

Compilation et assembleur

5d54626cc43ba7ce625903ffa39c049c

Groupe de sécurité de l'information, École normale supérieure

2 octobre 2019



Définition (Compilateur)

Programme qui transforme un programme écrit dans un langage **source** vers un programme en langage **objet**.

Exemple

gcc : du C vers le langage machine ghc : du Haskell vers le langage machine

En particulier, il est possible de compiler vers un langage intermédiaire très proche de la machine, mais néanmoins compréhensible par des humains : l'**assembleur** (e.g. *x86*, *mips*, *risc-v*, *armv7*).

Le programme qui traduit un programme assembleur en langage machine s'appelle un **assembleur**.

Il est possible de :

- Compiler vers un autre langage de haut-niveau (Haskell vers C)
- Compiler vers un pseudo-langage machine indépendant de la machine (bytecode,) → ces programmes nécessitent pour leur exécution un autre programme appelé interpréteur (JVM, `ocamlrun`, LLVM, ...)
- Compiler un programme pendant son exécution → *Just-in-time* ou *compilation à la volée* (souvent effectué par les interpréteurs sus-cités)
- Utiliser le compilateur d'un langage hôte pour compiler un code dans un autre langage → *embedded domain-specific language*

Structure d'un compilateur

Deux phases successives:

- Analyse : traduction du langage source en représentation abstraite mathématique
 - ▶ Analyse lexicale : traduction du source vers des **lexèmes** (tokens) → `lex`, `bison`
 - ▶ Analyse syntaxique : traduction des lexèmes vers un **arbre de syntaxe abstraite** (AST) → `yacc`, `menhir`
 - ▶ Typage, ...
- Synthèse : traduction de la représentation mathématique en langage objet
 - ▶ Optimisations
 - ▶ Traduction dans des langages intermédiaires plus proches de la machine (*RTL*, *Cminor*, *LLVM bytecode*, ...)
 - ▶ Production de code assembleur
 - ▶ Assemblage
 - ▶ Édition des liens (linking)

Definition (Lexer)

Programme qui transforme une suite de caractères (écrit dans un langage source) en une suite de lexèmes (ou tokens).

Exemple :

```
int b =  
    /* Here is my comment */  
    (int) c + 2 * d;
```

Sera transformé par le lexer en:

`int` `b` `=` `(` `int` `)` `c` `+` `2` `*` `d` `;`

On reconnaît les lexèmes grâce à des expressions régulières (reconnaissables par automate fini, théorème de Kleene, 1956).

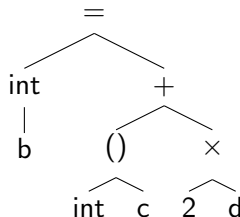
$$\langle \textit{digit} \rangle ::= 0 - 9$$
$$\langle \textit{alpha} \rangle ::= a - z | A - Z$$
$$\langle \textit{ident} \rangle ::= \langle \textit{alpha} \rangle (\langle \textit{alpha} \rangle | \langle \textit{digit} \rangle | _)^*$$
$$\langle \textit{integer} \rangle ::= \langle \textit{digit} \rangle^+$$

Definition (Parser)

Programme qui transforme une suite de lexèmes en un arbre de syntaxe abstraite.

Exemple :

`int b` \equiv `(int) c + 2 * d ;`



Analyse syntaxique (instructions)

On construit l'arbre de syntaxe abstraite grâce à une grammaire (*Backus–Naur form*) :

<i>path</i>	$::=$	<i>ident</i> <i>path</i> . <i>ident</i> <i>path</i> . all	variable record field pointer dereference
<i>expr</i>	$::=$	<i>path</i> <i>integer</i> True False ... <i>expr</i> (+ - < = ...) <i>expr</i> <i>path</i> 'Access null	l-value scalar value binary operator address of an l-value null pointer
<i>stmt</i>	$::=$	<i>path</i> := <i>expr</i> <i>path</i> := new type if <i>expr</i> then <i>stmt</i> * else <i>stmt</i> * end while <i>expr</i> loop <i>stmt</i> * end <i>ident</i> (<i>expr</i> *)	assignment allocation conditional "while" loop procedure call

