



UNIVERSITÉ DE REIMS CHAMPAGNE-ARDENNE

École Doctorale Sciences Technologie Santé

Rapport français

Pour obtenir le grade de:

Docteur de l'Université de Reims Champagne-Ardenne

Discipline : Informatique

Spécialité : Calcul Haute Performance

Présentée et soutenue par:

Julien LOISEAU

le 18 Mars 2018

Les Architectures Hybrides pour Atteindre l'Exascale

Sous la direction de :

Michaël KRAJECKI, Professeur des Universités

JURY

Pr. Michaël Krajecki	CRéSTIC, Université de Reims Champagne-Ardenne	Directeur
Dr. François Alin	CRéSTIC, Université de Reims Champagne-Ardenne	Co-directeur
Dr. Christophe Jaillet	CRéSTIC, Université de Reims Champagne-Ardenne	Co-directeur
Pr. William Jalby	Université de Versailles Saint-Quentin	Rapporteur
Dr. Benjamin Bergen	Ingénieur au Los Alamos National Laboratory, USA	Rapporteur
Pr. Zineb Habbas	LCOMS, Université de Lorraine	Rapporteur
Pr. Françoise Baude	I3S, Université de Nice Sophia Antipolis	Examineur
Guillaume Colin de Verdier	Ingénieur au CEA	Invité

Ce document comprend un résumé substantiel de ma thèse rédigée en anglais. Il présente tout d'abord la traduction de la table des matières de la thèse. Une traduction de l'introduction générale, une synthèse des trois parties de la thèse ainsi que de la conclusion générale.

1 Table des matières traduite en français

Cette traduction a été réduite aux parties, chapitres et sections de la thèse.

Introduction

I Exascale et calcul haute performance

1 Modèles pour le HPC

- 1.1 Introduction
- 1.2 Modèle de Von Neumann
- 1.3 Taxonomie de Flynn et modèles d'exécution
- 1.4 Mémoires
- 1.5 Caractérisation des performances dans le HPC
- 1.6 Conclusion

2 Hardware dans le HPC

- 2.1 Introduction
- 2.2 Évolution du modèle de Von Neumann
- 2.3 Architectures du 21ème siècle
- 2.4 Architectures distribuées architectures
- 2.5 Le supercalculateur ROMEO .
- 2.6 Conclusion

3 Software dans le HPC

- 3.1 Introduction
- 3.2 Modèles parallèles et distribués
- 3.3 Software/API
- 3.4 Benchmarks
- 3.5 Conclusion

II Métrique de problème complexes

4 Calcul : Le problème de Langford

- 4.1 Introduction
- 4.2 Algorithme de Miller
- 4.3 Méthode algébrique de Godfrey
- 4.4 Conclusion

5 Communication : Graph500

- 5.1 Introduction
- 5.2 Méthodes existantes
- 5.3 Environnement
- 5.4 Parcours en largeur
- 5.5 Résultats
- 5.6 Conclusions

III Application

6 Modélisation et systèmes complexes

- 6.1 Introduction
- 6.2 Modélisation physique et limitations
- 6.3 Cas appliqués
- 6.4 Conclusion

7 Simulations complexes sur architecture hybride

- 7.1 Introduction
- 7.2 FleCSI
- 7.3 SPH distribué sur architecture multi-core
- 7.4 SPH distribué sur architecture hybride
- 7.5 Résultats
- 7.6 Conclusion

Conclusion

Bibliographie

2 Introduction

Le monde du calcul haute performance (HPC) va prochainement atteindre une puissance de calcul inégalée avec l'Exascale. Les États Unis D'Amérique et l'Europe devraient l'atteindre aux horizons 2020-2021 mais la Chine pourrait proposer une telle machine dès 2019. Ces superordinateurs seront 100 fois plus rapides que l'estimation actuelle des performances du cerveau humain avec 10^{16} calculs à virgules flottante par secondes (FLOPS)[?, kurzweil2010singularity]t atteindre une puissance sans précédent d'un milliard de milliard (10^{18}) de FLOPS. Cette aventure a commencée avec les premiers ordinateurs à tube à vides et les besoins de la balistique pour la guerre. De nos jours les supercalculateurs étendent leur domaine d'application et sont un élément phare pour représenter la puissance d'une nation. Les applications se sont répandus dans tous les secteurs de la science et de la technologie.

