

Calcul – Cours 1: Introduction, Définition et Matériel

Jonathan Rouzaud-Cornabas

LIRIS / Insa de Lyon – Inria Beagle

Cours inspiré de ceux de Frédéric Desprez (DR Inria – Grenoble)

Définitions

- Parallèle VS Distribué
- Compute intensive VS data intensive VS Big Data
- High Performance Computing (HPC) VS High Throughput Computing (HTC)
- Durée de vie d'un code VS durée de vie d'une machine

- Parallel Programming – For Multicore and Cluster System (T. Rauber, G. Rünger)
- Sourcebook of Parallel Computing (J.J. Dongarra, I. Foster, G. Fox, W. Gropp, K. Kennedy, L. Torczon, A. White)
- Parallel Algorithms (H. Casanova, A. Legrand, Y. Robert)
- Parallel Computer Architecture (D.E. Culler, J. Pal Singh)
- Advanced Parallel Architecture - Parallelism, Scalability, Programmability (K. Hwang)
- Super-ordinateurs – Aux extrêmes du calcul (Numéro spécial de la recherche, Nov. 2011)
- Cours en ligne :
 - Why parallel, why now, Dr Clay Breshears, Intel
 - Architecture et Système des Calculateurs Parallèles, F. Pellegrini, LaBRI
 - Applications of Parallel Computers, J. Demmel, U.C. Berkeley CS267
 - K. Yellick
 - U.E. Parallélisme (Lyon 1), Frédéric Desprez

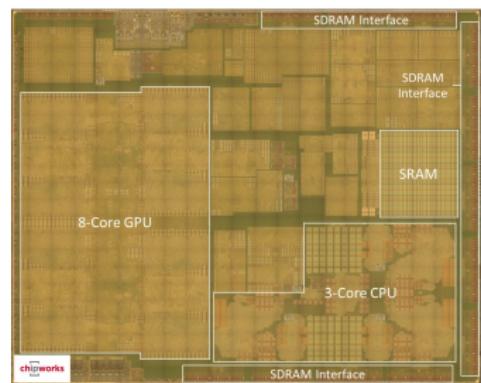
Pourquoi paralléliser et Loi de Moore

Pourquoi un cours sur le parallélisme ?

Le parallélisme est partout !



©RIKEN



Mais pourquoi en a-t-on besoin ?

- Résoudre des problèmes plus rapidement
 - Traiter plus de requêtes à la seconde (exemple Google)
 - Améliorer les temps de réponse des applications interactives
- Obtenir des meilleurs résultats dans le même temps
 - Rafiner les modèles (exemple Météo France)
 - Compliquer les modèles (multi-échelles)
- Traiter des problèmes de plus grande taille
 - Maillages plus importants
 - Données des réseaux de capteurs, des réseaux sociaux, ...
 - Recherche dans les pages web (Google)
 - LHC (CERN)

Pourquoi ne pas accélérer les processeurs ?

- Faisons un calcul rapide, pour avoir un processeur séquentiel à 1Tflop
 - Les données doivent aller de la mémoire au CPU (distance r)
 - Récupérer une donnée par cycle (10^{12} fois par seconde) à la vitesse de la lumière ($c = 299,792,458m/s \approx 3 \times 10^8 m/s$)
 - Donc $r < \frac{c}{10^{12}} = 0.3mm^2$
- Conséquence: Mettre 1 To de données dans $0.3mm^2$
 - Un mot occupe $\approx 3\text{Angstroms}^2$ (la taille d'un petit atome)
- Impossible à réaliser avec la technologie actuelle (gravure et interférence)
- Attention aussi à la chaleur que dégagerait un processeur comme ça !

Densité et puissance

- Les processeurs séquentiels puissants gachent de la puissance électrique
 - Spéculation, vérification de dépendances dynamique
 - Découverte de parallélisme
- Les systèmes concurrents sont plus efficaces d'un point de vue énergétique
 - L'augmentation de la fréquence (f) augmente le voltage (V)
 - L'augmentation du nombre de cores augmente la capacité (C) mais linéairement
 - La puissance dynamique est proportionnelle à $V^2 * f * C$
 - Tendance: Economiser de la puissance en baissant la fréquence

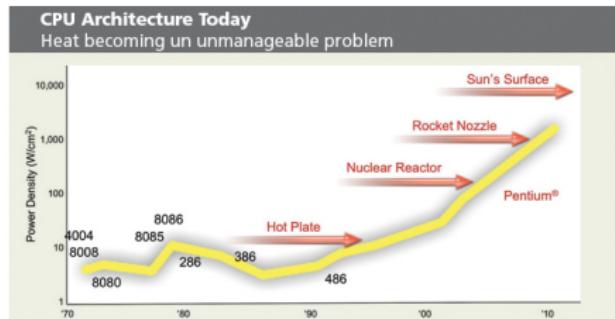
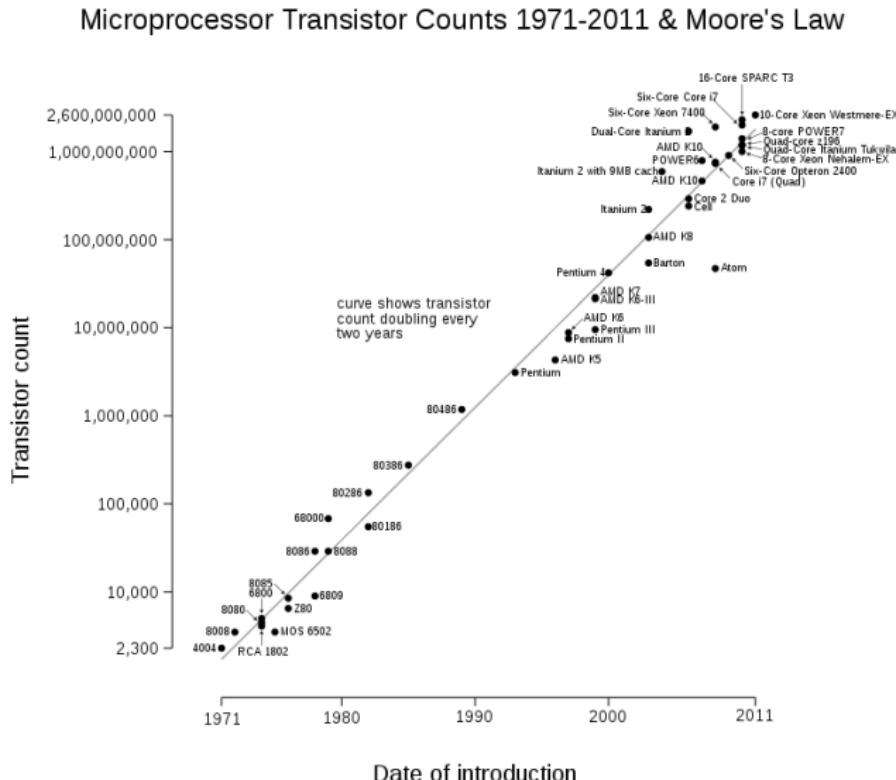


Figure 1. In CPU architecture today, heat is becoming an unmanageable problem.
(Courtesy of Pat Gelsinger, Intel Developer Forum, Spring 2004)

Loi de Moore

Gordon Moore (co-fondateur d'Intel) a prédit en 1965 que la densité des processeurs doublerait environ tous les 18 mois.