Analyse syntaxique (déclarations)

procedure ::= *procedure ident (param*) is*
 local begin stmt* end*

param ::= *ident : (in | in out | out) type*

local ::= *ident : type*

<i>type</i> ::=	<i>Integer Real Boolean</i>	scalar type
	<i>access type</i>	access type (pointer)
	<i>ident</i>	record type

Analyse syntaxique (exemple)

https://cs.wmich.edu/~gupta/teaching/cs4850/sumII06/The_syntax_of_C_in_Backus-Naur_form.htm

The syntax of C in Backus-Naur Form

```
<translation-unit> ::= {<external-declaration>}*

<external-declaration> ::= <function-definition>
                          | <declaration>

<function-definition> ::= {<declaration-specifier>}* <declarator> {<declaration>}* <compound-statement>

<declaration-specifier> ::= <storage-class-specifier>
                          | <type-specifier>
                          | <type-qualifier>

<storage-class-specifier> ::= auto
                          | register
                          | static
                          | extern
                          | typedef

<type-specifier> ::= void
                  | char
                  | short
                  | int
                  | long
                  | float
                  | double
                  | signed
                  | unsigned
                  | <struct-or-union-specifier>
                  | <enum-specifier>
                  | <typedef-name>

<struct-or-union-specifier> ::= <struct-or-union> <identifier> { {<struct-declaration>}+ }
                             | <struct-or-union> { {<struct-declaration>}+ }
                             | <struct-or-union> <identifier>

<struct-or-union> ::= struct
                  | union

<struct-declaration> ::= {<specifier-qualifier>}* <struct-declarator-list>
```

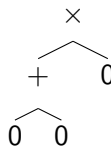
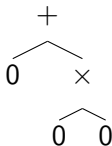
En pratique :

- La grammaire est **non-contextuelle** (ou algébrique) \rightarrow reconnaissable par automate à pile non déterministes, **algorithme CYK**, en temps $\mathcal{O}(n^3 \cdot |G|)$.
- En pratique, on se restreint aux langages reconnaissables par automates à pile **déterministes**, c'est-à-dire que la grammaire ne contient pas d'ambiguïtés. Ce sont les langages **LR(1)**, reconnaissables en temps $\mathcal{O}(n \cdot |G|)$.

Exemple d'ambiguïté (reduce/reduce) :

$expr$	$::=$	$integer$	scalar value
		$expr + expr$	binary operator
		$expr * expr$	binary operator

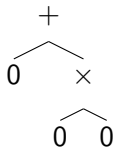
Parsons l'expression $0+0*0$:



Solution:

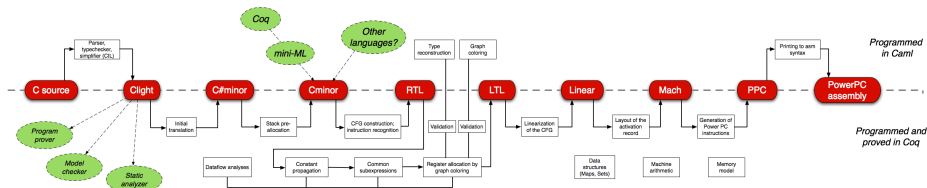
$$\begin{aligned} \text{expr} &::= \text{expr} + \text{texpr} \\ &\quad | \text{texpr} \\ \text{texpr} &::= \text{texpr} * \text{fexpr} \\ &\quad | \text{fexpr} \\ \text{fexpr} &::= \text{integer} \end{aligned}$$

Il n'y a plus d'ambiguïté :



