



RAPPORT DU PROJET DE PROGRAMMATION NUMÉRIQUE

ENCADRÉ PAR : JEAN-BAPTISTE BESNARD

RÉALISÉ PAR :

KARIM SMAIL

BOUZAHER SOFIANE

KREDDIA ASMA

DORAI ATTEF

6 JANVIER 2020

INTRODUCTION

- Le but est la mise en place d'un plugin dans le compilateur CLANG (LLVM) pour identifier les appels de fonctions candidats à la « taskification » (fonction pures)
- Pour ce faire, il faudra mettre en place la structure de base d'un plugin clang
- Définir ce qui distingue les fonctions candidates à identifier
- Implémenter ce support dans le plugin
- Le valider sur un cas de parallélisme producteur / consommateur.

TITRE DU SUJET

Analyse Statique Pour la Classification des Procédures Candidate à la « Taskification »

MOT CLEF

LLVM, Analyse statique, compilation, MPI + X, parallélisation automatique

DESCRIPTION GENERALE

Les architecture hybrides convergées à venir posent la question des modèles de programmation. En effet MPI depuis l'avènement des architectures many-core a dû être combiné avec du parallélisme intra-noeud en OpenMP (MPI + X). Le mélange de ces modèles se traduit nécessairement par une complexité accrue de l'expression des codes de calcul. Dans ce travail nous proposons de prendre cette tendance à contre-pied en posant la question de l'expression de tâche de calcul en pur MPI. Les étudiants se verront fournir une implémentation de Remote Procedure Calls (RPC) implémentés en MPI, le but du travail et de détecter quelles fonctions sont éligibles à la sémantique RPC statiquement lors de la phase de compilation (c.a.d. les fonction dites « pures » : indépendantes du tas, des TLS, etc ...). Le travail visera le compilateur LLVM dans lequel une passe sera rajoutée pour lister l'ensemble des fonctions éligibles à la sémantique RPC. Pour exemple, une implémentation d'un algorithme de cassage de mot de passe en MPI sera fournie avec pour but sa conversion en RPC producteur/consommateur (github.com/besnardjb/MPI_Brute/) avec l'outil.

PRESENTATION GENERALE DE L'OBJECTIF

Le but d'un programme est d'exécuter une tâche. Pour réaliser celle-ci, on donne à l'ordinateur une liste d'instructions qu'il va effectuer. Il existe plusieurs manières de traiter ces instructions, parmi ces manières, on trouve la programmation parallèle.

1-Pourquoi le parallélisme ?

L'exécution de certaines fonctions d'un programme de manière parallèle, nous permet un gain de temps d'exécution du programme, ce qui veut dire rendre le programme plus performant qu'avant, mais ça ne marche pas avec toutes les fonctions de tous les programmes, Donc les fonctions qui seront exécuter de manière parallèle doit être connu à la compilation du programme.

2-C'est quoi un compilateur ?

Un compilateur est un programme qui transforme un code source (écrit dans un langage de programmation de haut niveau d'abstraction) en un code objet (écrit dans langage de programmation de bas niveau) afin de créer un programme exécutable par une machine.

3-Le But de ce travail :

Le but est la mise en place d'un plugin dans le compilateur CLANG (LLVM) pour identifier les appels de fonctions candidats à la « taskification » (fonction pures) , Pour ce faire, il faudra mettre en place la structure de base d'un plugin clang ,définir ce qui distingue les fonctions candidates à identifier , et enfin implémenter ce support dans le plugin et Le valider sur un cas de parallélisme producteur / consommateur.

DEFINITIONS UTILES

1) LLVM (Low Level Virtual Machine) :

est une infrastructure de compilateur conçue pour l'optimisation du code à la compilation (c'est une infrastructure qui ne contient pas des outils nécessaires pour compiler du code source C ou C++ mais uniquement des outils d'optimisations et de génération de codes machine à partir d'un format intermédiaire).

