

Université Paris-Sud

# Compilation et Langages

Sylvain Conchon

Cours 1 / 14 janvier 2016

- **Cours** le jeudi, 8h30–10h30 dans le grand amphi (PUIO)
  - pas de polycopié, mais transparents disponibles
- **TD**
  - le jeudi, 10h45–12h45
- **MCC**
  - 1ère session : 0.5 CC1 + 0.5 EX1
  - 2ème session : 0.5 CC1 + 0.5 EX2

Toutes les infos sur le site web du cours

<http://www.lri.fr/~conchon/compilation/>

2

## Cette semaine

**Cours** de **compilation** : 8h30 - 10h30 (Grand Amphi)

**Cours** de **remise à niveau en OCaml** : 10h45 - 12h45 (Grand Amphi)

**TP** de **remise à niveau en OCaml** : 14h - 17 h (salles E202 et E205)

## Remerciements

à **Jean-Christophe Filliâtre** pour le matériel de son cours de compilation  
à l'ENS Ulm

Maîtriser les mécanismes de la **compilation**, c'est-à-dire de la transformation d'un langage dans un autre

Comprendre les différents aspects des **langages de programmation** par le biais de la compilation

... utiles pour concevoir des outils qui manipulent les programmes de manière symbolique, comme les outils de

- **preuve de programmes** (VCC, Dafny, Spec#, Frama-C, Spark, GNATProve, Why3, Boogie, etc.)
- **vérification par model checking** (Slam, Spin, CBMC, Murphi, Cubicle, Blast, Uppaal, Java Pathfinder, etc.)
- **analyse par interprétation abstraite** (Astrée, polyspace, etc.)
- **démonstration automatique** (z3, cvc4, Alt-Ergo, etc.),
- **test formel** (Pex, PathCrawler, etc.)
- etc.

Toutes ces thématiques seront abordées dans le **M2 FIIL** (Fondements de l'Informatique et Ingénierie du Logiciel)

5

6

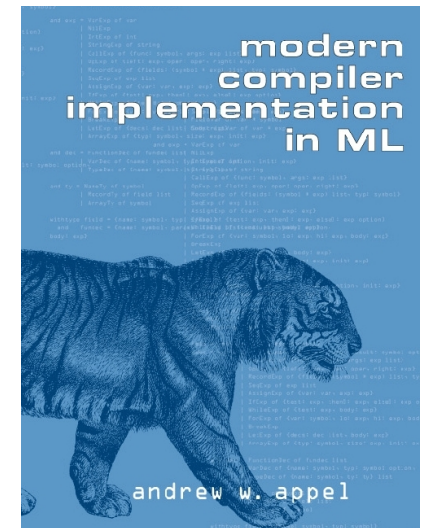
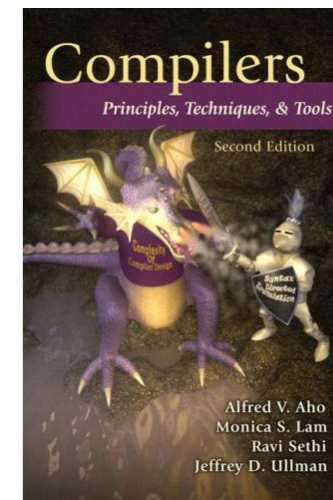
## Programmation

## Un peu de lecture

ici on programme

- en cours
- en TD/TP
- à l'examen

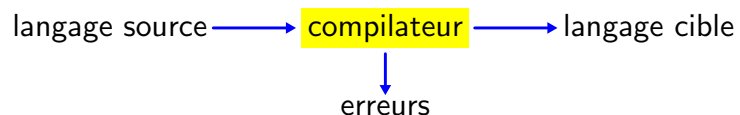
on programme en **Objective Caml**



7

8

Schématiquement, un compilateur est un programme qui traduit un « programme » d'un langage **source** vers un langage **cible**, en signalant d'éventuelles erreurs