- Par ailleurs, le compilateur génère des erreurs si des violations sont repérées (erreurs lexicales, erreurs syntaxiques, ...)
- Un lexer ne retourne **jamais** en arrière : il construit à chaque fois le plus long lexème valide à partir de la position courante.
- D'autres analyses (*sémantiques*) peuvent être effectuées : **typage statique**, **vérification des noms** (noms de variables uniques, ...), **data analysis**, **control-flow analysis**

Phase de synthèse



Le but est de transformer l'arbre de syntaxe abstraite en code objet (assembleur). Plusieurs étapes :

- Sélection d'instructions / Optimisations (propagation de constantes, élimination de code mort, factorisation de code, ...)
- Register transfer language (RTL) / Single static assignment (SSA) / Continuation passing style (CPS)
- Location transfer language (LTL)
- Assembleur (code linéarisé)

Un peu d'architecture matérielle

Computer Organization and Design: the HW/SW interface, Patterson et Hennessy

Rappel : les nombres sont représentés en binaires dans une machine.

1	0	0	1	0	1	1	0	0x96
---	---	---	---	---	---	---	---	------

Plusieurs représentations :

- Entier non signé :

$$1 \times 128 + 0 \times 64 + 1 \times 32 + 0 \times 16 + 0 \times 8 + 1 \times 4 + 1 \times 2 + 0 \times 1 = 166$$

- Entier signé :

$$1 \times (-128) + 0 \times 64 + 1 \times 32 + 0 \times 16 + 0 \times 8 + 1 \times 4 + 1 \times 2 + 0 \times 1 = -90$$

Un peu d'architecture matérielle

Mais pour un même nombre 0x96134f6f (en 32 bits, 4 octets) :

x+0	x+1	x+2	x+3
0x96	0x13	0x4f	0x6f

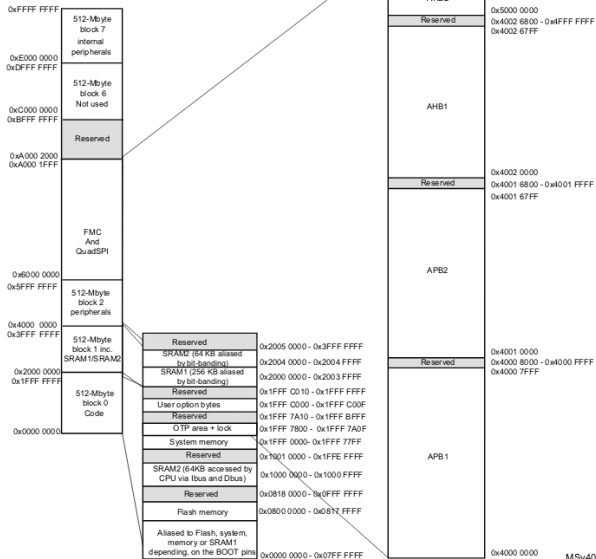
Figure : Big-endian

x+0	x+1	x+2	x+3
0x6f	0x4f	0x13	0x96

Figure : Little-endian

Un ordinateur contient deux types de mémoires :

- Les registres : ils contiennent exactement un mot, et sont localisés dans le processeur. Très rapides d'accès, mais leur nombre est restreint. On les désigne par leur nom de registre.
- La mémoire volatile : elle peut contenir beaucoup de données, mais avec un temps d'accès plus long. On y accède grâce à une adresse (pointeur en C). Les adresses ont une taille de 1 mot (i.e. un registre).
- Autres : ROM, périphériques, ... Leur accès se fait aussi grâce à une adresse.



Une architecture définit le nombre de registres et leur taille :

- x86:32 : 4 general purpose registers (eax, ebx, ecx, edx), 6 segment registers (cs, ds, es, fs, gs, ss)
- rv64gc : 31 general purpose registers de 64 bits (x1-x31) + 2 registres spéciaux (x0 et pc) + ...
- aarch64 : 31 general purpose registres de 64 bits (x0-x30) + 4 registres spéciaux de 64 bits (zr, sp, pc, elr) + ...

