

# Langage de Programmation et Compilation

Jean-Cristophe Filiâtre

29 septembre 2023

## Table des matières

<b>I</b>	<b>Cours 1 29/09</b>	<b>1</b>
0.1	Un Compilateur . . . . .	1
0.2	Le Bon et le Mauvais Compilateur . . . . .	1
0.3	Le Travail d'un Compilateur . . . . .	2
<b>1</b>	<b>L'assembleur</b>	<b>2</b>
1.1	Arithmétique des ordinateurs . . . . .	2
1.2	Architecture . . . . .	2
1.3	L'architecture x86-64 . . . . .	3
1.4	L'assembleur x86-64 . . . . .	3
1.5	Le Défi de la Compilation . . . . .	3

## Première partie

## Cours 1 29/09

### Introduction

Maîtriser les mécanismes de la compilation, transformation d'un langage dans un autre. Comprendre les aspects des langages de programmation.

#### 0.1 Un Compilateur

Un compilateur est un traducteur d'un langage source vers un langage cible. Ici le langage cible sera l'assembleur.

Tous les langages ne sont pas compilés à l'avance, certains sont interprétés, transpilés puis interprétés, compilés à la volée, transpilés puis compilés... Un compilateur prend un programme  $P$  et le traduit en un programme  $Q$  de sorte que :  $\forall P, \exists Q, \forall x, P(x) = Q(x)$ . Un interpréteur effectue un travail simple mais le refait à chaque entrée, et donc est moins efficace.

Exemple : le langage *lilypond* va compiler un code source en fichier .pdf.

#### 0.2 Le Bon et le Mauvais Compilateur

On juge un compilateur à :

1. Sa correction
2. L'efficacité du code qu'il produit
3. Son efficacité en tant que programme

« Optimizing compilers are so difficult to get right that we dare say that no optimizing compiler is completely error-free! Thus, the most important objective in writing a compiler is that it is correct »- *Dragon Book, 2006*

### 0.3 Le Travail d'un Compilateur

Le travail d'un compilateur se compose :

- d'une phase d'analyse qui :
  1. reconnaît le programme à traduire et sa signification
  2. signale les erreurs et peut donc échouer
- d'une phase de synthèse qui :
  1. produit du langage cible
  2. utilise de nombreux langages intermédiaires
  3. n'échoue pas

Processus : source  $\rightarrow$  analyse lexicale  $\rightarrow$  suite de lexèmes (tokens)  $\rightarrow$  analyse syntaxique  $\rightarrow$  Arbre de syntaxe abstraite  $\rightarrow$  analyse sémantique  $\rightarrow$  syntaxe abstraite + table des symboles  $\rightarrow$  production de code  $\rightarrow$  langage assembleur  $\rightarrow$  assembleur  $\rightarrow$  langage machine  $\rightarrow$  éditeur de liens  $\rightarrow$  exécutable.

## 1 L'assembleur

### 1.1 Arithmétique des ordinateurs

On représente les entiers sur  $n$  bits numérotés de droite à gauche. Typiquement,  $n$  vaut 8, 16, 32 ou 64. On peut représenter des entiers non signés jusqu'à  $2^n - 1$ . On peut représenter les entiers en définissant  $b_{n-1}$  comme un bit de signe, on peut alors représenter  $[-2^{n-1}, 2^{n-1} - 1]$ . La valeur d'une suite de bits est alors :  $-b_{n-1}2^{n-1} + \sum_{k=0}^{n-2} b_k 2^k$ . On ne peut pas savoir si un entier est signé sans le contexte.

La machine fournit des opérations logiques (bit à bit), de décalage (ajout de bits 0 de poids fort, 0 de poids faible ou réplication du bit de signe pour interpréter une division), d'arithmétique (addition, soustraction, multiplication).

### 1.2 Architecture

Un ordinateur contient :

- Une unité de calcul (CPU) qui contient un petit nombre de registres et des capacités de calcul
- Une mémoire vive (RAM), composée d'un très grand nombre d'octets (8 bits), et des données et des instructions, indifférenciables sans contexte.

L'accès à la mémoire coûte cher : à 1B instructions/s, la lumière ne parcourt que 30cm entre deux instructions.

En réalité, il y a plusieurs (co)processeurs, des mémoires cache, une virtualisation de la mémoire... Principe d'exécution : un registre (%rip) contient l'adresse de l'instruction, on lit (fetch) un ou plusieurs octets dans la mémoire, on interprète ces bits (decode), on exécute l'instruction (execute), on modifie (%rip) pour l'instruction suivante. En réalité, on a des pipelines qui branchent plusieurs instructions en parallèle, et on essaie de prédire les sauts conditionnels.

Deux grandes familles d'Architectures : CISC (complex instruction set), qui permet beaucoup d'instructions différentes mais avec assez peu de registres, et RISC (Reduced Instruction Set) avec peu d'instruction effectuées très régulièrement et avec beaucoup de registres. Ici, on utilisera l'architecture *x86-64*.