Quand on parle de compilation, on pense typiquement à la traduction d'un langage de haut niveau (C, Java, Caml, ...) vers le langage machine d'un processeur (Intel Pentium, PowerPC, ...)

```
% gcc -o sum sum.c
```



```
int main(int argc, char **argv) {
    int i, s = 0;
    for (i = 0; i <= 100; i++) s += i*i;
    printf("0*0+...+100*100 = %d\n", s);
}
```

```
001001111011110111111111111100000
101011111011111110000000000010100
101011111010010000000000001000000
101011111010010100000000000100100
101011111010000000000000000110000
101011111010000000000000000111000
100011111010111100000000000111000
...
```

9

10

Dans ce cours, nous allons effectivement nous intéresser à la compilation vers de **l'assembleur**, mais ce n'est qu'un aspect de la compilation

Un certain nombre de techniques mises en œuvre dans la compilation ne sont pas liées à la production de code assembleur

Certains langages sont d'ailleurs

- interprétés (Basic, COBOL, Ruby, etc.)
- compilés dans un langage intermédiaire qui est ensuite interprété (Java, Caml, etc.)
- compilés vers un autre langage de haut niveau
- compilés à la volée

un **compilateur** traduit un programme  $P$  en un programme  $Q$  tel que pour toute entrée  $x$ , la sortie de  $Q(x)$  soit la même que celle de  $P(x)$

un **interprète** est un programme qui, étant donné un programme  $P$  et une entrée  $x$ , calcule la sortie  $s$  de  $P(x)$

Dit autrement,

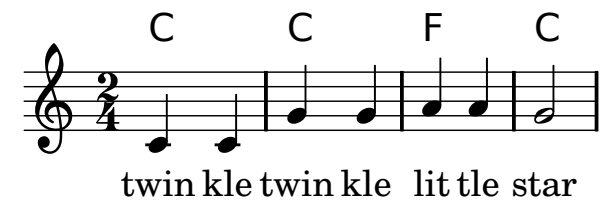
Le compilateur fait un travail complexe **une seule fois**, pour produire un code fonctionnant pour n'importe quelle entrée

L'interprète effectue un travail plus simple, mais le refait sur chaque entrée

Autre différence : le code compilé est généralement bien plus efficace que le code interprété

source → **lilypond** → fichier PostScript → **gs** → image

```
<<
  \chords { c2 c f2 c }
  \new Staff \relative c' { \time 2/4 c4 c g'4 g a4 a g2 }
  \new Lyrics \lyricmode { twin4 kle twin kle lit tle star2 }
>>
```



13

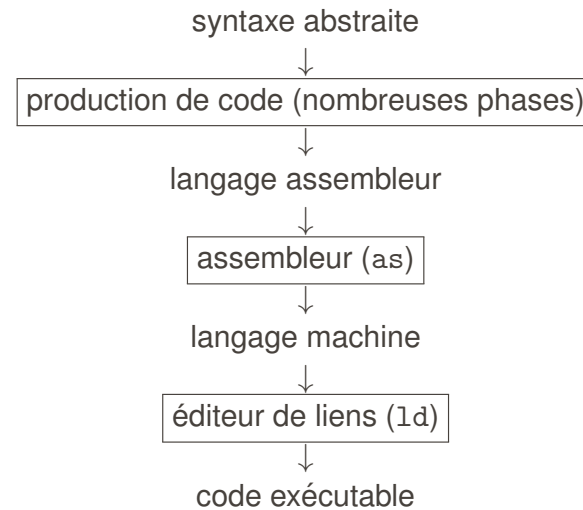
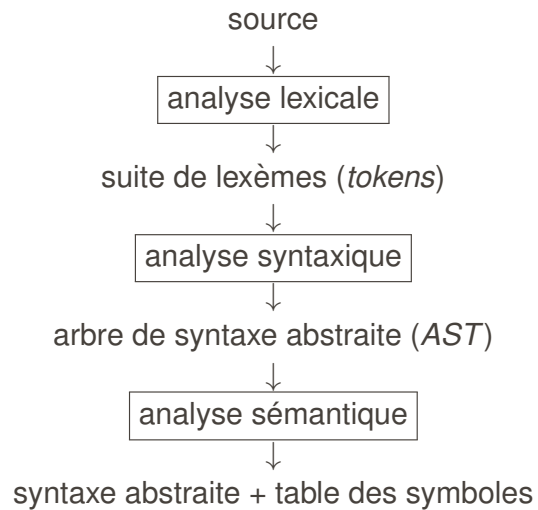
14

À quoi juge-t-on la qualité d'un compilateur ?

