

Analyse de performance et optimisation de code

Pierre AYOUB

5 février 2019



**INSTITUT DES SCIENCES ET
TECHNIQUES DES YVELINES**

Résumé

La simulation numérique est un procédé informatique visant à modéliser un phénomène par ordinateur, s'agissant le plus souvent d'un phénomène physique. Cette modélisation prend forme par des systèmes d'équations décrivant l'état du système physique représenté à chaque instant. De nombreux domaines scientifiques convergent vers la simulation informatique, tel que certaines branches de la physique, de l'analyse et de l'optimisation mathématique, ou encore le calcul haute performance en informatique. Enfin, la simulation trouve naturellement de nombreuses applications concernant des sujets variés, tel que la simulation du climat et des événements météorologiques, la simulation d'essais nucléaires, de l'effet d'un médicament sur un corps, ou encore des astres et de l'univers. Ce rapport s'articulera donc autour de la simulation de fumée, phénomène impliquant les lois de la mécanique des fluides. Notre travail portera sur l'aspect du calcul haute performance de cette simulation.

Table des matières

1	Introduction	4
2	Analyse du code	5
3	Protocole expérimental	7
3.1	Théorie	7
3.2	Pratique	7
4	Optimisations et mesures	9
4.1	Mise en place de la compilation	10
4.2	Inlining	10
4.3	Déroulage de boucle	10
4.4	Vectorisation	10
5	Conclusion	13

1 Introduction

Le projet que nous vous présentons aujourd'hui consiste à analyser puis, grâce à nos mesures, optimiser un code de simulation numérique. Ce dernier nous offre une interface graphique permettant d'ajouter de la fumée dans un espace confiné et, ainsi, d'en observer le comportement. Nous pouvons influencer la quantité de fumée et sa vitesse dans l'espace. De plus, l'application nous donne le contrôle sur la résolution de la simulation, cela revient à dire sur sa précision, qui détermine principalement la performance du programme.

Le déroulement du projet s'est effectué en plusieurs étapes distinctes :

Analyse du code Cette phase consiste à analyser le programme d'un point de vue mathématique et informatique. De cette première approche, il s'agira de comprendre les opérations du programme sur les équations qui régissent le système physique. De l'autre approche, il convient d'étudier l'architecture logicielle de l'application, ainsi que les choix mis en œuvre afin d'implémenter le ou les algorithmes nécessaires.

Protocole expérimental Une fois l'analyse effectuée, nous pouvons en déduire le moyen le plus adapté afin de mesurer les performances de notre implémentation. Nous allons donc mettre en avant les critères théoriques à atteindre dans nos mesures, puis nous exposerons la manière dont nous avons mis ceci en pratique.

Optimisations et mesures Grâce au protocole mis en place, nous pouvons quantifier la performance du programme. De ce fait, nous serons en mesure d'expérimenter différentes techniques d'optimisation sur le programme et d'en calculer l'accélération.

2 Analyse du code

L'analyse du code est la première étape à effectuer afin d'optimiser un code. Il s'agit d'identifier les différentes sections et fonctions qui pourraient être critiques, ainsi que de comprendre mathématiquement quelles fonctions remplissent quels rôles.

L'interface du programme et l'interaction avec l'utilisateur est géré par les fichiers « demo.html », « FluidSolver.java » et « WebStart.java ». À priori, ces fichiers nous importe peu dans notre processus d'analyse et d'optimisation. Le cœur de la simulation se déroule dans le fichier « fluid.c », dans lequel se situe les fonctions de calcul des phénomènes physique.