Loi de Moore

- En 1965, raisonnement empirique basé sur une relation entre la complexité des circuits et le temps
- Loi qui a perduré à travers les années
- Croissance due à plusieurs facteurs
 - Augmentation de la complexité des processeurs (densité en transistors, augmentation de la taille)
 - Ajout de fonctionnalités (caches internes, buffers d'instructions plus grands, lancement de plusieurs instructions par cycles, multithreading, profondeur des pipelines, ré-arrangement des instructions)

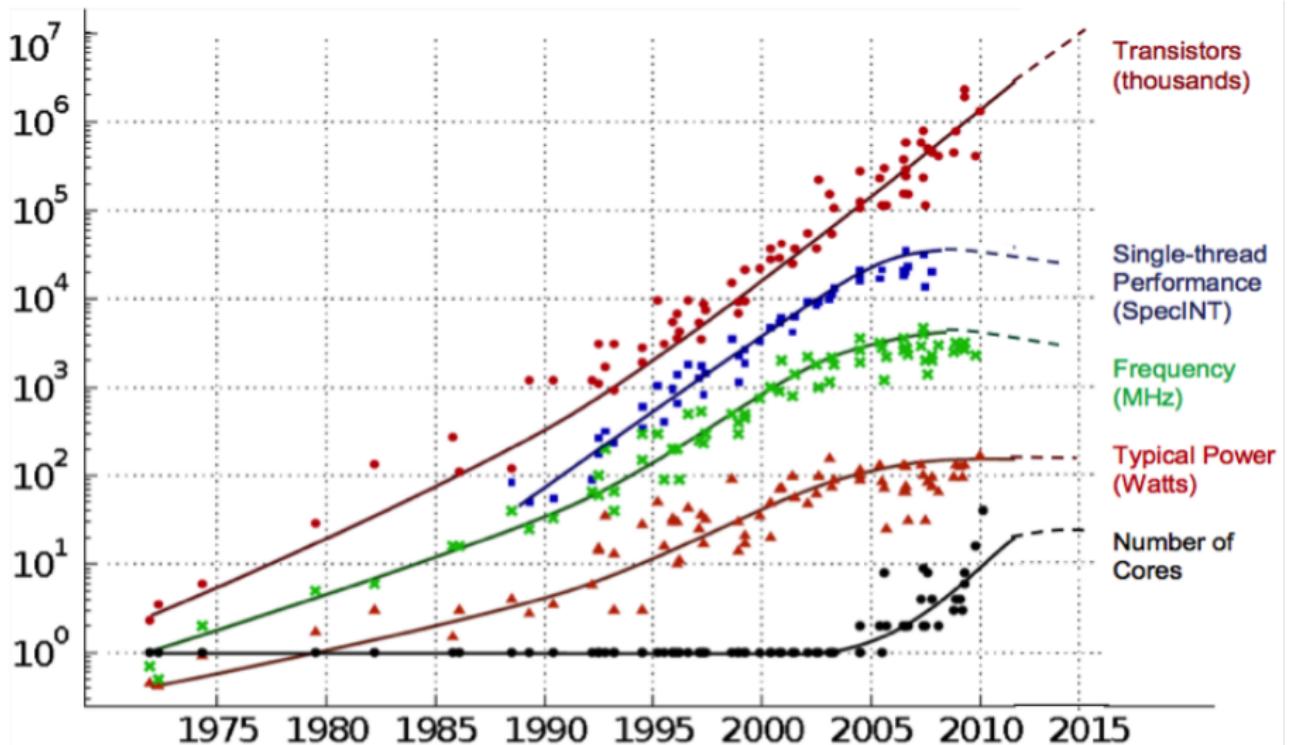
“The free lunch is over” – Herb Sutter

- La vitesse d'horloge ne va plus continuer à doubler ...
- ... mais on veut que les applications continuent à accélérer
- Quelques problèmes dues à l'augmentation de la vitesse d'horloge
 - Consommation électrique
 - Dissipation de la chaleur
 - Fuites
- Mais aussi
 - Limite physique due à la vitesse de la lumière (propagation des signaux)

The Free Lunch Is Over, A Fundamental Turn Toward Concurrency in Software, Herb Sutter,
Dr. Dobb's Journal, 30(3), March 2005.

<http://www.gotw.ca/publications/concurrency-ddj.htm>

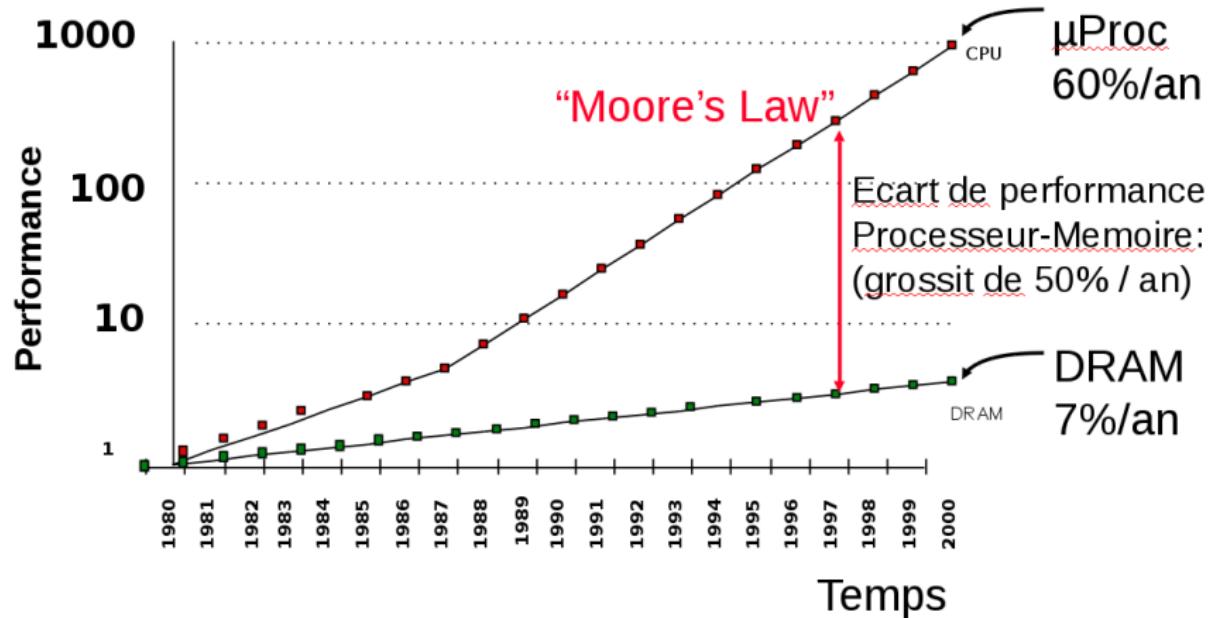
Loi de Moore : La réalité



Loi de Moore: Conséquence

- Le seul moyen d'augmenter les performances est l'augmentation du nombre d'unité de traitement travaillant en parallèle
- L'argument du coût des transferts des données tient toujours !
 - Il faut se débrouiller avec le rapport entre le grain de calcul et le grain de communication
 - Faire en sorte que le volume total des calculs dépasse de loin le volume des communications
- On va voir que c'est loin d'être simple !

Ecart entre le processeur et la mémoire



But: trouver des algorithmes qui minimisent les transferts de données, pas forcément les calculs

Problème de chargement des données

- La vitesse d'exécution des calculs ne dépend pas uniquement de la vitesse du processeur
- Il faut prendre en compte la vitesse de chargement des données de la mémoire vers le processeur
- La vitesse d'accès mémoire n'a augmenté que de 10% par an
 - Goulot d'étranglement important qui tend à augmenter
- Performance du système de mémoire dépendent de la fraction de la mémoire totale qui peut être stockée dans le cache
- Meilleures performances pour les unités parallèles car
 - Caches agrégés plus grands
 - Bande-passante agrégée plus importante
 - Attention à la localité des données (on verra ça plus tard)

Loi de Moore et Parallélisme: Conclusion

- La Loi de Moore tient toujours mais pas comme on s'y attendait
- Le parallélisme est obligatoire et présent partout
- Le chargement des données est aussi important que la rapidité de calcul
- Même problème sur les grands systèmes parallèles
 - Clusters
 - Super-calculateurs
 - Clouds
 - ...

Applications Parallèles

Quelques exemples d'applications

- Sans ordinateurs, soit
 - on étudie les phénomènes sur papier (théorie),
 - on construit un instrument et effectuer des expérimentations
- Limitations de la théorie et expérimentale
 - Trop compliqué (ex. : modéliser un Tsunami)
 - Trop cher (ex. : test de crash d'avion)
 - Trop lent (ex. : évolution climatique, des planètes, des organismes)
 - Trop dangereux (ex. : nucléaire, médicaments)
- Utiliser l'informatique pour simuler un phénomène
 - Basé sur les lois de la physique et l'utilisation de méthodes numériques
 - Expérience d'étude des lois de l'évolution
 - 30 ans d'expérience sur E. Coli → 60,000 générations
 - Tous les jours, avoir une personne pour manipuler
 - Simuler le même phénomène est faisable en quelques heures !

