et des mesures :

123456789012 . 1234567890123456

Capacité maximale de l'outil de mesure

Précision théorique (cf. doc)

Précision expérimentale, vu la fluctuation des exécutions

- Ne pas tenir compte de trop de décimales!
- Faire attention à ne pas déborder la capacité de mesure!

Mesure des temps d'exécution

Méthodologie de mesures

Problème fréquent :

Test en mode exclusif (mono-user).
Outil de mesure à 1ms de précision.

Fluctuation de 500ms d'une exécution à l'autre !!

Et plus encore avec la montée en fréquence automatique des pros. (effet de "chauffe")

Démarche conseillée :

- Mesurer les fluctuations, ne pas les ignorer (le *warm up* des processeurs peut perturber les premières)
- Ne pas donner que les valeurs moyennes
- Mesurer des temps > 10s si possible

Mesure des temps d'exécution

Méthodologie de mesures

Stocker des meta-données sur les conditions de mesure :

- Date de l'exécution**
- Auteur(s) du test**
- Outil(s) de mesure utilisé(s)
- Caractéristiques de la machine : RAM, Cache, Processeurs, ...
- OS utilisé (nom et version)
- Compilateur utilisé (nom et version)
- Options de compilation utilisées
- Test en multi-user/mono-user ?
- Présence d'IO dans le test ?
- Configuration du programme de test : taille des données, ...

On oublie souvent (et rapidement) à quel benchmark se réfère une série de mesures !

On manque souvent de détail sur les conditions de réalisation d'une série de mesures !

Mesure des temps d'exécution

Choix de la représentation

Quelle courbe présenter ?

- Avoir une idée de l'allure de la courbe attendue et/ou de son expression
- Choisir une représentation qui permet de visualiser "des droites" ou des formes géométriques simples (droite, cercle, angle droit...)

Exemple d'un temps d'exécution parallèle :

Cas idéal : $T(P) = T(1)/P \rightarrow$ une hyperbole

- l'hyperbole est mal identifiée par l'œil...
- ...on la confond facilement avec une courbe qui décroît selon une autre loi !

Les mesures vérifient-elles la théorie ?

Mesure des temps d'exécution

Choix de la représentation

Quelle courbe présenter ?

- Avoir une idée de l'allure de la courbe attendue et/ou de son expression
- Choisir une représentation qui permet de visualiser "des droites" ou des formes géométriques simples (droite, cercle, angle droit...)

Exemple d'un temps d'exécution parallèle :

Cas idéal : $\log(T(P)) = \log(T(1)) - \log(P)$

- en échelle log une hyperbole est très rapidement identifiée : droite de pente -1
- on détecte facilement :
 - un écart complet à la théorie
 - un écart à partir d'un seuil...

Mesure, métriques et analyse de performances

2 – Métriques de performances (S,e)

- Observation de la courbe du temps
- Métrique d'accélération (*Speedup*)
- Métrique d'efficacité (*Efficiency*)
- Métrique de passage à l'échelle (*Scaling*)
- Influence de la référence séquentielle

Métriques de performances

Observation de la courbe du temps

Difficultés pour accélérer une application (de taille fixée)

Quand l'algorithme parallèle a une complexité beaucoup plus grande que le meilleur algorithme séquentiel...

...le gain final est faible !

Quand les surcoûts de gestion du parallélisme (*Tsynchro*, *Tcomm*...) sont grands...

...l'accélération reste faible !

Métriques de performances

Observation de la courbe du temps

Courbe de temps prometteuse

Quand l'algorithme parallèle a :

- un temps d'exécution qui décroît significativement
- un surcoût limité sur une seule ressource

...le gain final sera important !

Exec Time (s)

Seq. Algo

Nb of resources

...et on poursuit l'analyse

Speedup Efficiency Size up Scalability...