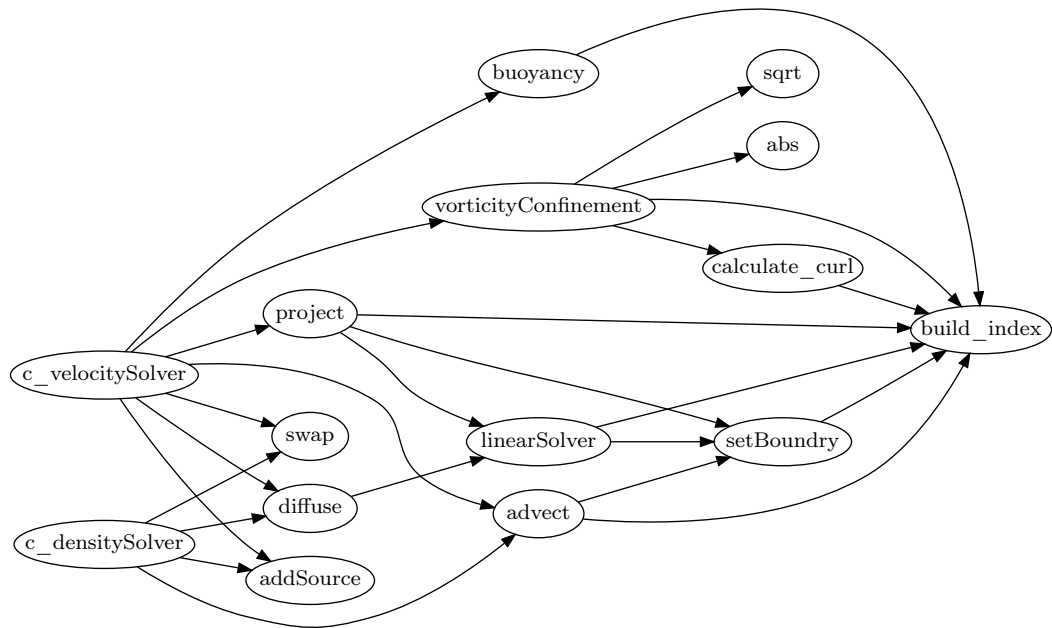


FIGURE 1 – Graphe d'appel du fichier *fluid.c*

Avant d'observer et d'analyser attentivement le code du fichier « fluid.c », nous avons utilisé **Cflow** afin d'obtenir un graphe d'appel pour avoir une vue d'ensemble et observer les liens entre les fonctions. Le résultat est présenté dans la figure 1. Nous décrivons ici les fonctions principales. On remarque très rapidement les deux points d'entrées du programme : « *c_densitySolver()* » et « *c_velocitySolver()* ». Ces deux fonctions sont utilisées conjointement pour assurer les deux fonctionnalités de notre simulation : la première permet de calculer la position d'une particule de fumée ajouté par la souris, la deuxième permet de calculer leur déplacement en fonction des différentes vitesses dans l'espace. Nous observons ensuite deux fonctions qui sont ici pour assurer des opérations physique : « *project()* » et « *diffuse()* ». Le nom laisse sous entendre que la première permet de faire une projection (des coordonnées d'un système à un autre ?) et la seconde permet de diffuser (une particule dans l'espace ?), mais nos connaissances en mécanique des fluides nous arrêtent ici pour les suppositions. Cependant, il est intéressant d'observer que ces deux fonctions font appel à une nouvelle fonction : « *linearSolver()* ». Cette fonction, comme le nom le laisse entendre, est le cœur de calcul de l'application : elle permet de résoudre

un système d'équation linéaire. Il ne fait aucune doute que c'est dans cette fonction que de nombreuses optimisations seront possibles ! Enfin, on observe quelques autres fonctions de physique (« `vorticityConfinement()` », « `advect()` », « `buoyancy()` ») ou des fonctions de mathématiques (« `sqrt()` », « `abs()` ») qui gravitent autour des fonctions précédemment décrites. Enfin, une fonction semble être appelée par quasiment toutes les autres de manière récursive : la fonction « `build_index()` ». Cette fonction, très courte, permet de passer d'un couple de coordonnées (i, j) identifiant une donnée à un offset correspondant à son emplacement mémoire par rapport à la base de la matrice. Cette fonction, présente à de nombreux endroits dans le code, devrait aussi avoir droit à une attention particulière.

3 Protocole expérimental

La mise en place d'un protocole expérimental de mesure est une étape nécessaire et cruciale dans toute optimisation de code. D'une part, le but de ce protocole est de mettre en lumière les points chauds du programme, c'est-à-dire les parties du code qui ralentissent considérablement l'exécution de la simulation. Ces points chauds seront les cibles de nos optimisations. D'autre part, après chaque tentative d'optimisation, le protocole doit nous permettre de mesurer l'impact de cette dernière, qu'il soit positif ou négatif, et enfin de le quantifier.