CLANG-LLVM :

la structure générale de cette infrastructure à l'échelle microscopique est constitué d'une façon similaire à tout compilateur moderne .

2) Analyse statique :

définie comme étant l'ensemble des méthodes utilisées pour obtenir des informations sur le comportement d'un programme lors de son exécution sans réellement l'exécuter.

3) Analyse dynamique (dynamic program analysis) :

contrairement à l'analyse statique, est une analyse qui nécessite l'exécution du programme et étudie son comportement +les effets de son exécution sur son environnement.

4) Compilation :

définie comme l'ensemble des étapes qui transforment un code (. C) en un code objet(. O). cela se fait à l'aide d'un programme appelé 'Compilateur '.

5) MPI :

Message Passing interface, est une norme conçue pour le passage de messages entre ordinateurs distants ou dans un ordinateur multiprocesseur (donc c'est un moyen de transfert de message créer pour obtenir de meilleures performances), Elle est devenue un standard de communication pour des nœuds exécutant des programmes parallèles sur des systèmes à mémoire distribuée. Elle définit une bibliothèque de fonctions, utilisable avec les langages C, C++ et Fortran.

6) Open MPI (MPI+X) :

(une bibliothèque MPI) est une bibliothèque de projet qui sert à combiner l'expertise, les techniques et les ressources de toute la communauté du calcul haute performance afin de créer la meilleure bibliothèque MPI disponible.

7) La parallélisation automatique :

est une étape de la compilation d'un programme qui consiste à transformer un code source écrit pour une machine séquentielle en un exécutable parallélisé pour ordinateur à un multiprocesseur symétrique.

8) Un multiprocesseur symétrique (ou symmetric shared memory multiprocessor -SMP) :

est une architecture parallèle qui consiste à multiplier les processeurs identiques au sein d'un ordinateur, de manière à augmenter la puissance de calcul.(Cela se fait en exécutant simultanément plusieurs processus du système).

9) Remonte Procédure Calls (RPC) :

est un protocole réseau permettant de faire des appels de procédures sur un ordinateur distant à l'aide d'un serveur d'applications (utilisé par exemple dans le modèle client- serveur).

10) Mémoire partagée :

(communication interprocessus) : est un moyen de partager des données entre différents processus donc la même zone mémoire peut être accédée par plusieurs processus .

11) Mémoire distribuée :

est une mémoire répartie en plusieurs nœuds, chaque portion n'étant accessible qu'à certains processus(la communication entre nœuds se fait à l'aide de MPI).

12) Open Mp (Open multi-processing) :

est une interface de programmation pour le calcul parallèle sur architecture à mémoire partagée. Cette API est prise en charge par de nombreuses plateformes, incluant GNU/Linux, OS X et Windows, pour les langages de programmation C, C++ et Fortran. Il se présente sous la forme d'un ensemble de directives, d'une bibliothèque logicielle et de variables d'environnement.

13) Cilk (Cilk ++ Cilk Plus) :

sont des langages de programmation à usage général conçus pour le calcul parallèle multithread. Ils sont basés sur les langages de programmation C et C ++, qu'ils étendent avec des constructions pour exprimer des boucles parallèles.

14) Le message actif :

est un objet de messagerie capable d'effectuer seul le traitement (contrairement aux systèmes de messagerie informatiques traditionnels dans lesquels les messages sont des entités passives sans puissance de traitement). donc Il s'agit d'un protocole de messagerie utilisé pour optimiser les communications réseau en mettant l'accent sur la réduction de la latence (le délai de transmission dans les communications)...

15) RDMA :

comme son nom l'indique, c'est un accès direct à la mémoire de la mémoire d'un ordinateur à celle d'un autre sans impliquer le système d'exploitation de l'un ou de l'autre (cela se manifeste à l'aide de la carte réseau).

16) DMA(direct memory accès) :

est une propriété du computer systems qui permet à certains hardwares d'accéder à la mémoire indépendamment au CPU .