Métriques de performances

Métrique d'accélération

Speedup :

$$S(P) = \frac{T(1)}{T(P)}$$

$S(P) < 1$: on ralentit !
 mauvaise parallélisation
 $1 < S(P) < P$: "normal"
 $P < S(P)$: hyper-accélération
 analyser & justifier

hyper-accélération

accélération normale

ralentissement

$S(P) = P$: accélération idéale

Métriques de performances

Métrique d'accélération

Speedup :

$$S(P) = \frac{T(1)}{T(P)}$$

$S(P) < 1$: on ralentit !
 mauvaise parallélisation
 $1 < S(P) < P$: "normal"
 $P < S(P)$: hyper-accélération
 analyser & justifier

Cas standard :

accélération idéale

accélération mesurée

Il y a toujours des sources de perte de performance...

Métriques de performances

Métrique d'accélération

Cas d'hyper-accélération :

Ce n'est pas magique, et ce n'est pas normal

- On doit analyser le phénomène et l'expliquer
- Corriger une erreur ou exploiter une optimisation

Exemples d'explications :

- on ne fait plus les bonnes opérations (résultat faux)
- les données tiennent dans le cache total des P processeurs
- on a modifié l'algorithme de départ et on converge plus vite (ex. de l'algorithme génétique optimisé !)
- on cherche une solution dans un arbre et on stoppe le pgm

P_0 P_1 $S(2) = 3$

Métriques de performances

Métrique d'efficacité

Efficacité :

$$e(P) = \frac{S(P)}{P}$$

Taux d'utilisation des ressources, ou fraction obtenue de l'accélération idéale

- $e(P) \in [0;1]$, $\in [0\%;100\%]$
- $e(P) > 100\%$ \Leftrightarrow hyper-accélération

Efficiency: $e(p)$ %

Number of computing nodes : p

L'utilisateur s'intéresse à l'accélération obtenue

L'acheteur de la machine s'intéresse à l'efficacité des applications exécutées

Le développeur s'intéresse aux deux

Métriques de performances

Choix de la référence séquentielle

A quel programme et exécution séquentielle se comparer ?

- Même programme lancé sur un seul processeur ?
Même algorithme implanté en séquentiel ?
Meilleur algorithme séquentiel connu ?
- Compilation séquentielle avec le même compilateur ?
Compilation avec le meilleur compilateur séquentiel ?
- Optimisations séquentielles autorisées par la parallélisation ?
Optimisations séquentielles maximales ?
- Exécution sur un seul processeur de la machine parallèle ?
Exécution sur la meilleure machine séquentielle ?

Métriques de performances

Choix de la référence séquentielle

Tous les choix sont plausibles :

Chaque choix de référence séquentielle correspond à :

- un point de vue différent,
- une préoccupation différente,
- un objectif d'analyse différent

→ Faire le choix correspondant à sa problématique
→ Énoncer clairement ce choix

Exemples :

Utilisateur final →
SON pgm séquentiel sur SA machine séquentielle

Développeur de code parallèle →
SON pgm parallèle sur UN proc de SA machine parallèle

Métriques de performances

Choix de la référence séquentielle

La référence séquentielle peut être obtenue avec :

Même algo, même langage, même optim seq., même proc → $S_1(P)$ **Bonnes perf faciles à obtenir**

↕

Meilleur algo, meilleur langage, meilleur optim seq., meilleur proc → $S_2(P)$ **Bonnes perf très difficiles à obtenir**

Métriques de performances

Sources de pertes de performances

- Sous-optimisation séquentielle ↔ Sous-utilisation de chaque coeur
- Impossibilité de vectoriser les boucles
- Fraction séquentielle ↔ Sous-utilisation des noeuds de calcul
- False-sharing (en mémoire partagée)
- Surcoût des synchronisations
- Surcoût des communications
- Déséquilibre de charge entre tâches
- IO séquentielles/re-séquentialisées ↔ Plate-forme trop faible
- Réseau d'interconnexion trop faible

Rmq : certains algorithmes nécessitent des plates-formes très performantes (d'autres non...)