3.1 Théorie

Le protocole se doit d'être le plus représentatif possible d'une utilisation typique du programme. Pour cela, on ne doit pas uniquement mesurer la simulation exécutant une fonctionnalité particulière, mais l'ensemble des fonctionnalités. Il est bon de noter que ceci est vrai dans le cas d'un unique protocole, comme ici : en effet, on pourrait imaginer plusieurs protocoles différents pour optimiser chaque fonctionnalité indépendamment des autres. Ici, le programme nous offre deux fonctionnalités : ajouter de la fumée dans un espace confiné et ajouter une vitesse (c'est-à-dire un vecteur vitesse) modifiant ainsi la trajectoire de la fumée. Notre protocole doit donc faire intervenir ces deux possibilités dans notre expérience.

Lors de nos expériences, il ne faut pas oublier que le hasard ou l'aléa des mesures peuvent biaiser un résultat. Afin d'éviter cela, il faut donc utiliser une valeur moyenne ou une valeur médiane.

3.2 Pratique

Pour effectuer nos mesures, il nous faut utiliser un outil de d'analyse de performance. Plusieurs possibilités s'offrent à nous, parmi lesquels nous pouvons notamment citer : Performance Counters for Linux (perf), OProfile, Gprof, MAQAO, Callgrind/Kcachegrind, et bien d'autre. Parmi ces choix-là, nous avons choisis d'utiliser « perf ». Nous pouvons lui citer plusieurs avantages qui nous ont intéressés : intégration en ligne de commande, outil bas-niveau, développement stable, accès aux compteurs hardware. Nous pouvons donc en déduire que c'est un outil tout à fait adapté à l'analyse de code bas-niveau (C, assembleur).

Le programme de simulation utilise une interface graphique pour rentrer en interaction avec l'utilisateur. De ce fait, nous ne pouvons pas lui donner un jeu de données défini à l'avance qui serait constant entre les différentes mesures. Nous devons alors définir des actions à réaliser de la manière la plus similaire possible entre les exécutions. De manière « arbitraire », nous définissons un run de la manière suivante : une fois l'applet lancé, nous traçons une ligne de fumée en diagonale, du point en bas à droite au point en haut à gauche [clic gauche enfoncé], puis nous refaisons la même opération en appliquant une vitesse sur ce trajet [clic droit enfoncé], dans le même sens.

Cependant, même en appliquant les mêmes actions entre chaque expérience, les temps d'exécution seront irrémédiablement différents. Dès lors, pour mesurer la performance d'une fonction, nous n'allons pas regarder le nombre de cycles passés dans cette fonction mais son pourcentage de cycle par rapport à l'ensemble du programme. En effet, étant donné que le pourcentage est un rapport sur le temps de l'ensemble du programme, ce dernier peut varier tout en gardant un pourcentage similaire pour une même fonction, car le nombre de cycles passé dans cette fonction augmentera proportionnellement.

Pendant nos expériences, nous choisirons une grille de 225x225 car c'est le premier seuil où l'application devient trop lente pour être utilisée sur notre machine de test. Cela nous permet de mieux nous rendre compte de nos optimisations.

Lors de nos mesures préliminaires, on observe au maximum un écart de $\pm 5\%$ entre les différentes mesures successives. Nous choisissons donc, pour chaque résultat, d'effectuer 5 tests et de prendre la valeur médiane : cela nous semble suffisant.

4 Optimisations et mesures

Le compilateur utilisé pour ce projet est GNU Compiler Collection (GCC) v8.2.0. La mise en place de la compilation consiste à déterminer les bons flags de compilation sur lesquels nous allons partir pour optimiser notre code. Nous ajoutons tout d'abord « -g3 » afin d'inclure dans le binaire les symboles et informations sur le code source. Nous nous assurons que ce flag n'as pas d'incidence sur les performances. Ci-dessous les résultats des analyses du projet au départ, avec GCC :

Samples: 79K of event 'cycles:ppp', Event count (approx.): 63306420665

Overhead	Command	Shared Object	Symbol
46,86%	appletviewer	libfluid.so	[.] linearSolver
23,99%	appletviewer	libfluid.so	[.] build_index
7,86%	appletviewer	libfluid.so	[.] build_index@plt
2,91%	appletviewer	libfluid.so	[.] advect
2,09%	appletviewer	libfluid.so	[.] project
1,11%	appletviewer	libfluid.so	[.] vorticityConfinement
0,92%	appletviewer	libfluid.so	[.] SWIG_JavaArrayInFloat
0,71%	appletviewer	libfluid.so	[.] SWIG_JavaArrayArgoutFloat
0,70%	appletviewer	libc-2.28.so	[.] __memmove_avx_unaligned_erms
0,68%	appletviewer	libfluid.so	[.] calculate_curl
0,65%	appletviewer	libfluid.so	[.] addSource
0,65%	appletviewer	libfluid.so	[.] swap