Revenons aux registres

Une architecture définit le nombre de registres et leur taille :

- x86:32 : 4 general purpose registers (eax, ebx, ecx, edx), 6 segment registers (cs, ds, es, fs, gs, ss)
- rv64gc : 31 general purpose registers de 64 bits (x1-x31) + 2 registres spéciaux (x0 et pc) + ...
- aarch64 : 31 general purpose registres de 64 bits (x0-x30) + 4 registres spéciaux de 64 bits (zr, sp, pc, elr) + ...

Mais aussi registres un peu spéciaux : les flags

- ARM CPSR : NZCV qui renseignent sur l'opération précédente
- ...

En réalité c'est plus compliqué que cela :

- Opérations flottantes
- Extensions au jeu d'instruction (SSE, AVX, ...)
- Caches
- Mémoire virtualisée
- Niveaux de privilège

Et les instructions dans tout ça ?

Plusieurs types d'instructions en assembleur :

- Opérations arithmétiques et logiques : addition, soustraction, ...
- Opérations mémoire : load, store.
- Opérations de flot de contrôle : unconditional and conditional branches, function calls.

Les instructions sont stockées en mémoire, et leur longueur/format est spécifié dans l'architecture.

Le pointeur vers le code se nomme le **program counter** (ou **pc**).

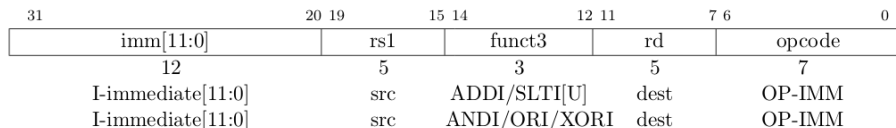
31	30	25	24	21	20	19	15	14	12	11	8	7	6	0				
funct7				rs2			rs1		funct3		rd			opcode		R-type		
imm[11:0]						rs1		funct3		rd			opcode		I-type			
imm[11:5]				rs2			rs1		funct3		imm[4:0]			opcode		S-type		
imm[12]		imm[10:5]			rs2			rs1		funct3		imm[4:1]		imm[11]		opcode		B-type
imm[31:12]										rd			opcode		U-type			
imm[20]		imm[10:1]			imm[11]		imm[19:12]				rd			opcode		J-type		