Problèmes complexes

- Intelligence Artificielle / Machine Learning
- Science
 - Comprendre la matière depuis les particules jusqu'à le cosmos
 - Prévoir la météo et le climat
 - Comprendre la bio-chimie des organismes vivants
- Ingénierie
 - Combustion et conception de moteurs
 - Modélisation de structures et des tremblements de terre
 - Nanotechnologies
- Entreprise
 - Finance numérique (HFT: marché très fort)
 - Recherche d'information
 - Exploration de données
- Défense
 - Essai nucléaire (uniquement numérique en France)
 - Cryptographie

Des besoins grandissants

- Croissance exponentielle de la puissance de calcul demandée (et acquise)
 - La simulation devient le troisième pilier de la science (à côté de la théorie et de l'expérimentation)
 - AI/ML utilisée partout et pour tout
- Croissance exponentielle des données expérimentales
 - Amélioration des techniques et technologies en acquisition, analyse, visualisation de données,
 - Aspects transports et stockage
 - Outils collaboratifs
 - Big data

Animation

- Le rendu utilisé pour appliquer des lumières, des textures et des ombres sur des modèles 3D afin de générer des images 2D pour un film
- Utilisation massive du calcul parallèle pour générer le bon nombre d'images pour un film complet (24 images/seconde)
- Quelques exemples
 - Pixar, 1995, Toy Story: premier film entièrement généré par ordinateur ("renderfarm" constituée de 100 machines dualprocs)
 - Pixar, 1999, Toy Story 2: utilisation d'un système à 1400 processeurs pour une meilleure qualité d'image
 - 2001, Monstres et compagnie: 250 serveurs à 14 procs pour un total de 3500 processeurs
 - Transformers 2, Industrial Light and Magic: render farm avec 5700 cores

Vidéo Massive Software

Biologie

- Augmentation importante des calculs avec l'arrivée massive de séquences d'ADN pour un grand nombre d'organismes, y compris des humains
- Celera corp. : Comparaison de génomes complets
 - Division du génome en petits segments, trouver les séquences ADN des segments expérimentalement, utiliser un calculateur pour construire la séquence complète en trouvant les zones de recouvrements
 - Nombres de comparaisons énorme
- Human Brain Project: Simuler le cerveau humain

Astrophysique

Explorer l'évolution des galaxies, processus thermonucléaires, traiter les données issues des télescopes

- Analyse d'ensembles de données très larges
 - Ensemble de données "Sky surveys"
 - Sloan Digital Sky Surveys, <http://www.sdss.org/>
 - Analyse de ces ensembles de données pour trouver des nouvelles planètes, comprendre l'évolution des galaxies

Vidéo Supernova

Modéliser le climat

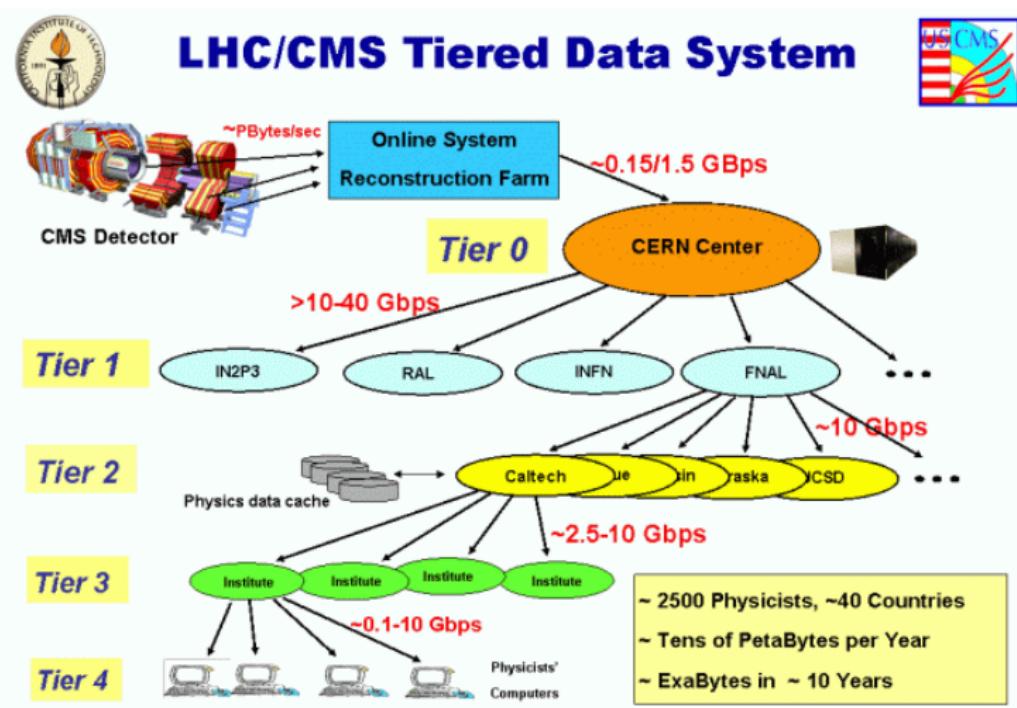
- Calculer (température, pression, humidité, vitesse du vent) = $f(\text{latitude}, \text{longitude}, \text{hauteur}, \text{temps})$
- Approche
 - Discréteriser le domaine (un point de mesure tous les 10 kms)
 - Exécuter un algorithme qui prédit le temps à $t + 1$ en fonction du temps à t
- Utilisation
 - La météo
 - Prédire les catastrophes naturelles
 - Evaluer les changements climatiques
 - Compétitions sportives

Vidéo Climat et Tsunami

LHC: Large Hadron Collider – Grand collisionneur de hadrons

- Les accélérateurs de particules sont la clef pour découvrir des particules massives ($E = mc^2$)
- En particulier, le LHC avait pour but de découvrir/prouver le boson de Higgs
- Le LHC est l'accélérateur de particule le plus puissant jamais construit
- Quelques données:
 - 40 millions de collisions par seconde
 - Après filtrage, 100 collisions intéressantes par seconde
 - Un megabyte de données numérisé à chaque collision (soit 0.1Gb/s)
 - 10^{10} collisions enregistrées par an (soit 10Pb/year)

LHC: Large Hadron Collider – Grand collisionneur de hadrons



IA et HPC : Exemple sur Summit

- Comprendre l'être humain : comprendre les causes génétiques d'Alzheimer et pouvoir prédire les risques
- Combattre le cancer : chercher des relations cachées entre les gènes, des marqueurs biologiques et l'environnement
- De plus en plus de données : Réussir à les intégrer / classifier *i.e.* besoin de bande passante
- Des modèles (*e.g.* taille des réseaux de neurones) de plus en plus complexe : nécessité de puissances de calcul

Data Science: Vers un quatrième paradigme

- Les données (scientifiques et autres) augmentent exponentiellement
- La capacité de générer des données a excédé notre capacité à les stocker et les analyser
- Les systèmes de calcul et les capteurs suivent la loi de Moore
- Les données de plusieurs Petabytes sont maintenant communs
 - Climat: quelques dizaines de Pb
 - Genome: Séquenceurs haute débit (5Gb/jour VS 1genome/10ans)
 - Physique nucléaire: Prévision LHC à $16Pb/an$
 - Astrophysique: LSST et SKA pourront produire des données $Pb/nuit$
 - Entreprise: Un ordre plus faible mais croissance rapide

Qu'est ce que le parallélisme ?

Qu'est-ce que le calcul parallèle ?

- Pouvoir accélérer une application en
 - ① Divisant cette application en sous-tâches
 - ② Exécuter ces sous-tâches en parallèles sur des unités différentes
- Pour réussir, il faut être capable de
 - ① Trouver le parallélisme dans l'application
 - ② Trouver le bon grain de calcul/échange de données
 - ③ Avoir des connaissances pour concevoir une solution efficace sur la machine cible

Focus du cours

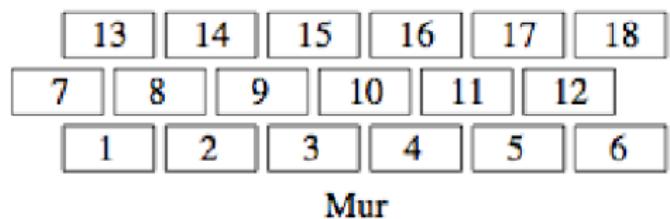
Le parallélisme est utilisé partout

- Recouvrement d'E/S sur un processeur
- Chargement et préparation d'instructions suivantes avec l'exécution d'instructions courantes dans un processeur
- Utilisation d'unités différentes (unités arithmétiques entière et flottante, unités flottantes multiples, unités pour le traitement graphiques,...)
- Multitasking, recouvrement de chargement de pages mémoire et calcul
- Multiprogrammation horizontale, processeurs VLIW (Very Long Instruction Word)

Dans ce cours, on va regarder le parallélisme en général (modèles, architecture, algorithmique) avec un focus plus important sur l'utilisation des CPUs et les GPUs (NVidia)

Aperçu du parallélisme

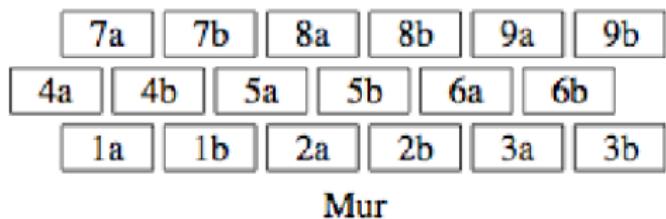
- Travail d'un maçon montant un mur de briques



- Seul il procède par rangées → C'est lent !