FIGURE 2 – Profilage global avec GCC en « -O0 »

Comme nous pouvions l'imaginer, la Figure 2 nous confirme que les fonctions les plus chaudes sont « linearSolver() » et « build_index() », la première de par sa nature de calcul intensif, la deuxième du fait qu'elle soit appelée énormément de fois. On remarque que toutes les autres fonctions sont d'une importance mineure dans le temps d'exécution du programme : inférieure à 3% pour la plus chaude. À elle seule, la fonction « linearSolver() » prend presque 50% du temps d'exécution du programme : c'est la plus importante.

La Figure 3 nous permet d'analyser quelles instructions de la fonction « linearSolver() » ont un impact important sur le temps d'exécution. La première image nous montre un lot successif d'instructions prenant chacune 1 à 4% du temps d'exécution : des *add*, *mul*, *mov*, *lea*. La deuxième image met en évidence une instruction *div* qui prends à elle seule 16% du temps de la fonction. On identifie donc parfaitement les deux sources de ralentissement de la fonction : la première étant un trop grand nombre d'instructions pour effectuer les calculs et la seconde étant la division.

Le temps passé dans la fonction « build_index() » est divisé équitablement en *mov/add* d'après la Figure 4, permettant de faire les deux additions et la multiplication que doit faire la fonction. On remarque tout de même que 32% du temps de la fonction est passé dans la mise en place de la stack frame de la fonction (2 premières instructions *mov*). Enfin, on remarquera que 7% du temps de la simulation est passée dans la fonction « build_index@plt ». Pour comprendre d'où vient cette fonction, il faut revenir sur les bases du fonctionnement du chargement dynamique de code. Afin d'avoir des bibliothèques partagées qui puissent être chargées qu'une seule fois en mémoire mais à des emplacements indépendants pour deux processus différents, on utilise un mécanisme de code à position indépendante (Position-Independent

Code (PIC)). Ce système fait recours à une Global Offset Table (GOT) contenant les adresses des variables et fonctions dont l'emplacement n'est pas connu au moment de la compilation. Cette GOT est consultée par les fonctions contenues dans la Procedure Linkage Table (PLT), une table contenant des fonctions *stubs* (wrapper) qui sont appelées à la place de la vraie fonction. Cette fonction stub (@plt) permet d'aller consulter la GOT, de remplir la GOT lors du premier appel (*lazy loader*), enfin d'exécuter la vraie fonction. Toutes les prochaines fois où la fonction stub sera appelée, alors on va directement chargé l'adresse contenu dans la *GOT* pour appeler la fonction désirée.

4.1 Mise en place de la compilation

Premièrement, nous avons testé les flags suivants activant des options de lot d'optimisation : « -O1 », « -O2 », « -O3 » (agressivité au détriment de la taille du code généré), « -Os » (optimisation de la taille du code). L'utilisation d'un de ces flags nous apporte une légère amélioration de fluidité visuelle, cependant il se trouve que nous n'avons pas observé de différence majeure entre ces 4 différentes flags, que ce soit en termes de code généré ou de performances et répartition du temps entre les fonctions. Ci-dessous les images des résultats avec GCC et « -O2 » :

Nous avons aussi testé les performances du programme avec « Clang/Low Level Virtual Machine (LLVM) », et les flags d'optimisation les plus agressifs (« -O3 »). Ci-dessous les images des résultats avec Clang et « -O3 » :

Nous observons donc un net gain à utiliser Clang par rapport à GCC. Cependant, pour nous permettre d'appliquer des optimisations vu en cours à la main, nous allons opter pour la version que produit GCC comme base de travail.