17) un tas (heap) :

est une structure de données de type arbre permettant de retrouver directement l'élément que l'on veut traiter en priorité. Aussi définie comme étant l'un des deux segments de mémoire (on cite aussi le Pile d'exécution où Stack) utilisé lors de l'allocation dynamique de mémoire durant l'exécution d'un programme informatique.

18) Transport Layer Security (TLS) :

sont des protocoles de cryptage qui garantissent pleinement la sécurité des communications pour toutes les données échangées. Ces systèmes sont largement utilisés pour garantir la sécurité des communications sur internet.

Première partie

Fonctions pures - Fonctions impures

Définition d'une fonction pure :

une fonction pure ne dépend pas et ne modifie pas l'état de variables hors de sa portée.

En pratique, cela signifie qu'une fonction pure retourne toujours le même résultat avec des paramètres identiques.

Son exécution ne dépend pas de l'état du système.

C'est-à-dire elle possède les propriétés suivantes :

1. Sa valeur de retour est la même pour les mêmes arguments (pas de variation avec des variables statiques locales, des variables non locales, des arguments mutables de type référence ou des flux d'entrée).

2. Son évaluation n'a pas d'effets de bord :

En informatique, une fonction est dite à effet de bord (effet secondaire) si elle modifie un état en dehors de son environnement local, c'est-à-dire a une interaction observable avec le monde extérieur autre que retourner une valeur.

Par exemple, les fonctions qui modifient une variable locale statique, une variable non locale ou un argument mutable passé par référence,

les fonctions qui effectuent des opérations d'entrées-sorties ou les fonctions appelant d'autres fonctions à effet de bord.

Souvent, ces effets compliquent la lisibilité du comportement des programmes et/ou nuisent à la réutilisabilité des fonctions et procédures.

L'avantage principal d'une fonction pure :

-l'appel à cette fonction avec les mêmes paramètres renverra toujours le même résultat.

-On simplifie également la mise en place des tests automatiques, ce qui sécurise notre application.

-Les fonctions pures ont pour avantage d'être prédictibles.

Ce qui permet de les tester plus facilement et surtout de mettre leur résultat en cache pour ne pas avoir à refaire le calcul pour des valeurs qu'on a déjà traitées.

Les fonctions pures sont souvent utilisées pour générer d'autres fonctions. Dans ce cas, elles sont appelées "Higher Order Functions" ou "Fonctions de rang supérieur".

Note : Une fonction de rang supérieur peut ne pas être une fonction pure.

Exemples de Fonctions pures en C :

Les fonctions arithmétiques sont l'archétype des fonctions pures.

Les fonctions suivantes sont pures :

1.La fonction floor : retournant la partie entière par défaut d'un nombre :

Cette fonction retourne la valeur minimale d'un nombre, soit l'entier le plus proche inférieur ou égal au nombre.

Voici un exemple montrant une utilisation plus classique de cette fonction :

```
double floor( double value );
```

2.La fonction max(resp.min) : retournant le maximum(resp.minimum) de deux variable :

```
double MAX(double X, double Y) if (X>Y) return X; else return Y;
```

Définition d'une fonction impure :

Une fonction impure est une fonction qui peut avoir des effets de bords(elle peut aussi accidentellement ne pas en avoir).

le résultat de la fonction peut dépendre du contexte, et son exécution peut le modifier .

une fonction de comptage (qui rend le nombre de fois où elle a été appelée, nécessitant donc la modification d'une variable externe) ,ou les fonctions NOW du paquetage STANDARD qui rendent l'heure qui est dans le monde simulé, et donc une valeur différente a chaque appel, sont des fonctions impures.

des cas qui rendent une fonction impure en C :

Les fonctions C suivantes sont impures car elles ne vérifient pas la propriété 1 ci-dessus :

1/à cause de la variation de la valeur de retour avec une variable non locale :

```
int f(int x) return x;
```

2/à cause de la variation de la valeur de retour avec un argument mutable de type référence :

```
(int *) f() // la fonction f renvoie un pointeur sur entier int x = 42; return x; // l'adresse de x est bien un pointeur sur entier
```