### 1.3 L'architecture x86-64

Extension 64 bits d'une famille d'architectures compatibles Intel par AMD adoptée par Intel. Architecture à 16 registres, avec adressage sur 48 bits au moins et de nombreux modes d'adressage. On ne programme pas en langage machine mais en assembleur, langage symbolique avec allocation de données globales, qui est transformé en langage machine par un assembleur qui est en réalité un compilateur. On utilise l'assembleur GNU avec la syntaxe AT&T (la syntaxe Intel existe aussi).

### 1.4 L'assembleur x86-64

Pour assembler un programme assembleur, appeler `as -o file.o` puis appeler l'édition de lien avec `gcc -no-pie file.s -o exec-name`. On peut déboguer en ajoutant l'option `-g`. La machine est petite boutiste (little-endian) si elle stocke les valeurs dans la RAM en commençant par le bit de poids faible, gros boutiste (big-endian) pour le poids fort.

Commandes : Dans cette liste, `%(r)` désigne l'adresse mémoire stockée dans `r`

- `movq $a %b` permet de mettre la valeur `a` dans le registre `b`
- `movq %a %b` permet de copier le registre `a` dans le registre `b`
- `movq $label %b` permet de changer l'adresse de l'étiquette dans le registre `b`
- `addq %a %b` permet d'additionner les registres `a` et `b`.
- `incq %r` permet d'incrémenter le registre `r`, de même pour `decq`.
- `negq %r` permet de modifier la valeur de `r` en sa négation
- `notq %r` permet de modifier la valeur de `r` en sa négation logique.
- `orq %r1 %r2` (resp. `andq` et `xorq`) permet d'affecter à `r2`, `or(r1, r2)` (resp. `and`, `xor`)
- `salq $n %r/salq %cl %r` décale la valeur de `r` de `n` (ou `%cl`) zéros à gauche.
- `sarq` est le décalage à droite arithmétique, `shrq` le décalage à droite logique.
- Le suffixe `q` désigne une opération sur 64 bits. `b` désigne 1 octet, `w` désigne 2 octets, `l` désigne 4 octets. Il faut préciser les deux extensions si celles-ci diffèrent.
- `jmp label` permet de jump à une étiquette.

La plupart des opérations positionnent des drapeaux selon leur résultat.

Certaines instructions : `j(suffixe)` (jump), `set(suffixe)` et `cmov(suffixe)`(move) permettent de tester des drapeaux et d'effectuer une opération selon leur valeur.

On ne sait pas combien il y a d'instructions en x86-64.

### 1.5 Le Défi de la Compilation

C'est de traduire un programme d'un langage de haut niveau vers ce jeu d'instruction.

Constat : les appels de fonctions peuvent être arbitrairement imbriqués et les registres ne suffisent pas  $\Rightarrow$  on crée alors une pile car les fonctions procèdent majoritairement selon un mode LIFO.

La pile est stockée tout en haut, et croît dans le sens des adresses décroissantes, `%rsp` pointe sur le sommet de la pile. Les données dynamiques sont allouées sur le tas, en bas de la zone de données. Chaque programme a l'illusion d'avoir toute la mémoire pour lui tout seul, illusion créée par l'OS. En assembleur on a des facilités d'utilisation de la pile :

- `pushq $a` push `a` dans la pile
- `popq %rdi` dépile

Lorsque `f` (caller) appelle une fonction `g` (callee), on ne peut pas juste `jmp g`. On utilise `call g` puis une fois que c'est terminé `ret`.

Mézalor tout registre utilisé par `g` sera perdu par `f`. On s'accorde alors sur des **conventions d'appel**. Des arguments sont passés dans certains registres, puis sur la pile, la valeur de retour est passée dans `%rax`. Certains registres sont *callee-saved* i.e. l'appelé doit les sauvegarder pour qu'elle survive aux appels. Les autres registres sont dit *caller-saved* et ne vont pas survivre aux appels.

Il faut également qu'en entrée de fonction `%rsp + 8` doit être multiple de 16, sinon des fonctions peuvent planter.

Il y a quatre temps dans un appel :

1. Pour l'appelant, juste avant l'appel :
2. Pour l'appelé, au début de l'appel :
  - (a) Sauvegarde **%rbp** puis le positionne.
  - (b) Alloue son tableau d'activation.
  - (c) Sauvegarde les registres *callee-saved*.
3. Pour l'appelé, à la fin de l'appel :
  - (a) Placer le résultat dans **%rax**
  - (b) Restaure les registres sauvegardés
  - (c) Dépile son tableau d'activation
  - (d) Exécute **ret**
4. Pour l'appelant, juste après l'appel :
  - (a) Dépile les éventuels arguments
  - (b) Restaure les registres *caller-saved*