4.2 Inlining

TODO

4.3 Déroulage de boucle

TODO

4.4 Vectorisation

TODO

Percent	
	mov -0x8(%rbp),%eax
	lea -0x1(%rax),%ecx
2,18	mov -0x34(%rbp),%edx
0,23	mov -0x4(%rbp),%eax
0,00	mov %ecx,%esi
	mov %eax,%edi
2,58	→ callq build_index@plt
0,02	cltq
2,48	lea 0x0(,%rax,4),%rdx
0,01	mov -0x20(%rbp),%rax
0,05	add %rdx,%rax
0,53	movss (%rax),%xmm0
	x[build_index(i,
2,31	movaps %xmm0,%xmm4
0,89	addss -0x38(%rbp),%xmm4
0,39	movss %xmm4,-0x38(%rbp)
	mov -0x8(%rbp),%eax
2,25	lea 0x1(%rax),%ecx
0,00	mov -0x34(%rbp),%edx
0,39	mov -0x4(%rbp),%eax
	mov %ecx,%esi
2,31	mov %eax,%edi
0,35	→ callq build_index@plt
0,03	cltq
2,43	lea 0x0(,%rax,4),%rdx
0,01	mov -0x20(%rbp),%rax
0,06	add %rdx,%rax
2,47	movss (%rax),%xmm0
4,32	addss -0x38(%rbp),%xmm0
	x[build_index(i,
4,13	mulss -0x18(%rbp),%xmm0
1,34	movss %xmm0,-0x38(%rbp)
	mov -0x34(%rbp),%edx
1,23	mov -0x8(%rbp),%ecx
0,01	mov -0x4(%rbp),%eax
1,36	mov %ecx,%esi
0,00	mov %eax,%edi
1,09	→ callq build_index@plt
1,27	cltq
1,18	lea 0x0(,%rax,4),%rdx
	mov -0x28(%rbp),%rax
0,01	add %rdx,%rax
4,36	movss (%rax),%xmm0
4,23	addss -0x38(%rbp),%xmm0

		x0[build_index(i, j, grid_size))] / c;
1,54	movss -0x38(%rbp),%xmm0	
16,60	divss -0x2c(%rbp),%xmm0	

FIGURE 3 – Profilage de « linearSolver() » avec GCC en « -O0 »

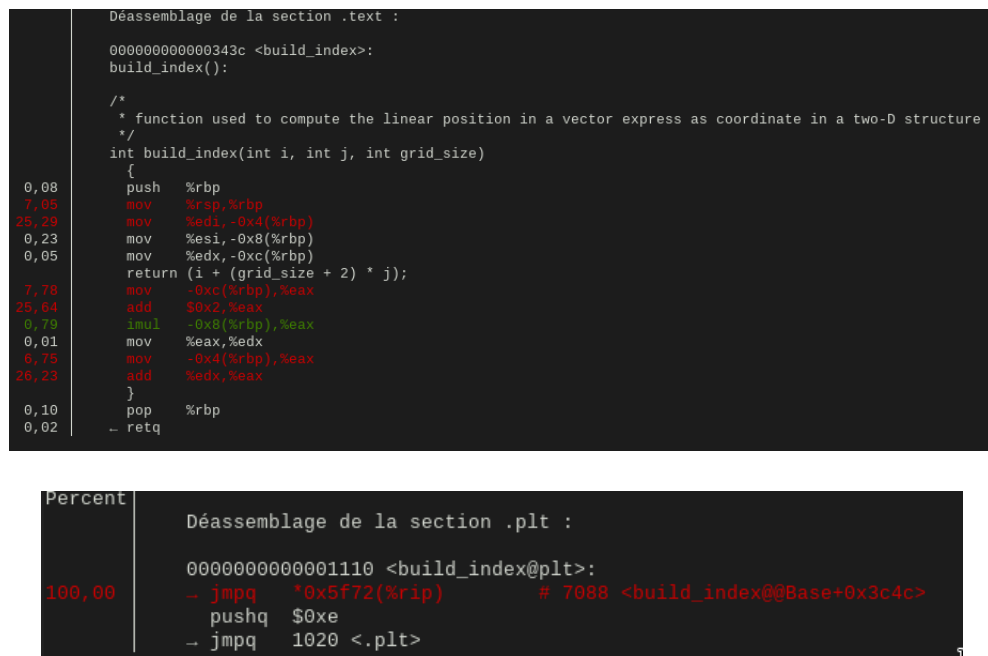


FIGURE 4 – Profilage de « build_index() » et « build_index()@plt » avec GCC en « -O0 »

5 Conclusion

TODO

Acronymes

GCC GNU Compiler Collection

LLVM Low Level Virtual Machine

PIC Position-Independent Code

GOT Global Offset Table

PLT Procedure Linkage Table