Mesure, métriques et analyse de performances

3 – Métriques de Size Up

- Difficultés à supporter un size up
- Métrique de size up en temps d'exécution
- Métrique de size up en temps et ressources
- Exemple expérimental

Métriques de Size Up

Difficultés à supporter un Size Up

1^{er} objectif : pouvoir traiter des problèmes plus gros sur plus de rsracs

Un code qui **réplique** la plupart de ses données sur toutes les machines sera toujours limité par la taille mémoire d'une machine...
... et ne sera **pas apte au size up**

Métriques de Size Up

Difficultés à supporter un Size Up

1^{er} objectif : pouvoir traiter des problèmes plus gros sur plus de rsracs

Un code qui **répartit** la plupart de ses données sur toutes les machines pourra stocker plus de données sur plus de machines...
... et sera **apte au size up**

Métriques de *Size Up*

Difficultés à supporter un *Size Up*

1^{er} objectif : pouvoir traiter des problèmes plus gros sur plus de rsrc

→ Conception d'une répartition initiale des données, et d'un schéma de communication à volume minimal

- Besoin de faire circuler les données initiales entre les nœuds : pour qu'un nœud puisse poursuivre ses calculs sur d'autres données que les siennes
- Besoin de faire circuler les résultats intermédiaires entre les nœuds : pour qu'un nœud puisse poursuivre les calculs d'un autre, avec ses propres données

Métriques de *Size Up*

Métrique de *Size Up* en temps d'exec.

2^{ème} objectif : maintenir constant T_{exec}

$$T(1 \times n_1, p_1) = T(2 \times n_1, p_2) = T(k \times n_1, p_k) = C^{te}$$

avec : T (taille pb, nb rsrc)

→ Objectif pas toujours atteint !

Métriques de *Size Up*

Métrique de *Size Up* en temps et rsrc.

3^{ème} objectif : maintenir constant T_{exec} avec le nombre min de rsrc

$$T(1 \times n_1, p_1) = T(2 \times n_1, p_2) = T(k \times n_1, p_k) = C^{te}$$

avec : T (taille pb, nb rsrc)
avec : $O(p_k) = O(\text{calcul}(k \times n_1))$

→ Objectif pas toujours atteint !

Métriques de *Size Up*

Métrique de *Size Up* en temps et rsrc.

3^{ème} objectif : maintenir constant T_{exec} avec le nombre min de rsrc

$$T(1 \times n_1, p_1) = T(2 \times n_1, p_2) = T(k \times n_1, p_k) = C^{te}$$

avec : T (taille pb, nb rsrc)
avec : $O(p_k) = O(\text{calcul}(k \times n_1))$

→ Objectif pas toujours atteint !

Métriques de *Size Up*

Bilan : *Size Up* en 3 objectifs successifs

1^{er} objectif : pouvoir traiter des problèmes plus gros sur plus de rsrc

2^{ème} objectif : maintenir $T_{exec} = C^{te}$

3^{ème} objectif : $T_{exec} = C^{te}$ avec le nombre min de rsrc

Métriques de *Size Up*

Exemple expérimental

Co-simulation d'échanges thermiques au sein de buildings faiblement couplés (calculs quasi-indépendants)

Objectif : *size up* à temps constant avec nombre minimal de rsrc

$$T(1 \times 1 \text{ bldg}, p_1) = T(2 \times 1 \text{ bldg}, p_2) = T(k \times 1 \text{ bldg}, p_k) = C^{te}$$

avec : $O(p_k) = O(\text{calcul}(k \times 1 \text{ bldg})) \approx O(k)$

Version 1 : **hcp** de calculs et peu de comm. → **Size up dès 1Gb/s**

Version 2 : **peu** de calculs et peu de comm. → **Size up avec 10Gb/s**

Mesure, métriques et analyse de performances

4 – Critères de passage à l'échelle

- Démarche de *Size Up* et de *Speedup*
- Graphique de passage à l'échelle
- Pourquoi une métrique en T-exec ?