Aperçu du parallélisme

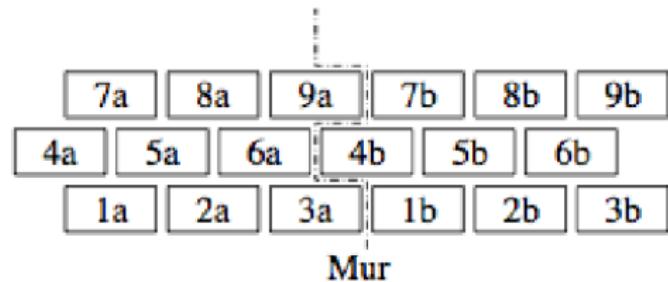
- Travail de 2 maçons (a et b) montant un mur de briques



- Une brique après l'autre → ils se gènent pour prendre des briques et les mettre en place !

Aperçu du parallélisme

- Travail d'un maçon montant un mur de briques



- Chacun s'attribue une portion de mur pour travailler
- Plus efficace mais
 - b a plus de chemin à faire pour récupérer les briques
 - ils se gènent pour prendre les briques
- Variantes
 - a travaille de droite à gauche et b commence plus vite mais problèmes de synchronisation

Aperçu du parallélisme

- Quelques réflexions sur l'exemple
 - Plus efficace que le travail d'un seul maçon mais
 - Plus de travail en général car organisation entre les maçons
- En général
 - Pour avoir une application parallèle, il faut que l'application soit décomposable en sous-problèmes suffisamment indépendants
 - Il faut pouvoir organiser le travail à répartir.
 - Surcoût dû à la répartition du travail (transmission des briques)
 - Trouver le meilleur algorithme parallèle → Pas forcément celui qui est le plus efficace en séquentiel !

Qu'espère-t-on ?

- Avoir un bon speed-up !
 - Idéalement, on espère avoir une accélération de p sur p processeurs !
- Malheureusement, c'est rarement le cas
 - Parties séquentielles d'un algorithme
 - Problèmes de surcoût (overhead) dus à des calculs redondants, les coûts de transfert de données (mémoire, disque, réseau)
- Parfois le gain est supérieur à p
 - Appelé speed-up superlinéaire
 - Grâce à des différences de vitesse mémoire (mémoire vive vs caches), moins de calcul grâce au parallélisme (recherche dans des arbres)
 - Applications pour lesquelles l'exécution sur un processeur est impossible (durée infinie d'exécution)

Somme préfixé

- Entrée : un opérateur \otimes et n éléments x_0, x_1, \dots, x_{n-1}
- Sortie : n éléments s_0, s_1, \dots, s_{n-1} tel que $s_i = x_0 \otimes x_1 \otimes \dots \otimes x_i$
- Exemple: Opérateur +

éléments	16	23	7	31	9
	16	16 +23	16 +23 +7	16 +23 +7 +31	16 +23 +7 +31 +9
prefix sums	16	39	46	77	86

Algorithme séquentielle

procedure PREFIX SUM(X, n)

$s_0 \leftarrow x_0$

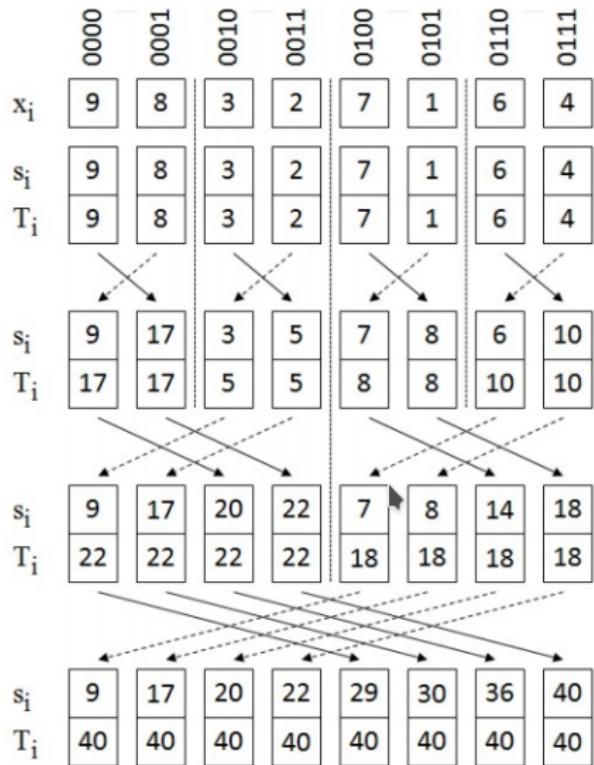
for $i \leftarrow 1$ to $n - 1$ **do**

$s_i \leftarrow s_{i-1} \otimes x_i$

return S

- Runtime $O(n)$
- Il y a une dépendance séquentielle i.e. pour calculer s_i il faut s_{i-1} .
Comment parallélise-t-on ?

- Nombre d'éléments n
- Nombre de processeurs p
- Considérons le cas suivant
 $n = p$
 - Élement sur P_i : x_i
 - Somme préfixé sur P_i : s_i
 - Somme totale sur P_i : T_i
- Temps de calcul $O(\log p)$
- Temps de communication $O(\log p)$



Algorithme parallèle

```
procedure PARALLEL PREFIX SUM( $id, X; d, p$ )
    prefix_sum  $\leftarrow X; d$ 
    total_sum  $\leftarrow prefix\_sum$ 
     $d \leftarrow \log_2 p$ 
    for  $i \leftarrow 1$  to  $d - 1$  do
        Envoyer total_sum au processeur  $id2 = id \otimes 2^i$ 
         $total\_sum + receivedtotal\_sum$ 
        if then  $id2 < id$ 
            prefix_sum  $\leftarrow total\_sum + receivedtotal\_sum$ 
    return prefix_sum
```

- Runtime $O(\log p) \neq \frac{\text{sequential runtime}}{p} = O(1)$

Solution générale

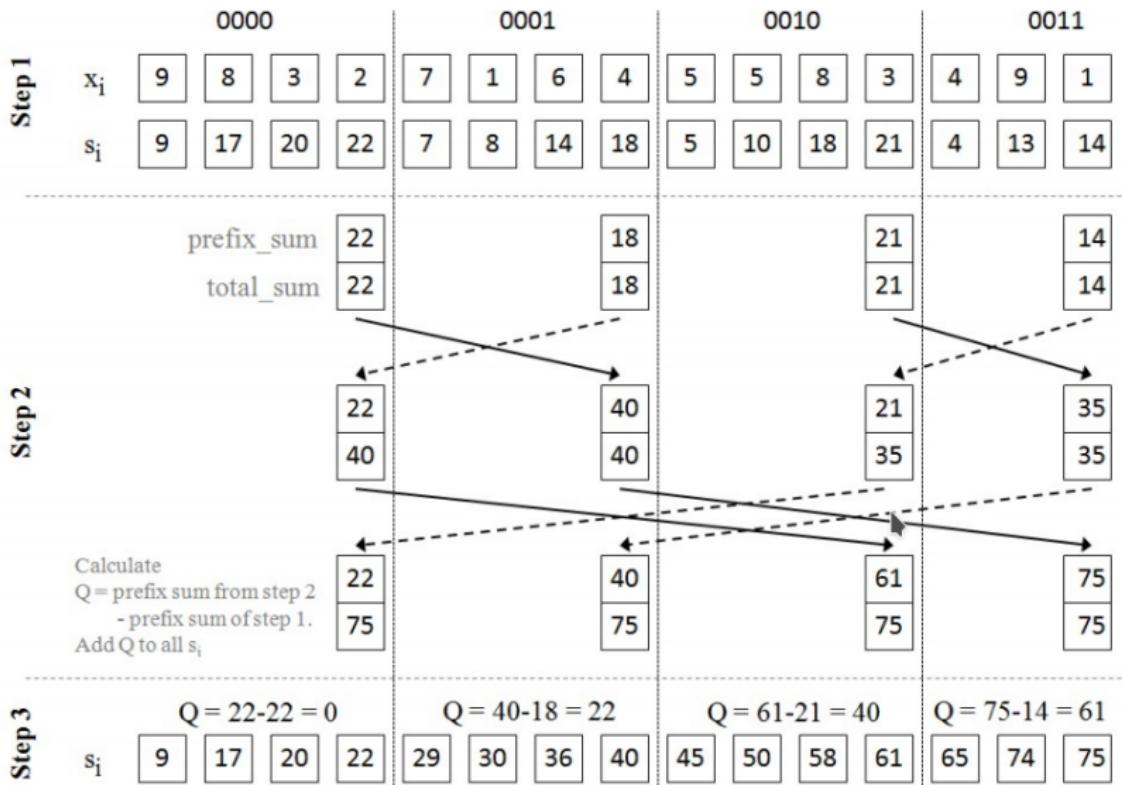
Cas réaliste:

- $n > p$
- n n'est pas un multiple de p
- p n'est pas une puissance de 2

Solution générale

Etape

- ① Chaque processeur calcule la somme préfixée pour $\frac{n}{p}$ éléments qu'il a localement
- ② En utilisant la dernière version calcul de la somme préfixé de chaque processeur, exécuter l'algorithme parallèle avec p -éléments
- ③ Sur chaque processeur, combiner le résultat de l'algorithme précédent avec la somme préfixée calculé localement (Etape 1)

Exemple ($n=15, p=4$)

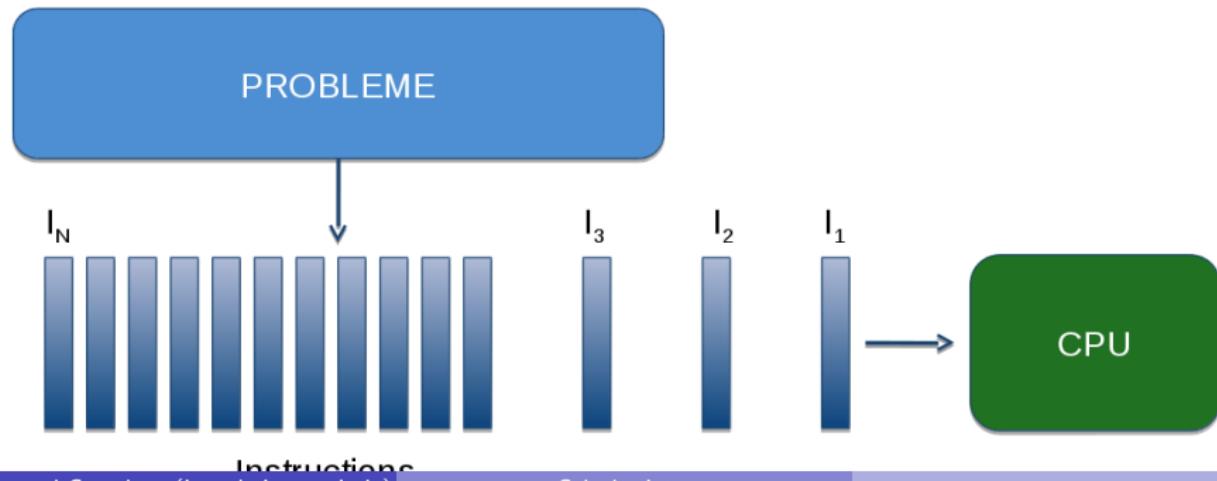
Complexité

- Etape 1 : Calcul des sommes préfixées localement sur $\frac{n}{p}$ éléments
 - Complexité de calcul : $O(\frac{n}{p})$
 - Complexité de communication : 0
- Etape 2 : Parallèle préfix utilisant la précédente somme préfixé sur chaque processeur
 - Complexité de calcul : $O(\log p)$
 - Complexité de communication : $O(\log p)$
- Etape 3 : Mettre à jour $\frac{n}{p}$ somme préfixée de l'étape 1 avec les résultats de l'étape 2
 - Complexité de calcul : $O(\frac{n}{p})$
 - Complexité de communication : 0
- Au totale
 - Complexité de calcul : $O(\frac{n}{p} + \log p)$
 - Complexité de communication : $O(\log p)$

En pratique sur une machine

Les programmes sont généralement conçus pour s'exécuter sur des processeurs séquentiels

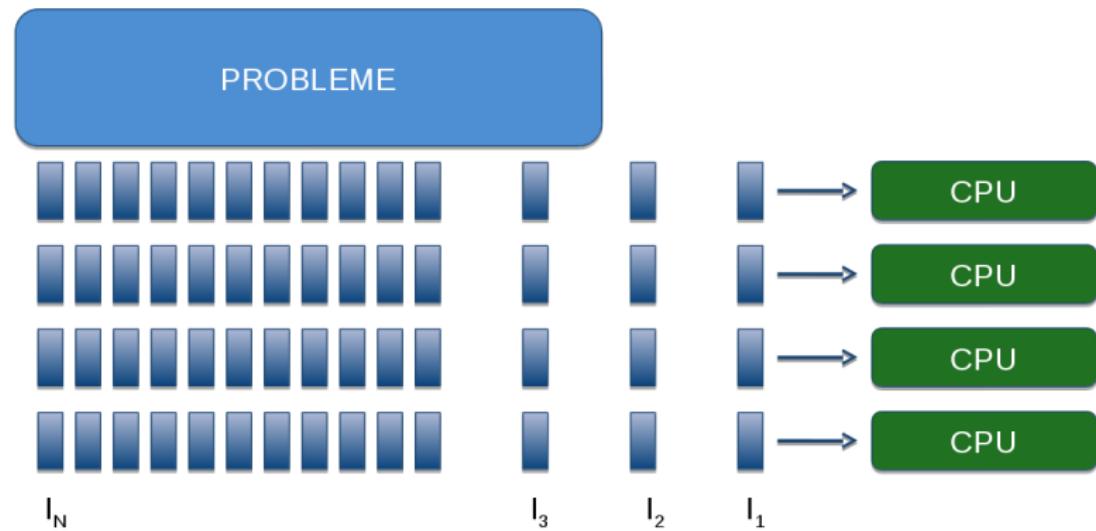
- Unité centrale (CPU) unique
- Application découpée en suite d'instructions exécutées l'une après l'autre
- Une seule instruction s'exécute à un instant donné



En pratique sur une machine

Sous la forme la plus simple, on utilise plusieurs ressources pour résoudre un problème

- Problème divisé en sous-parties indépendantes (si possible)
- Utilisant plusieurs CPU



Qu'est-ce qu'une machine parallèle?

- Une collection d'éléments de calcul capables de communiquer et de coopérer dans le but de résoudre rapidement des problèmes de grande taille
- Une collection d'éléments de calcul
 - Combien ?
 - De quelle puissance ?
 - Que sont-ils capables de réaliser?
 - Quelle est la taille de leur mémoire associée ?
 - Quelle est l'organisation ?
 - Comment les entrées/sorties sont-elles réalisées ?
- ... capables de communiquer ...
 - Comment sont-ils reliés entre eux ?
 - Que sont-ils capables d'échanger ?
 - Quel est leur protocole d'échange d'information ?

Qu'est-ce qu'une machine parallèle?