- à sa correction
- à l'efficacité du code qu'il produit
- à sa propre efficacité

Typiquement, le travail d'un compilateur se compose

- d'une phase d'**analyse**
  - reconnaît le programme à traduire et sa signification
  - signale les erreurs et peut donc échouer (erreurs de syntaxe, de portée, de typage, etc.)
- puis d'une phase de **synthèse**
  - production du langage cible
  - utilise de nombreux langages intermédiaires
  - n'échoue pas



17

18

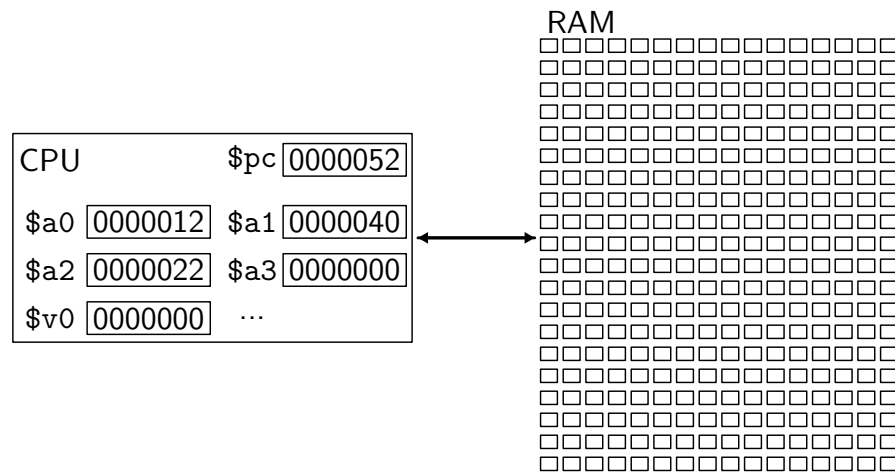
## Un peu d'architecture

### L'assembleur MIPS

(voir la documentation sur la page du cours)

Très schématiquement, un ordinateur est composé

- d'une unité de calcul (CPU), contenant
  - un petit nombre de registres entiers ou flottants
  - des capacités de calcul
- d'une mémoire vive (RAM)
  - composée d'un très grand nombre d'octets (8 bits)  
par exemple, 1 Go =  $2^{30}$  octets =  $2^{33}$  bits, soit  $2^{2^{33}}$  états possibles
  - contient des données et des instructions



L'accès à la mémoire coûte cher (à un milliard d'instructions par seconde, la lumière ne parcourt que 30 centimètres entre 2 instructions !)

21

La réalité est bien plus complexe

- plusieurs (co)processeurs, dont certains dédiés aux flottants
- un ou plusieurs caches
- une virtualisation de la mémoire (MMU)
- etc.

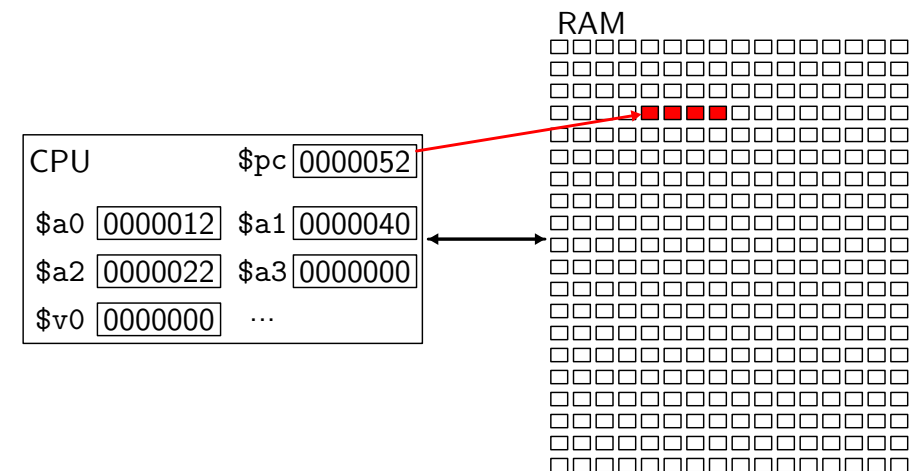
22

## Principe d'exécution

L'exécution d'un programme se déroule ainsi

- un registre (\$pc) contient l'adresse de l'instruction à exécuter
- on lit les 4 (ou 8) octets à cette adresse (*fetch*)
- on interprète ces bits comme une instruction (*decode*)
- on exécute l'instruction (*execute*)
- on modifie le registre \$pc pour passer à l'instruction suivante (typiquement celle se trouvant juste après, sauf en cas de saut)