Critères de passage à l'échelle

Démarche de *Size Up* & de *Speedup*

Définition et critères d'un passage à l'échelle « simple » :

- Permettre un *Size Up* pour traiter de plus gros pb sur plus de rsrsc
- Et de manière économiquement supportable

→ *Size Up* avec temps d'exéc. constant et nbr minimal de rsrsc

Critères de passage à l'échelle

Démarche de *Size Up* & de *Speedup*

Définition et critères d'un passage à l'échelle « complet » :

- Permettre un *Size Up* pour traiter de plus gros pb sur plus de rsrsc
- Et de manière économiquement supportable
- Permettre un *Speedup* pour toutes tailles de pb
- Avec toujours le même profil de décroissance du temps d'exéc.

Critères de passage à l'échelle

Graphique de passage à l'échelle

Constitution et exploitation d'un graphique de passage à l'échelle :

- on mesure $T(n,p)$ pour différentes tailles de pb (n) et nbr de rsrsc (p)
- on trace les courbes $T(n,p)$ en échelle log

Cas idéal : pour chaque taille de pb on obtient une droite parallèle aux autres, simplement traduite vers la droite!

“ Pouvoir traiter des problèmes sans limite de taille dans le futur ”

Critères de passage à l'échelle

Graphique de passage à l'échelle

Constitution et exploitation d'un graphique de passage à l'échelle :

- on mesure $T(n,p)$ pour différentes tailles de pb (n) et nbr de rsrsc (p)
- on trace les courbes $T(n,p)$ en échelle log

Validation du passage à l'échelle : par comparaison des mesures aux courbes attendues (pentes et positions)

Critères de passage à l'échelle

Graphique de passage à l'échelle

Constitution et exploitation d'un graphique de passage à l'échelle :

- on mesure $T(n,p)$ pour différentes tailles de pb (n) et nbr de rsrsc (p)
- on trace les courbes $T(n,p)$ en échelle log

Utilisation en « abaque » : Pour une taille de pb donnée on peut identifier le nb de rsrsc à utiliser pour respecter un T-exec maximum

Une solution qui *passse à l'échelle complètement* permet de :

- répondre aux besoins de calculs
- nécessiter le minimum de rsrsc
- quantifier les rsrsc nécessaires
- planifier les dépenses associées

Critères de passage à l'échelle

Pourquoi une métrique en T-exec ?

Quand la taille du pb croît on ne peut plus faire d'exéc séquentielle :

- pas assez de RAM, pas assez de disque...
- exécution trop longue, rsrc non monopolisable...

→ on ne peut pas mesurer de *référence séquentielle* !
 → on ne peut pas calculer de *Speedup* !!
 → on doit bâtir une analyse seulement sur des mesures du T-exec

Mesure et analyse de performances

5 – Loi de Amdahl

- Motivations
- Définitions
- Impact sur le speedup
- Confrontation à la réalité

Performances

Loi de Amdahl : Motivations

Fraction séquentielle du programme ...
 ... quel est son impact sur les performances ?

Performances

Loi de Amdahl : Définitions

Hypothèses de départ :

$$T(1) = T_{seq} = T_s^a + T_p^a$$

↓ Temps de la partie *parallélisable* du programme (au sens de Amdahl)
 ↓ Temps de la partie *séquentielle* du programme (au sens de Amdahl)

$$T(P) = T_s^a + T_p^a / P + \text{Overhead}$$

Hyp : machine parallèle idéale

- Synchro sans surcoût !
- Msgs instantanés !
- ...

$$T(P) = T_s^a + T_p^a / P$$

Performances

Loi de Amdahl : Définitions

Définition de la fraction séquentielle au sens de Amdahl :

$$f_s^a = \frac{T_s^a}{T(1)} = \frac{T_s^a}{T_s^a + T_p^a} \in [0;1]$$

Et la fraction parallèle: $f_p^a = \frac{T_p^a}{T(1)} = \frac{T_p^a}{T_s^a + T_p^a} \in [0;1]$