Opérations

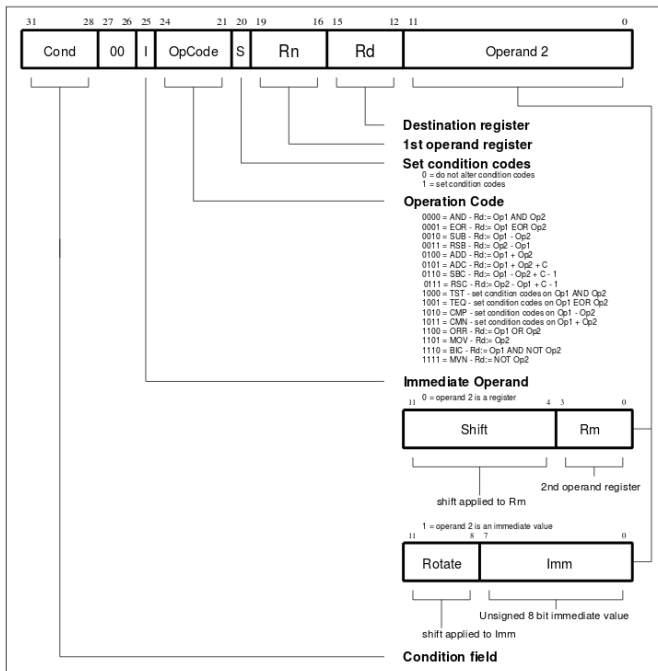
RV32I Base Instruction Set

imm[31:12]				rd	0110111	LUI
imm[31:12]				rd	0010111	AUIPC
imm[20 10:1 11 19:12]				rd	1101111	JAL
imm[11:0]		rs1	000	rd	1100111	JALR
imm[12 10:5]	rs2	rs1	000	imm[4:1 11]	1100011	BEQ
imm[12 10:5]	rs2	rs1	001	imm[4:1 11]	1100011	BNE
imm[12 10:5]	rs2	rs1	100	imm[4:1 11]	1100011	BLT
imm[12 10:5]	rs2	rs1	101	imm[4:1 11]	1100011	BGE
imm[12 10:5]	rs2	rs1	110	imm[4:1 11]	1100011	BLTU
imm[12 10:5]	rs2	rs1	111	imm[4:1 11]	1100011	BGEU
imm[11:0]		rs1	000	rd	0000011	LB
imm[11:0]		rs1	001	rd	0000011	LH
imm[11:0]		rs1	010	rd	0000011	LW
imm[11:0]		rs1	100	rd	0000011	LBU
imm[11:0]		rs1	101	rd	0000011	LHU
imm[11:5]	rs2	rs1	000	imm[4:0]	0100011	SB
imm[11:5]	rs2	rs1	001	imm[4:0]	0100011	SH
imm[11:5]	rs2	rs1	010	imm[4:0]	0100011	SW
imm[11:0]		rs1	000	rd	0010011	ADDI
imm[11:0]		rs1	010	rd	0010011	SLTI
imm[11:0]		rs1	011	rd	0010011	SLTIU
imm[11:0]		rs1	100	rd	0010011	XORI

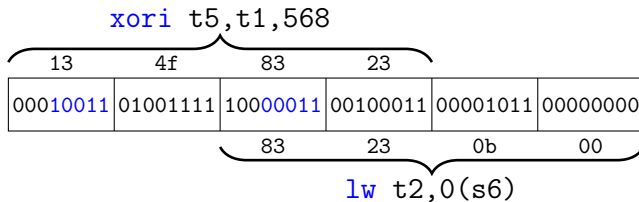
Regardons par exemple `addi` en risc-v (I-type)



ADDI adds the sign-extended 12-bit immediate to register *rs1*. Arithmetic overflow is ignored and the result is simply the low XLEN bits of the result. `ADDI rd, rs1, 0` is used to implement the `MV rd, rs1` assembler pseudo-instruction.



Overlapping code



Première abstraction : les fonctions

- But du jeu : réutiliser des morceaux de code.
- Besoin d'une interface qui permette d'en contrôler les effets de bord.
- Besoin de généricité : le même code doit fonctionner avec plusieurs données.
- Compilateurs différents, langages différents, ...
- Représentation mémoire des structures de données.

⇒ **Application binary interface (ABI)**

La fonction n'a pas le droit de faire n'importe quoi avec la mémoire ou les registres.

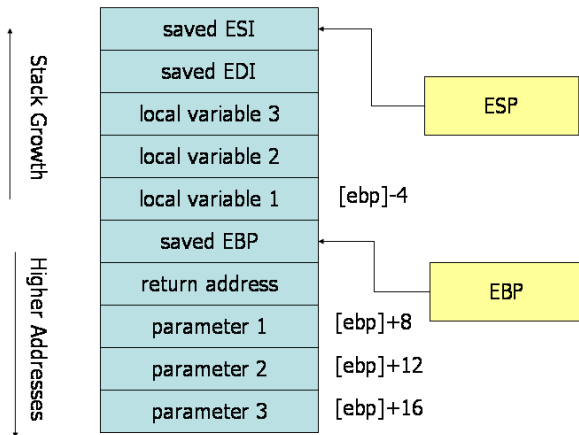
Premiers pas dans le reverse : l'ABI

Register	ABI Mnemonic	Meaning
x0	zero	Zero
x1	ra	Return address
x2	sp	Stack pointer
x3	gp	Global pointer
x4	tp	Thread pointer
x5-x7	t0-t2	Temporary registers
x8-x9	s0-s1	Callee-saved registers
x10-x17	a0-a7	Argument registers
x18-x27	s2-s11	Callee-saved registers
x28-x31	t3-t6	Temporary registers

f0-f7	ft0-ft7	Temporary registers
f8-f9	fs0-fs1	Callee-saved registers
f10-f17	fa0-fa7	Argument registers
f18-f27	fs2-fs11	Callee-saved registers
f28-f31	ft8-ft11	Temporary registers