## Principe d'exécution



instruction : 000000 00001 00010 0000000000001010  
 décodage : add \$a1 \$a2 10

i.e. ajouter 10 au registre \$a1 et stocker le résultat dans le registre \$a2

23

24

Là encore la réalité est bien plus complexe

- pipelines
  - plusieurs instructions sont exécutées en parallèle
- prédiction de branchement
  - pour optimiser le pipeline, on tente de prédire les sauts conditionnels

Deux grandes familles de microprocesseurs

- CISC (*Complex Instruction Set*)
  - beaucoup d'instructions
  - beaucoup de modes d'adressage
  - beaucoup d'instructions lisent / écrivent en mémoire
  - peu de registres
  - exemples : VAX, PDP-11, Motorola 68xxx, Intel x86
- RISC (*Reduced Instruction Set*)
  - peu d'instructions, régulières
  - très peu d'instructions lisent / écrivent en mémoire
  - beaucoup de registres, uniformes
  - exemples : Alpha, Sparc, MIPS, ARM

on choisit **MIPS** pour ce cours (et les TD/TP)

25

26

## L'architecture MIPS

- 32 registres, r0 à r31
  - r0 contient toujours 0
  - utilisables sous d'autres noms, correspondant à des conventions (zero, at, v0–v1, a0–a3, t0–t9, s0–s7, k0–k1, gp, sp, fp, ra)
- trois types d'instructions
  - instructions de transfert, entre registres et mémoire
  - instructions de calcul
  - instructions de saut

documentation : sur le site du cours

## Simulateurs MIPS

En pratique, on utilisera un simulateur MIPS, **MARS** (ou **SPIM**)

En ligne de commande

- `java -jar Mars_4_2.jar file.s`

En mode graphique et interactif

- `java -jar Mars_4_2.jar`
- charger le fichier et l'assembler
- mode pas à pas, visualisation des registres, de la mémoire, etc.

Documentation : sur le site du cours

27

28

- chargement d'une constante (16 bits signée) dans un registre

```
li    $a0, 42    # a0 <- 42
lui   $a0, 42    # a0 <- 42 * 2^16
```

- copie d'un registre dans un autre

```
move  $a0, $a1   # copie a1 dans a0 !
```

29

- addition de deux registres

```
add   $a0, $a1, $a2 # a0 <- a1 + a2
add   $a2, $a2, $t5  # a2 <- a2 + t5
```

de même, sub, mul, div

- addition d'un registre et d'une constante

```
addi  $a0, $a1, 42  # a0 <- a1 + 42
```

(mais pas subi, muli ou divi !)

- négation

```
neg   $a0, $a1      # a0 <- -a1
```

- valeur absolue

```
abs   $a0, $a1      # a0 <- |a1|
```

30

- NON logique ( $\text{not}(100111_2) = 011000_2$ )

```
not   $a0, $a1      # a0 <- not(a1)
```

- ET logique ( $\text{and}(100111_2, 101001_2) = 100001_2$ )

```
and   $a0, $a1, $a2 # a0 <- and(a1, a2)
andi  $a0, $a1, 0x3f # a0 <- and(a1, 0...0111111)
```

- OU logique ( $\text{or}(100111_2, 101001_2) = 101111_2$ )

```
or    $a0, $a1, $a2 # a0 <- or(a1, a2)
ori   $a0, $a1, 42  # a0 <- or(a1, 0...0101010)
```

31

- décalage à gauche (insertion de zéros)

```
sll   $a0, $a1, 2    # a0 <- a1 * 4
sllv  $a1, $a2, $a3  # a1 <- a2 * 2^a3
```