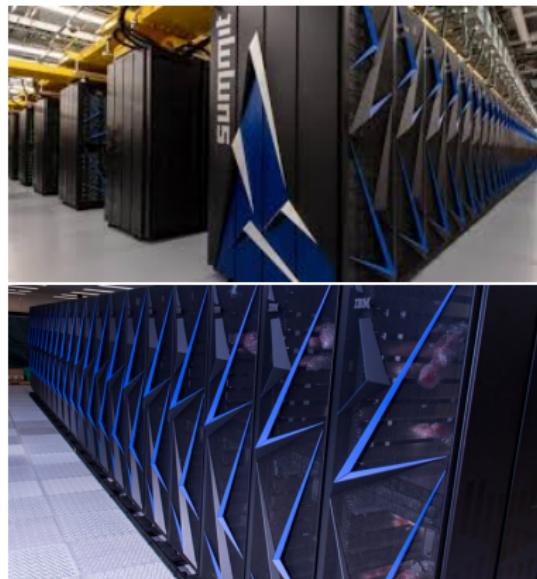
- ... et de coopérer ...
 - Comment les éléments de calcul synchronisent ils leurs efforts ?
 - Quel est leur degré d'autonomie ?
 - Comment sont-ils pris en compte par le système d'exploitation ?
- dans le but de résoudre des problèmes de grande taille.
 - Quels sont les problèmes à fort potentiel de parallélisme ?
 - Quel est le modèle de calcul utilisé ?
 - Quel est le degré de spécialisation des machines à un problème donné ?
 - Comment choisir les algorithmes ?
 - Quelle efficacité peut-on espérer ?
 - Comment ces machines se programment elles ?
 - Quels langages faut-il ?
 - Comment faut-il exprimer le parallélisme ?
 - Le parallélisme est il implicite ou explicite ?
 - L'extraction du parallélisme est-elle automatique ou manuelle ?

Unités de mesure des performances

- Les unités en HPC
 - Flop: floating point operation, généralement double précision
 - Flop/s: opérations flottantes par seconde
 - Octets: taille des données (8 pour un nombre en double précision)
- Des tailles typiques millions, milliards, trillions ...
 - Mega Mflop/s = 10^6 flop/sec
 - Giga Gflop/s = 10^9 flop/sec
 - Tera Tflop/s = 10^{12} flop/sec
 - Peta Pflop/s = 10^{15} flop/sec
 - Exa Eflop/s = 10^{18} flop/sec
 - Zetta Zflop/s = 10^{21} flop/sec
 - Yotta Yflop/s = 10^{24} flop/sec
- La machine la plus puissante au monde actuellement 187 Pflop/s et consomme 8.8MW (125 et 15.3MW pour la 2ème)
 - Liste mise-à-jour deux fois par an: www.top500.org

Summit (Sierra) – ORNL (LLNL) (US)

- 2,282,544 (1,572,480) cores
 - 9,216 (8,640) IBM POWER9 22C 3.07GHz
 - 27,648 (17,280) NVIDIA Volta GV100
- 122,300 (71,610) TFlop/s crête
- 187,659 (119,194) TFlop/s obtenus
- 2,801,664 (1,382,400) GB de mémoire
- 8,805 kW de conso électrique
- Linux comme OS



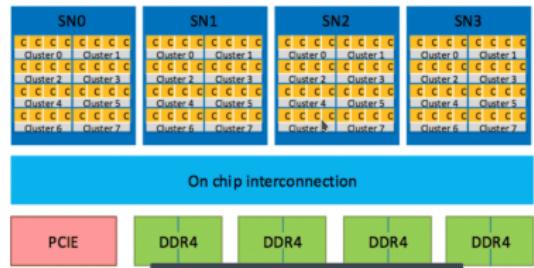
Tianhe-2 (MilkyWay-2) – National SuperComputer Center (CH)

- 3,120,000 cores
 - 32,000 Intel Xeon E5-2692 12C 2.200GHz
 - 48,000 Intel Xeon Phi 31S1P
- 61,444 TFlop/s crête
- 100,678 TFlop/s obtenus
- 1,024,000 GB de mémoire
- 18,482.00 kW de conso électrique
- Linux comme OS



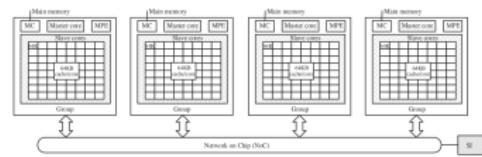
Tianhe-2A (MilkyWay-2A) – National SuperComputer Center (CH)

- 3,120,000 cores
 - 35,584 Intel Xeon E5-2692 12C 2.200GHz
 - 71,168 Matrix-2000 (260cores 1.2Ghz)
- ??? TFlop/s crête
- 94,97 TFlop/s obtenus
- 3,400,000 GB de mémoire
- 16,900.00 kW de conso électrique
- Linux comme OS



Sunway TaihuLight – National Supercomputing Center in Wuxi (CH)

- 10,649,600 cores
 - 40,960 Sunway SW26010 260C
1.45GHz
- 125,436 TFlop/s crête
- 93,014.6 TFlop/s obtenus
- 1,310,720 GB de mémoire
- 15,371.00 kW de conso électrique
- Linux comme OS



Piz Daint – Swiss National Supercomputing Centre (CSCS)

- 361,760 cores
 - Xeon E5-2690v3 12C 2.6GHz
 - NVIDIA P100
- 19,590 TFlop/s crête
- 25,326.3 TFlop/s obtenus
- 340,480 GB de mémoire
- 2,272.13 kW de conso électrique
- Linux comme OS



Mais avant le parallélisme, optimiser les codes !

- Quelques optimisations classiques
 - Préfetching d'instructions
 - Ré-ordonnancement d'instructions
 - Unités fonctionnelles pipelinées
 - Prédiction de branches
 - Allocation d'unités fonctionnelles
 - Hyperthreading
- Par contre, cela nécessite une complexification du matériel (unités fonctionnelles parallèles) et du logiciel (compilateurs, système d'exploitation) pour les supporter

Trouver le parallélisme

Comment trouver le parallélisme ?

- Partir d'un langage séquentiel ?
 - L'application a un parallélisme intrinsèque
 - Le langage de programmation choisi ne possède pas d'extension "parallèle"
- Le compilateur, le système d'exploitation et/ou le matériel doivent se débrouiller pour trouver le parallélisme caché!
 - Fonctionnement correct pour quelques applications trivialement parallèles (parallélisation de boucles imbriquées simples par exemple)
 - Mais en général, résultats décevants et pbs liés à la dynamicité (pointeurs en C par exemple)

Automatique dans les processeurs

- Parallélisme au niveau du bit (BLP, Bit Level Parallelism)
 - Dans les opérations flottantes
- Parallélisme d'instructions (ILP, Instruction Level Parallelism)
 - Exécuter plusieurs instructions par cycle d'horloge
 - Super-scalar, VLIW, EPIC, ThLP (Thread level parallelism: multithreading)
- Parallélisme des gestionnaires de mémoire
 - Recouvrir les accès mémoire avec le calcul (prefetch)
 - Opérations vectorielles en parallèle ($A[*] = 3 \times A[*]$)

Automatique dans les processeurs

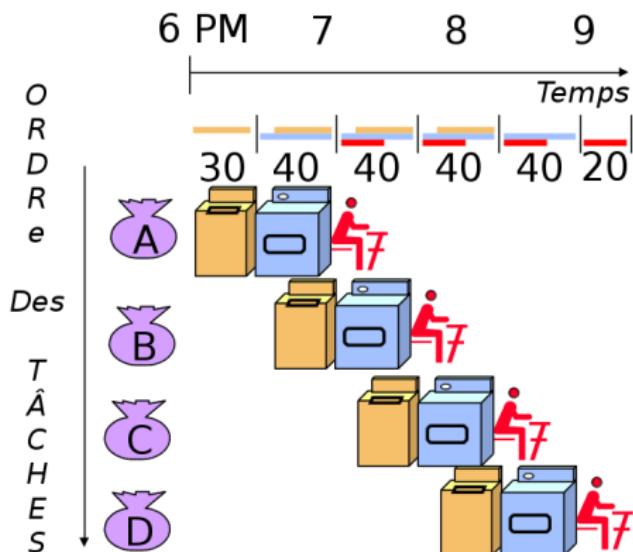
- Parallélisme au niveau du système
 - Exécuter des tâches différentes sur des processeurs (ou des cores) différents
 - fork [func1(), func2()], join [*]
- Limites à ce parallélisme “implicite”
 - Niveau d'intelligence des processeurs et des compilateurs
 - Complexité des applications
 - Nombre d'éléments en parallèle

Coopération

- Le programmeur et le compilateur travaillent ensemble
 - L'application a un parallélisme intrinsèque
 - Le langage possède des extensions permettant d'exprimer le parallélisme
 - Le compilateur va traduire le programme pour des unités multiples
- Le programmeur donne des conseils au compilateur sur les zones qu'il faut optimiser, quelles sont les boucles parallèles, ...
- Le compilateur, en partant des informations qu'il possède sur le matériel (taille des caches, nombre d'unités parallèles, informations sur les performances), va pouvoir générer un code performant.