Figure : Naming convention for registers, per RISC-V ELF psABI.

ABI : Pile d'appel



De façon plus générale

L'ABI définit une convention d'appel :

- Où se trouvent les paramètres ? (registres, pile, ...)
- Quels sont les registres dont la valeur reste inchangée ?
- Alignements des paramètres et de la pile ?

Exemples :

- <https://github.com/riscv/riscv-elf-psabi-doc/blob/master/riscv-elf.md>
- https://www.agner.org/optimize/calling_conventions.pdf
- web.archive.org/web/20080410091312/http://msdn2.microsoft.com/en-us/library/zxk0tw93%28vs.71%29.aspx
- <https://www.ohse.de/uwe/articles/gcc-attributes.html#func-fastcall>

Et dans gcc

```
int __attribute__((stdcall)) function(int a, int b);

int __attribute__((fastcall)) function(int a, int b);

int __attribute__((cdecl)) function(int a, int b);

int __attribute__((sysv_abi)) function(int a, int b);

int __attribute__((no_caller_saved_registers))
    function(int a, int b);

int __attribute__((thiscall)) function(int a, int b);
```


En pratique : exemple d'une fonction

```
addi    sp,sp,-16
sd      ra,8(sp)

jal     ra,dummy
lui     a0,0x9932
lui     a3,0x23371
lui     a2,0xa0212
mv      a1,zero
jal     ra,dummy4

ld      ra,8(sp)
mv      a0,zero
addi    sp,sp,16
ret
```

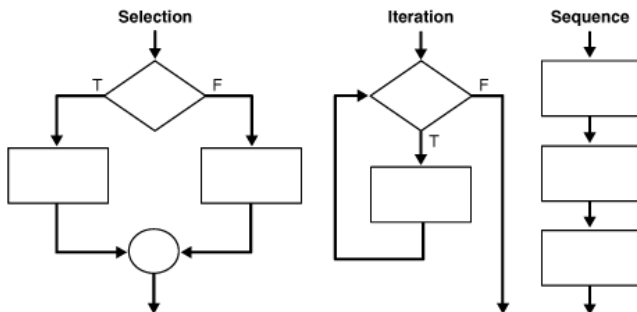
Pour en revenir à la compilation

La traduction de l'arbre de syntaxe abstraite peut se faire quasi-automatiquement vers un langage qui n'est pas tout à fait de l'assembleur :

- $expr1 + expr2 \Rightarrow x = expr1; y = expr2; \text{ADD } z, y, x;$
- $*expr1 \Rightarrow x = expr1; \text{LW } y, 0(x)$
- $fun(expr1) \Rightarrow x = expr1; \text{PUSH } x; \text{CALL } fun$

Programmation structurée (deuxième abstraction)

Il manque les instructions de flots de controle:



Ce qu'il reste à faire

- Prendre en compte l'analyse de durée de vie (defuse chain).
- Allocation des registres (coloriage d'un graphe d'interférence).
- Optimisation des instructions en trop.

Des questions ?

Programme de cet après midi : rétro-ingénierie statique, merci d'avoir sur votre machine :

- objdump/gcc/readelf/strings
- Ghidra 9.1

Pour approfondir/bibliographie :

<https://www.lri.fr/~filliatr/ens/compil/>

https://perso.telecom-paristech.fr/guilley/ENS/program_2019_2020.html/