Depuis 1962, et en considérant le Cray CDC 6600 tel que le premier superordinateur, la puissance de calcul des ces machines n'a fait qu'augmenter tout en suivant une observation faite par le co-fondateur de l'entreprise Intel, Gordon Moore. Mieux connue sous la nom de "Loi de Moore", cette loi de 1965 explique que, considérant l'évolution constante de la technologie, le nombre de transistors sur un circuit intégré va doubler environs tous les deux ans pour un prix similaire. Ceci va permettre une augmentation parallèle de la puissance de calcul que peux fournir les ordinateurs et superordinateurs. Plus important encore, comme "*L'argent est le nerf de la guerre*", le prix des puces pour des performances meilleures va diminuer.

Cette observation de Gordon Moore peut-être observée avec l'évolution de la puissance des supercalculeur, la liste TOP500¹. Présenté sur la figure 1, le loi de Moore se poursuit et reste vraie malgré des décennies.

La diminution de la taille des semi-conducteurs avec des transistors de plus en plus petits n'est pas la seule raison de cette évolution linéaire. Le premier processeurs étaient batis autour d'un unique coeur de calcul (CPU) avec une augmentation du nombre de transistors et une meilleure fréquence d'opération. Ils ont vite rencontré une limitation pour atteindre des fréquences toujours plus élevées du fait de l'énergie requise mais aussi de la chaleur générée devant être canalisée. C'est pourquoi, au début du vingtième siècle, IBM proposa le premier processeur multi-core, le Power4. Les constructeurs ont commencés à créer des puces avec plus de un coeur pour augmenter la puissance de calcul tout en continuant à réduire la taille des composants. Cela permis de répondre à la constante demande en puissance de calcul et la Loi de Moore fut conservée. Bien entendu cette nouvelle architecture avait plusieurs problèmes. Un coût supplémentaire en synchronisation entre les coeurs pour l'accès à la mémoire, le partage des tâches mais aussi une complexité des algorithmes parallèles. Des nos jours un CPU classique propose de deux à des dizaines de coeurs de calcul sur une seule carte.

In order to reach even more computational power some researchers started to use many-core approaches. By using hundreds of cores, these devices take advantage of very "simple" computing units, with slow frequency and low power consumption but add more complexity and requirement for their efficient programming with even more synchronizations needed between the cores. Typically, those many-core architectures are used coupled with a CPU that sends the data and drives them. Some accelerators like the Intel Xeon Phi can be driven or driver depending on their configuration. Usually called accelerators,