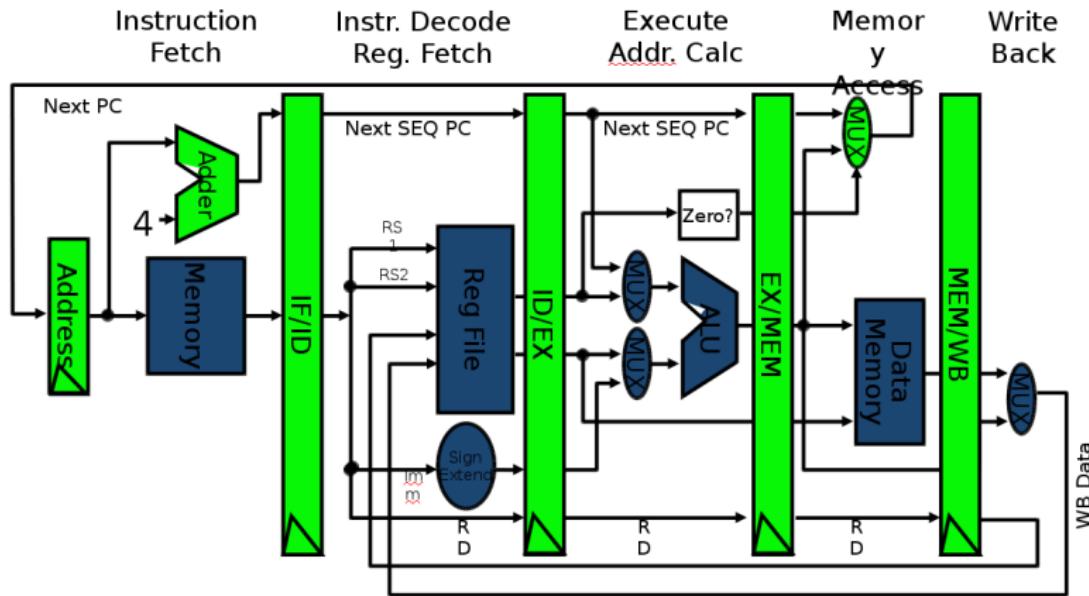
Pipelines

Exemple de Dave Patterson : 4 personnes qui lavent leur linge
 lavage (30min) + séchage (40min) + pliage (20min) = 90min



- Dans cet exemple:
 - exécution séquentielle = $4 \times 90\text{min} = 6\text{h}$
 - exécution pipelinée = $30 + 4 \times 40 + 20\text{min} = 3.5\text{h}$
- Bande passante = chargt/h
- Sans pipeline $4/6$
- Avec pipeline = $4/3.5$
- Le pipeline améliore la bande-passante mais pas la latence (90 min)
- La bande-passante est limitée par l'étage le plus lent
- Accélération potentielle = nombre d'étages du pipeline

Pipeline: Processeur



Pipeline utilisé par les unités arithmétiques

- Une multiplication flottante peut avoir une latence de 10 cycles et un débit de 1 cycle

SIMD: Single Instruction, Multiple Data

Calcul scalaire

- Mode "classique"
- 1 opération = 1 résultat

$$\begin{array}{c}
 \text{X} \\
 + \\
 \text{Y} \\
 \hline
 \text{X} + \text{Y}
 \end{array}$$

Calcul SIMD

- avec SSE (Streaming SIMD extensions) / SSE2
- 1 opération = n résultat

$$\begin{array}{c}
 \text{X} \quad \begin{matrix} \text{x3} & \text{x2} & \text{x1} & \text{x0} \end{matrix} \\
 + \\
 \text{Y} \quad \begin{matrix} \text{y3} & \text{y2} & \text{y1} & \text{y0} \end{matrix} \\
 \hline
 \text{X} + \text{Y} \quad \begin{matrix} \text{x3+y3} & \text{x2+y2} & \text{x1+y1} & \text{x0+y0} \end{matrix}
 \end{array}$$

SSE/SSE2 sur processeur Intel

- Types de données SSE2: tout ce qui rentre dans 16 octets, soit

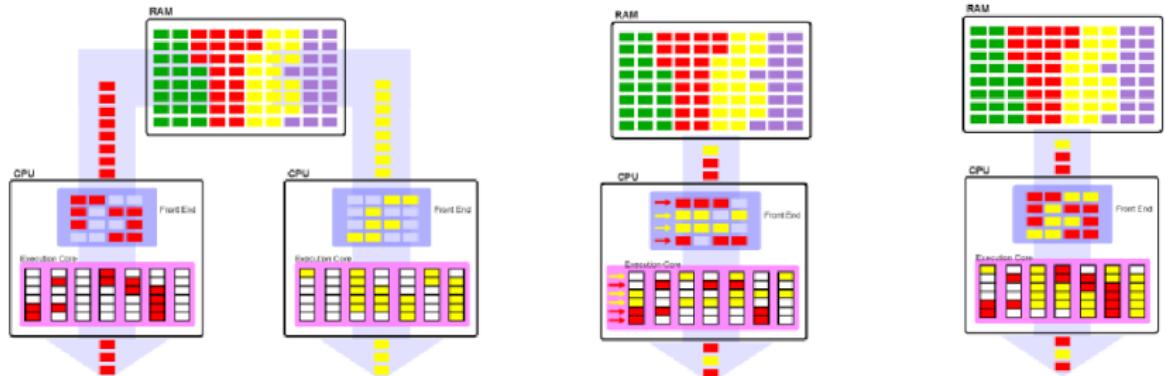


- Les instructions effectuent des additions, des multiplications, etc. sur toutes les données dans ces registres de 16 bits en parallèle
- Challenges
 - Doit être contigu en mémoire et aligné
 - Quelques instructions pour déplacer les données d'une partie d'un registre vers une autre
 - Similaire au GPU, processeurs vectoriels (mais plus d'opérations simultanées)

Instructions spéciales et compilateurs

- En plus des instructions SIMD, le processeur peut avoir aussi d'autres instructions
 - Instruction multiply-add (Fused Multiply-Add, FMA) : $x = y + x \times z$
 - Le processeur exécute ces instructions à la même fréquence qu'un \times ou un $+$
- En théorie les compilateurs connaissent ces instructions
 - Lors de la compilation, le compilateur va ré-arranger les instructions pour obtenir un bon ordonnancement des instructions qui va maximiser le pipeline (FMA et SIMD)
 - Il utilise des mélanges de telles instructions dans des boucles internes
- En pratique, le compilateur a besoin d'aide
 - Prise en compte de drapeaux de compilation
 - Ré-arranger le code pour qu'il trouve plus facilement les bons "pipelines"
 - Utiliser des fonctions spéciales
 - Ecrire en assembleur !

ILP: multithreading



SMP mono-threadé

(*Symmetric Multi-Processor*)

Chaque processeur est superscalaire

CPU super-threadé

CPU Hyper-threaded

=

Multithreading

Simultané (SMT)

Exemple: Produit de matrices optimisé en séquentiel

Motivation

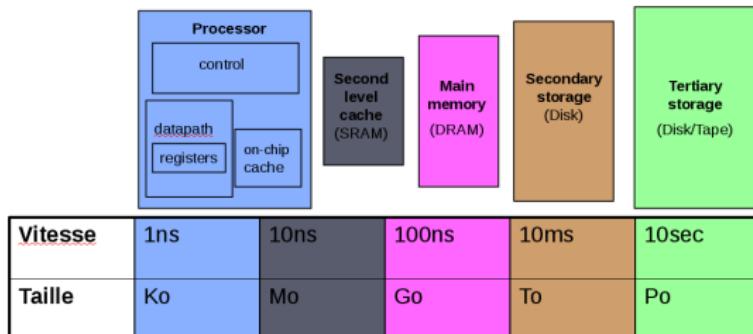
- La plupart des applications s'exécutent à une vitesse inférieure à 10% des performances en crête d'un processeur
- La plupart des pertes sont sur un processeur uniquement
 - Exécution du code avec des performances de 10 à 20% des performances en crête
 - Les autres pertes sont dues aux communications
- Les pertes de performances sont dues à la gestion des données en mémoire
 - Le déplacement des données coûte plus “cher” que les opérations arithmétiques et les branchements
- But de cette partie
 - Comprendre comment se comportent les algorithmes et pourquoi ils sont compliqués à implémenter