- décalage à droite arithmétique (copie du bit de signe)

```
sra   $a0, $a1, 2    # a0 <- a1 / 4
```

- décalage à droite logique (insertion de zéros)

```
srl   $a0, $a1, 2
```

- rotation

```
rol   $a0, $a1, 2
ror   $a0, $a1, 3
```

32



- comparaison de deux registres

```
slt  $a0, $a1, $a2    # a0 <- 1 si a1 < a2
                        #      0 sinon
```

ou d'un registre et d'une constante

```
slti $a0, $a1, 42
```

- variantes : sltu (comparaison non signée), sltiu
- de même : sle, sleu / sgt, sgtu / sge, sgeu
- égalité : seq, sne

33

- lire un mot (32 bits) en mémoire

```
lw   $a0, 42($a1)    # a0 <- mem[a1 + 42]
```

l'adresse est donnée par un registre et un décalage sur 16 bits signés

- variantes pour lire 8 ou 16 bits, signés ou non (lb, lh, lbu, lhu)

34

- écrire un mot (32 bits) en mémoire

```
sw   $a0, 42($a1)    # mem[a1 + 42] <- a0
                        # attention au sens !
```

l'adresse est donnée par un registre et un décalage sur 16 bits signés

- variantes pour écrire 8 ou 16 bits (sb, sh)

On distingue

- branchement** : typiquement un saut conditionnel, dont le déplacement est stocké sur 16 bits signés (-32768 à 32767 instructions)
- saut** : saut inconditionnel, dont l'adresse de destination est stockée sur 26 bits

35

36

- branchement conditionnel

```
beq $a0, $a1, label # si a0 = a1 saute à label
                        # ne fait rien sinon
```

- variantes : bne, blt, ble, bgt, bge (et comparaisons non signées)
- variantes : beqz, bnez, bgez, bgtz, bltz, blez

37

### Saut inconditionnel

- à une adresse (*jump*)

```
j label
```

- avec sauvegarde de l'adresse de l'instruction suivante dans \$ra

```
jal label # jump and link
```

- à une adresse contenue dans un registre

```
jr $a0
```

- avec l'adresse contenue dans \$a0 et sauvegarde dans \$a1

```
jalr $a0, $a1
```

38

Quelques appels système fournis par une instruction spéciale

```
syscall
```

Le code de l'instruction doit être dans \$v0, les arguments dans \$a0–\$a3 ;

Le résultat éventuel sera placé dans \$v0

Exemple : appel système `print_int` pour afficher un entier

```
li $v0, 1 # code de print_int
li $a0, 42 # valeur à afficher
syscall
```

de même `read_int`, `print_string`, etc. (voir la documentation)

39

On ne programme pas en langage machine mais en assembleur

L'assembleur fourni un certain nombre de facilités :

- étiquettes symboliques
- allocation de données globales
- pseudo-instructions

40

La directive

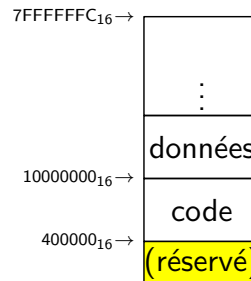
```
.text
```

indique que des instructions suivent, et la directive

```
.data
```

indique que des données suivent

Le code sera chargé à partir de l'adresse 0x400000  
et les données à partir de l'adresse 0x10000000



Unec étiquette symbolique est introduite par

```
label:
```

et l'adresse qu'elle représente peut être chargée dans un registre

```
la $a0, label
```

41

```
.text
main: li $v0, 4    # code de print_string
      la $a0, hw    # adresse de la chaîne
      syscall      # appel système
      li $v0, 10    # exit
      syscall
.data
hw: .asciiz "hello world\n"
```

(.asciiz est une facilité pour .byte 104, 101, ... 0)

42

## Le défi de la compilation

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

En particulier, il faut

- traduire les structures de contrôle (tests, boucles, exceptions, etc.)
- traduire les appels de fonctions
- traduire les structures de données complexes (tableaux, enregistrements, objets, clôtures, etc.)
- allouer de la mémoire dynamiquement