¹<https://www.top500.org>

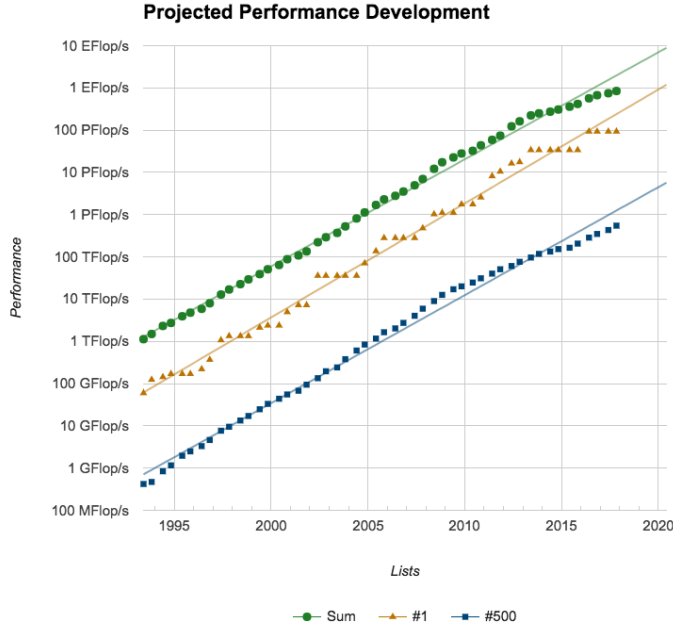


Figure 1: Évolution de la puissance de calcul, liste du TOP500

those devices are used in addition to the host processor to provide their efficient computational power in the key part of execution. The most famous accelerators are the Xeon Phi, the General-Purpose Graphics Processing Unit (GPGPU), initially used for graphic processing, Field Programmable Gates Array (FPGA) or dedicated Application-Specific Integrated Circuit (ASIC). The model using a host with additional device(s) appears and we will be referred to as "Hybrid Architecture". In fact, a cluster can be composed of CPUs, CPUs with accelerator(s) of the same kind, CPUs with heterogeneous accelerators or even accelerators like Xeon Phi driving different kinds of accelerators.

Since either 2013 or 2014 many companies, like the Gordon Moore's company *Intel* itself, stated that the Moore's law is over. This can be seen on figure 1: on the right part of the graph, the evolution is no longer linear and tends to decrease slowly in time. This can be contributed to two main factors. First, we slowly reach the maximal shrink size of the transistors implying hard to handle side effects. Second, the power wall implied by the power consumption required by so many transistors and frequency speed on the chip.

Even with all these devices, current supercomputers face several problems in their conception and utilization. The three main walls bounding the overall computational power of the machines are: the power consumption wall, the communication wall and the memory wall. Sub-problems like the interconnect wall, resilience wall or even the complexity wall also arise and make the task even more difficult.

In this context of doubts and questions about the future of HPC, this study proposes several points of view. We believe the future of HPC is made with these hybrid architectures or acceleration devices adapted to the need, using well suited API, framework and code. We consider that the classical benchmarks, like the TOP500, are not enough to target the main walls of those architectures, especially accelerators. Domain scien-

tists' applications like physics/astrophysics/chemist/biologist require benchmarks based on more irregular cases with heavy computation, communications and memory accesses.

In this document, we propose a metric that extracts the three main issues of HPC and apply them to accelerated architectures to determine how to take advantage of these architectures and what are the limitations for them. The first step of this metrics is obtained when merging two characteristic problems and then a third problem, combining all our knowledge. The first two are targeting computation and communication wall over very irregular cases with high memory accesses, using an academic combinatorial problem and the Graph500 benchmark. The last is a computational scientific problem that will cover both difficulties of the previous problems and appears to be hard to implement on supercomputers and even more on accelerated ones. The results obtain supports our thesis and show hybrid architectures as the best solution to reach Exascale.

Cette thèse se décompose en trois parties.

The first part explores the state of the art in HPC from the main laws to the hardware. We go through the basic laws from Amdahl's to Gustafson's laws and the specification of speedup and efficiency. Classical CPUs, GPGPUs and other accelerators are described and discussed regarding the state of the art. The main methods of ranking and the issues regarding them are presented.

In the second part we propose our metric based on characteristic problems to target classical and hybrid architectures. The Langford problem is described as an irregular and computationally heavy problem. This demonstrates how the accelerators, in this case GPUs, are able to support the memory and computation wall. This work leads to the publication of one journal paper [KLAJ16b] and many conferences, presentations and posters [DJK⁺14, LJAK16, JDK⁺14]. Our implementation of the Langford problem allowed us to beat a world record with the last instances of this academic problem.