Processeurs séquentiels réels

- Les vrais processeurs possèdent
 - Des registres et des caches
 - Zones de données rapides et de petite taille
 - Stockent les données utilisées récemment ou proches
 - Différentes opérations mémoire peuvent avoir des coûts très différents
 - Parallélisme
 - Plusieurs unités fonctionnelles qui s'exécutent en parallèle
 - Ordres différents, mélanges d'instructions avec des coûts différents
 - Pipelines
- Pourquoi est-ce que c'est notre problème ?
 - En théorie, les compilateurs et le matériel comprennent l'architecture et son fonctionnement pour optimiser le programme
 - En pratique, c'est faux
 - Ils ne connaissent pas l'algorithme qui tirera un plus grand bénéfice du processeur

Hiérarchie mémoire

- La plupart des programmes possèdent une forte localité dans leurs accès aux données
 - **Localité spatiale:** accéder à des données proches des accès précédents
 - **Localité temporelle:** ré-utiliser une donnée qui a été accédée auparavant
- Les hiérarchies mémoire tentent d'exploiter la localité pour améliorer les performances (et on verra plus tard que sur les multi-cores c'est pire)



Approche pour supporter la latence mémoire

- La bande-passante a augmenté plus vite que la latence a décrue
 - 23% par an contre 7% par an
- Quelques techniques
 - Eliminer les opérations en sauvegardant les données dans des mémoires petites (mais rapides, les caches) et en les ré-utilisant
 - Ajouter de la localité temporelle dans les programmes
 - Utiliser mieux la bande-passante en récupérant un morceau de mémoire et en le sauvegardant dans un cache et en utilisant le bloc entier
 - Ajouter de la localité spatiale dans les programmes
 - Utiliser mieux la bande-passante en permettant au processeur d'effectuer plusieurs lectures en même temps
 - Concurrence dans le flot d'instructions (chargement d'un tableau entier dans les processeurs vectoriels, prefetching)
 - Recouvrement calcul et opérations mémoire
 - prefetching

Quelques rappels sur les caches

- **Cache:** mémoire rapide (et coûteuse) qui conserve des copies des données dans la mémoire principale (gestion cachée à l'utilisateur)
 - Exemple simple: une donnée à l'adresse xxxx1101 est stockée dans le cache à l'adresse 1101
- **Cache hit:** accès à la donnée dans le cache, peu coûteux
- **Cache miss:** accès à une donnée non cachée, coûteux
 - Besoin d'accéder au niveau supérieur (plus lent)
- Longueur d'une ligne de cache: nombre d'octets chargés en même temps

Pourquoi avoir des niveaux de caches multiples ?

- Le cache rapide ça coûte très cher !
- On-chip vs off-chip
 - Les caches internes sont plus rapides mais limités en taille
- Les caches de grandes tailles sont plus lents
 - Le matériel met plus de temps à vérifier les adresses plus longues
 - L'associativité, qui permet d'avoir des ensembles de données plus important, a un coût

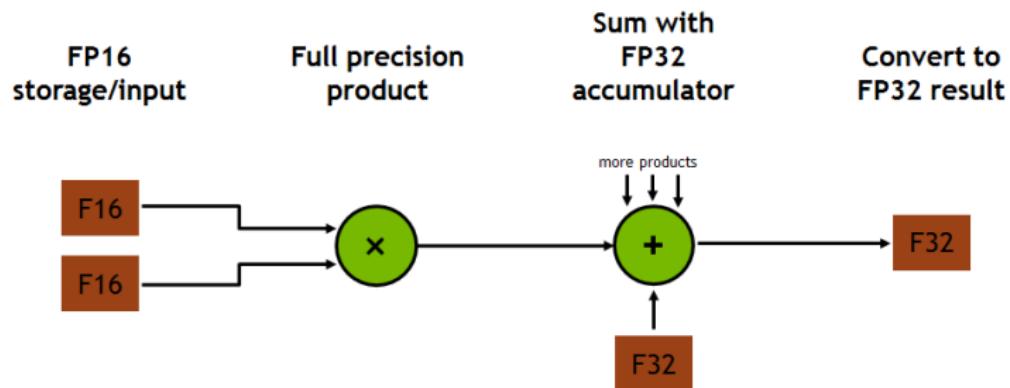
Conclusion sur les caches

- Les performances réelles d'un programme peuvent être difficiles à appréhender en fonction de l'architecture
 - La moindre modification de l'architecture et du programme a une grosse influence sur les performances
 - Pour écrire des programmes efficaces, il faut prendre en compte l'architecture
 - On voudrait des modèles simples pour concevoir des algorithmes efficaces
- Tenir compte des caches dans le programme
 - Utiliser un algorithme diviser-pour-régner pour que les données soient idéalement placées dans les caches L1 et L2
 - Utiliser une bibliothèque qui se charge de ça

Produit de matrices

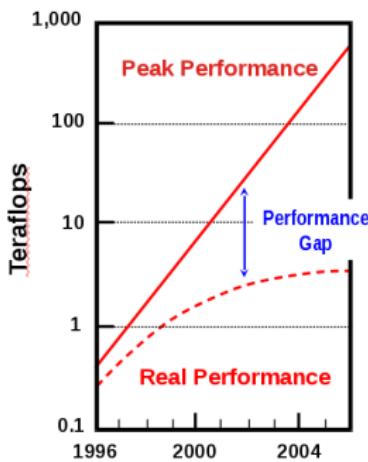
- Un noyau important de plusieurs applications numériques
 - Apparaît dans la plupart des algorithmes d'algèbre linéaire
 - Base du deep learning
 - Goulot d'étranglement de nombreuses applications
 - D'autres applications: algorithmes de graphes, réseaux de neurones, imagerie, ...
- Les optimisations peuvent être utilisées pour d'autres applications
- Assez facile à optimiser
- L'algorithme le plus étudié dans le monde du HPC

Tensor Core



Améliorer les performances réelles

- Les performances de crête augmentent exponentiellement
 - En 1990's, les performances de crête ont augmenté 100x; dans les années 2000, elles augmenteront 1000x (et on espère en 2010)
- Mais l'efficacité (les performances par rapport aux performances en crête) ont plutôt déclinées
 - 40-50% sur les supercalculateurs vectoriels des années 1990
 - Maintenant proches de 5-10% sur les supercalculateurs d'aujourd'hui
- Réduire l'écart..
 - Algorithmes qui obtiennent des performances sur un seul processeur et sont extensibles sur plusieurs milliers
 - Modèles de programmation plus efficaces et des outils pour les machines massivement parallèles



Le parallélisme de nos jours

- Tous les vendeurs de processeurs produisent des processeurs multicore
 - Toutes les machines seront bientôt parallèles
 - Pour continuer à doubler la puissance il faut doubler le parallélisme
- Quelles applications vont (bien) tirer partie du parallélisme ?
 - Est-ce qu'il faudra les redévelopper from scratch ?
 - Et les systèmes d'exploitation ?
- Est-ce que tous les programmeurs devront être des programmeurs de machines parallèles ?
 - Il faut des nouveaux modèles logiciels
 - Essayer de cacher le parallélisme au maximum
 - Mais il faut le comprendre !
- L'industrie parie sur ces changements...
- ... mais encore beaucoup de travail à effectuer

Challenges

- Les applications parallèles sont souvent très sophistiquées
 - Algorithmes adaptatifs qui nécessitent un équilibrage dynamique
- Le parallélisme multi-niveaux est difficile à gérer
- La taille des nouvelles machines donnent des problèmes d'efficacité
 - Problèmes d'extensibilité
 - Sérialisation et déséquilibrage de charge
 - Goulots d'étranglement et en communication et/ou en entrées/sorties
 - Parallélisation insuffisante ou inefficace
 - Fautes et/ou pannes
 - Gestion de l'énergie
- Difficulté d'obtenir les performances les meilleures sur les nœuds eux-mêmes
 - Contention pour la mémoire partagée
 - Utilisation de la hiérarchie mémoire sur les processeurs multi-cores
 - Influence du système d'exploitation

Conclusion

- L'ensemble des machines parallèles est constituée d'un ensemble (très) large d'éléments
 - depuis des unités parallèles dans les processeurs
 - Jusqu'à des datacenters connectés à travers le monde
- Champ d'étude du parallélisme
 - Architectures
 - Algorithmes
 - Logiciels, compilateurs,
 - Bibliothèques
 - Environnements
- Changements historiques importants
 - Jusque dans les années 90, réservées à des gros calculs de simulation
 - Aujourd'hui, parallélisme dans tous les processeurs (des téléphones aux supercalculateurs), parallélisme dans la vie courante (iPad, Google !)
 - Sans parallélisme, pas de Deep Learning