D'où : $f_s^a + f_p^a = 1$

Réécriture du temps d'exécution:

$$T(P) = T_s^a + T_p^a / P$$

$$T(P) = f_s^a T(1) + (1 - f_s^a) T(1) / P$$

Performances

Loi de Amdahl : Définitions

Conséquence sur le speedup :

Par définition du Speedup : $S^a(P) = \frac{T(1)}{T(P)}$

Donc : $S^a(P) = \frac{T(1)}{f_s^a T(1) + (1 - f_s^a) T(1) / P}$

Soit : $S^a(P) = \frac{1}{f_s^a + \frac{1 - f_s^a}{P}}$ $\lim_{P \rightarrow \infty} S^a(P) = \frac{1}{f_s^a}$ $S^a(P) \leq \frac{1}{f_s^a}$

$$S^a(P) = P \cdot \frac{1}{1 + f_s^a \cdot (P - 1)}$$

$$S^a(P) = S^{ideal}(P) \cdot \frac{1}{1 + f_s^a \cdot (P - 1)}$$

Performances

Loi de Amdahl : Impact sur Speedup

Allure du speedup :

Exemple (a.n.) :

$$f_s^a = 1\% \Rightarrow \begin{cases} S^a(P) < 100 \\ S^a(100) = 50.3 \end{cases}$$

Petite fraction séquentielle → grosse limitation pour P grand!

Performances

Amdahl : Confrontation à la réalité

En réalité les overhead ne sont pas négligeables :

$$T(P) = T_s^a + T_p^a / P + \text{Overhead}$$

⇒ $S^{\text{réel}}(P) < S^a(P)$

Tempes de communication
Tempes de synchronisation
Tempes d'ordonnancement
...

Pourtant on n'observe pas de si mauvais résultats (en général) !

"On sait obtenir de bon speedup sur de grosses machines parallèles"

Les fractions séquentielles sont généralement faibles ?
Le modèle ne correspond pas toujours à la réalité ?

Loi de Gustafson

Mesure et analyse de performances

6 – Loi de Gustafson

- Motivations
- Définitions
- Impact sur le speedup

Performances

Loi de Gustafson : Motivation

1 – Impact de la fraction séquentielle (selon Amdahl) ?

Programme séquentiel → Programme parallèle

Fraction séquentielle

2 – Mais on observe qu'Amdahl n'est pas réaliste (trop pessimiste) :

Construire un modèle plus proche de la réalité ...

Performances

Loi de Gustafson : Définitions

Idee/Constataion de base :

En pratique : on ne veut pas diminuer le temps d'exécution T_0 , mais augmenter la quantité de travail W traitée en T_0

- On traite W_0 en séquentiel en : $T(1, W_0)$
- On traite W en parallèle en : $T(P, W)$

Et on choisit P et W tels que : $T(P, W) = T(1, W_0) = T_0$

Performances

Loi de Gustafson : Définitions

Temps que durerait le traitement séquentiel de W :

$$T(1, W) = T_s^g(P, W) + P \times T_p^g(P, W) - \text{Overhead}$$

Temps de la partie parallélisée sur P processeurs (au sens de Gustafson)

Temps de la partie séquentielle du programme (au sens de Gustafson)