The Graph500 problem is then proposed as an irregular and communications heavy problem. We present our implementation, and moreover, the logic to take advantage of the GPUs computational power for characteristic applications. This work led to the publication of a conference paper [KLAJ16a] and many presentations and posters [LAJK15, LJAK15].

In the third part, we consider a problem that is substantial and irregular in regards to computation and communications. We analyze this problem and show that it combines all the previous limitations. Then we apply our methodology and show how modern supercomputers can overcome these issues. This computational science problem is based on the Smoothed Particle Hydrodynamics method. The former application began with the development of the FleCSI framework from the Los Alamos National Laboratory which allowed us to exchange directly with the LANL domain scientists on their needs. We intend to provide an efficient tool for physicists and astrophysicists, called FleCSPH, based on our global work to efficiently parallelize these types of production applications. This work led to several presentation and posters [DBGH⁺16, LLMB17].

The last part will summarize on this work and results to show some of the main prospects of this study and my future researches.

3	Modèles du HPC
4	Métrique de problème complexes
5	Application
6	Conclusion
	Bibliographie

Bibliography

- [DBGH⁺16] Nicola De Brye, Daniel George, Glenn Hordemann, Hyun Lim, Julien Loiseau, Jonah Miller, and Johnathan Sharman. Domain partitioning and problem space representations for compact binary mergers. In *Poster at the 2016 Co-Design Summer School, Los Alamos.*, 2016.
- [DJK⁺14] Hervé Deleau, Christophe Jaillet, Michaël Krajecki, Julien Loiseau, Luiz Angelo Steffenel, François Alin, and Lycée Franklin Roosevelt. Towards the parallel resolution of the langford problem on a cluster of gpu devices. In *CSC14: The Sixth SIAM Workshop on Combinatorial Scientific Computing*, page 66, 2014.
- [JDK⁺14] Christophe Jaillet, Hervé Deleau, Michaël Krajecki, Luiz-Angelo Steffenel, François Alin, and Julien Loiseau. Langford problem: Massively parallel resolution on a multigpu cluster. In *CSC14: The Sixth SIAM Workshop on Combinatorial Scientific Computing*, 2014.
- [KLAJ16a] Michaël Krajecki, Julien Loiseau, François Alin, and Christophe Jaillet. Bfs traversal on multi-gpu cluster. In *2016 IEEE Intl Conference on Computational Science and Engineering (CSE) and IEEE Intl Conference on Embedded and Ubiquitous Computing (EUC) and 15th Intl Symposium on Distributed Computing and Applications for Business Engineering (DCABES)*, 2016.
- [KLAJ16b] Michaël Krajecki, Julien Loiseau, François Alin, and Christophe Jaillet. Many-core approaches to combinatorial problems: case of the langford problem. *Supercomputing Frontiers and Innovations*, 3(2):21–37, 2016.
- [LAJK15] Julien Loiseau, François Alin, Christophe Jaillet, and Michaël Krajecki. Parcours de grands graphes sur architecture hybride cpu/gpu. In *Modèles et Analyses Réseau: Approches Mathématiques et Informatiques (MARAMI) 2015. 14-16 Octobre Nîmes. Journée Fouille de Grands Graphes (JFGG)*, 2015.
- [LJAK15] Julien Loiseau, Christophe Jaillet, François Alin, and Michaël Krajecki. Massively parallel resolution of combinatorial problems on multigpu clusters. In *GTC Technology Conference 2015.*, 2015.
- [LJAK16] Julien Loiseau, Christophe Jaillet, François Alin, and Michaël Krajecki. Résolution du problème de langford sur architecture massivement parallèle cpu/gpu. In *Compas' Conference 2016. Lorient, June 2016, France*, 2016.

- [LLMB17] Julien Loiseau, Hyun Lim, Nick Moss, and Ben Bergen. A parallel and distributed smoothed particle hydrodynamics implementation based on the flecsu framework. In *SuperComputing conference, Denver CO*, 2017.