43

## Appels de fonctions

**Constat :** les appels de fonctions peuvent être arbitrairement imbriqués  
⇒ les registres ne peuvent suffire pour les paramètres / variables locales  
⇒ il faut allouer de la mémoire pour cela

les fonctions procèdent selon un mode *last-in first-out*, c'est-à-dire de **pile**

44

pile
↓
↑
données dynamiques (tas)
données statiques
code

La **pile** est stockée tout en haut, et croît dans le sens des adresses décroissantes ; `$sp` pointe sur le sommet de la pile

les données dynamiques (survivant aux appels de fonctions) sont allouées sur le **tas** (éventuellement par un GC), en bas de la zone de données, juste au dessus des données statiques

ainsi, on ne se marche pas sur les pieds

lorsqu'une fonction  $f$  (l'appelant ou *caller*) souhaite appeler une fonction  $g$  (l'appelé ou *callee*), elle exécute

```
jal g
```

et lorsque l'appelé en a terminé, il lui rend le contrôle avec

```
jr $ra
```

problème :

- si  $g$  appelle elle-même une fonction, `$ra` sera écrasé
- de même, tout registre utilisé par  $g$  sera perdu pour  $f$

il existe de multiples manières de s'en sortir,

mais en général on s'accorde sur des **conventions d'appel**

45

46

## Conventions d'appel

utilisation des registres

- `$at`, `$k0` et `$k1` sont réservés à l'assembleur et l'OS
- `$a0–$a3` sont utilisés pour passer les quatre premiers arguments (les autres sont passés sur la pile) et `$v0–$v1` pour renvoyer le résultat
- `$t0–$t9` sont des registres **caller-saved** i.e. l'appelant doit les sauvegarder si besoin ; on y met donc typiquement des données qui n'ont pas besoin de survivre aux appels
- `$s0–$s7` sont des registres **callee-saved** i.e. l'appelé doit les sauvegarder ; on y met donc des données de durée de vie longue, ayant besoin de survivre aux appels
- `$sp` est le pointeur de pile, `$fp` le pointeur de *frame*
- `$ra` contient l'adresse de retour

47

## L'appel, en quatre temps

il y a quatre temps dans un appel de fonction

- 1 pour l'appelant, juste avant l'appel
- 2 pour l'appelé, au début de l'appel
- 3 pour l'appelé, à la fin de l'appel
- 4 pour l'appelant, juste après l'appel

s'organisent autour d'un segment situé au sommet de la pile appelé le **tableau d'activation**, en anglais **stack frame**, situé entre `$fp` et `$sp`

48

- ❶ passe les arguments dans \$a0–\$a3, les autres sur la pile s'il y en a plus de 4
- ❷ sauvegarde les registres \$t0–\$t9 qu'il compte utiliser après l'appel (dans son propre tableau d'activation)
- ❸ exécute

```
jal appelé
```

- ❶ alloue son tableau d'activation, par exemple

```
addi $sp, $sp, -28
```

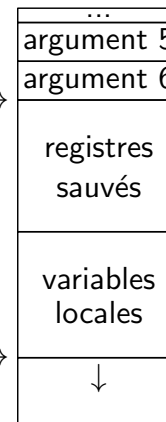
- ❷ sauvegarde \$fp puis le positionne, par exemple

```
sw    $fp, 24($sp)
addi  $fp, $sp, 24
```

- ❸ sauvegarde \$s0–\$s7 et \$ra si besoin

\$fp→

\$sp→



\$fp permet d'atteindre facilement les arguments et variables locales, avec un décalage fixe quel que soit l'état de la pile

49

50

- ❶ place le résultat dans \$v0 (voire \$v1)
- ❷ restaure les registres sauvegardés
- ❸ dépile son tableau d'activation, par exemple

```
addi $sp, $sp, 28
```

- ❹ exécute

```
jr    $ra
```

- ❶ dépile les éventuels arguments 5, 6, ...
- ❷ restaure les registres caller-saved

51

52

**exercice** : programmer la fonction factorielle

- une machine fournit
  - un jeu limité d'instructions, très primitives
  - des registres efficaces, un accès coûteux à la mémoire
- la mémoire est découpée en
  - code / données statiques / tas (données dynamiques) / pile
- les appels de fonctions s'articulent autour
  - d'une notion de tableau d'activation
  - de conventions d'appel