Performances

Loi de Gustafson : Définitions

Hypothèse : partie séquentielle constante

$$T_s^g(P_1, W_1) = T_s^g(P_2, W_2) = T_s^g(P, W) = Cte$$

Vrai notamment pour des calculs scientifiques itératifs

Performances

Loi de Gustafson : Définitions

Introduction de la fraction séquentielle de Gustafson :

$$f_s^g / T_s^g = f_s^g \times T(P, W), \quad f_s^g \in [0, 1]$$

$T(P, W) = Cte$, par définition
 avec : $P = P_g(W)$

$T_s^g = Cte$, par hypothèse

$$f_s^g = Cte$$

Nouvelle expression du temps séquentiel (estimé) :

$$T(1, W) = T_s^g(P, W) + P \times T_p^g(P, W)$$

$$T(1, W) = f_s^g \times T(P, W) + P \times (1 - f_s^g) \times T(P, W)$$

Performances

Loi de Gustafson : Définitions

Speedup fonction de la fraction séquentielle :

$$S(P, W) = \frac{T(1, W)}{T(P, W)}$$

$$S^g(P, W) = \frac{f_s^g T(P, W) + P(1 - f_s^g) T(P, W)}{T(P, W)}$$

$$S^g(P, W) = f_s^g + P(1 - f_s^g)$$

$S^g(P, W) = f_s^g + P(1 - f_s^g)$ Speedup au sens de Gustafson

Mais la modélisation de Gustafson considère uniquement des couples (P_i, W_i) tels que $T(P_i, W_i) = T(1, W_0)$

$$S^g(P(W)) = f_s^g + P(W) \cdot (1 - f_s^g)$$

Performances

Loi de Gustafson : Impact sur speedup

Speedup non borné :

$$\lim_{P \rightarrow \infty} S^g(P) = \infty$$

Allure du speedup :

Bilan de la loi de Gustafson :

- Adopte le **point de vue utilisateur** : augmenter W à T-exec constant
- Hypothèse optimiste** de partie séquentielle constante ($T_s^g(P, W) = Cte$)

Mesure et analyse de performances

7 – Liens Amdahl-Gustafson

- Comparaison Amdahl – Gustafson
- Relation entre les fractions séquentielles
- Relation entre les courbes de speedup
- Généralisation des modélisations
- Bilan

Performances

Amdahl – Gustafson : Comparaison

Amdahl et Gustafson aboutissent à des conclusions différentes :

Amdahl

Etudie l'**extensibilité forte** (*strong scalability*) : capacité à décroître T-exec qd seule la taille de la machine croit.

Gustafson

Etudie l'**extensibilité faible** (*weak scalability*) : capacité à maintenir T-exec qd les tailles de la machine et du pb croissent.

Performances

Amdahl – Gustafson : Comparaison

Amdahl et Gustafson aboutissent à des conclusions différentes :

Amdahl

Gustafson

En fait Amdahl et Gustafson :

- Font des hypothèses différentes
- Définissent des fractions séquentielles différentes
- **Ne sont pas incompatibles**

Performances

Amdahl – Gustafson : Comparaison

Différence de modélisation de la fraction séquentielle :

Gustafson-Barsis:

T_s^g remains constant when w increases

Amdahl:

$T_s^a(w)$ evolution unspecified when w increases

Performances

Amdahl – Gustafson : Relation entre S

Expressions complètes des speedup

$$S^a(p, w) = \frac{1}{f_s^a(w) + \frac{1 - f_s^a(w)}{p}}$$

$$S^g(p_g(w)) = f_s^g + P_g(w) \cdot (1 - f_s^g)$$

Allures des speedup

En augmentant P et W (selon Gustafson) à f_s^g constant :

→ $S^g(P)$ intersekte différent $S^a(P)$, avec des $f_s^a(w)$ décroissant

Performances

Amdahl – Gustafson : Bilan

Amdahl : W Cte, P ↑, T-exec ↓
 Réaliste quand on travaille à W Cte ★ Informaticien
 Réaliste quand on augmente P pour diminuer T-exec

Gustafson : W ↑, P ↑, T-exec Cte
 Réaliste quand on augmente W sur *certaines* pbs ★ Utilisateur
 Ex : plus de cycles de calculs lors de simulations

En général :

Valeurs des courbes réelles moins bonnes à cause des *overhead*
 Allures des courbes réelles entre Amdahl et Gustafson

Mesure et analyse de performances

FIN