3/à cause de la variation de la valeur de retour avec un flux d'entrée :

```
int f(int x) x = 0; scanf("return x;
```

Les fonctions C suivantes sont impures car elles ne vérifient la propriété 2 ci-dessus :

1/L'effet de bord de la fonction f ici est de modifier la valeur de la variable globale x :

```
void f(int x) x = 1;
```

2/à cause de la mutation d'une variable statique locale :

```
void f(void) static int i = 0; /* i sera initialisée à 0 à la compilation seulement */  
i++;  
printf("i vaut
```

3/à cause de la mutation d'un argument mutable de type référence :

```
void f(int* a) *a = 2;
```


4/à cause de la mutation d'un flux de sortie

```
void f(void) printf("Hello.");
```

Deuxième partie

L'infrastructure CLANG – LLVM

Définition de L'infrastructure clang - LLVM :

D'une façon macroscopique, elle est construite d'une manière similaire à tout compilateur moderne, elle ne contient pas les outils nécessaires pour compiler du code source C ou C++ mais uniquement des outils d'optimisation et de génération de codes machines à partir d'un format intermédiaire.

- Pour mieux comprendre l'ensemble des étapes de compilation par cette infrastructure, on doit définir en détails :

1. le Frontend :

C'est le premier bloc de tout compilateur, son objectif est de valider que le programme est syntaxiquement et sémantiquement correct puis de le traduire vers une représentation intermédiaire (IR pour Intermediate Representation) l'un des objectifs de cette représentation intermédiaire étant de simplifier le travail des autres blocs qui ne peuvent pas travailler avec la complexité de code source C ou encore pire C++.

2. les Passes :

Sont en charge d'analyser et/ou de transformer l'IR en optimisant certaines choses tout en préservant la sémantique du code. Son objectif était très souvent la maximisation des performances du code (par exemple en jouant sur la taille du code).

3. le backend :

Est en charge de transformer l'IR vers du code machine pour une architecture donnée.

Afin d'intervenir aux différents niveaux de la chaîne de compilation, plusieurs outils vont venir intervenir tels que :

CLANG (appelé Driver de compilation) :

quand on dit clang on dit frontend C/C++ de l'infrastructure LLVM En tant que driver de compilation, l'outil clang peut-être arrêté à différents niveaux de la chaîne de compilation, et pour mieux comprendre ce point je vous montre des exemples d'utilisation de l'outil clang :

Cas 1 :

`clang -S -emit-llvm -o test.ll test.c =>` dans ce cas clang est utilisé comme étant un frontend c'est-à-dire générer IR textuel à partir du code source

Cas 2 :

`clang -o test.bin test.s =>` clang est utilisé comme assembleur et linker

OPT :

cet outil permet d'appliquer un ensemble de passes LLVM, l'entrée d'OPT est un fichier au format IR (bit code ou textuel) et la sortie produite est également un fichier au format IR (bit code ou textuel).

Pour le choix des passes a appliqué, l'outil OPT peut-être utilisé avec les options usuelles -O1, -O2, -O3 (si aucune des options -Ox n'est pas spécifiée, OPT n'applique aucune passe).

On peut aussi spécifier individuellement les passes que nous souhaitons appliquer et on prend comme exemple la commande suivante :

`opt -S -mem2reg -constprop -o test-after-cp.ll test.ll` telque `mem2reg` et `constprop` sont les passes qu'on a choisi.

LLC :

Il permet de compiler du code au format LLVM IR (bitcode ou textuel) vers du code assembleur pour une architecture donnée.

LLVM-AS et LLVM-DIS :

`llvm-as` et `llvm-dis` permettent respectivement de passer du format LLVM IR textuel au format LLVM IR bitcode et inversement.

LLI :

Cet outil permet d'exécuter du code au format LLVM IR (bitcode ou textuel) non pas en le compilant vers du code machine mais en l'interprétant directement.