

# Génération de code-I

## Constructions élémentaires

28 septembre 2011

- 1 Introduction
  - Modèle d'exécution
- 2 Génération de code utilisant une pile
  - Modèle d'exécution à pile
  - Expressions arithmétiques simples
  - Variables
  - Instructions conditionnelles

# Où en est-on ?

- L'analyse **lexicale et syntaxique** permet de construire un premier **arbre de syntaxe abstraite** et d'éliminer des programmes incorrects.
- L'analyse **sémantique** travaille par récurrence sur l'arbre de syntaxe abstraite pour calculer de nouvelles informations (portée, types, surcharge . . . ) qui sont conservées dans l'arbre ou ailleurs. De nouvelles erreurs sont détectées.
- D'autres représentations comme le **graphe de flots** collecteront des informations liées à l'exécution du programme.
- La **génération de code** produit un nouveau code exécutable en machine.

## Génération de code pour de l'assembleur MIPS

- *Introduction*: modèle d'exécution à l'aide d'une pile
- *Code intermédiaire: modèle à pile*
  - Instructions;
  - Cas de base;
  - Compilation des expressions conditionnelles.
- *Appel de procédures*
- *Langages objets*
- *Langages fonctionnels*

## 1 Introduction

- Modèle d'exécution

## 2 Génération de code utilisant une pile

- Modèle d'exécution à pile
- Expressions arithmétiques simples
- Variables
- Instructions conditionnelles

- la taille du code est connue à la compilation
- une partie de la mémoire est réservée aux instructions du programme
- un pointeur (*pc*) indique où on en est dans l'exécution
- le code s'exécute de manière séquentielle :
  - sauf instruction explicite de saut, les instructions du programme sont exécutées l'une après l'autre.

- Les données du programme sont stockées soit dans la **mémoire** soit dans les **registres** de la machine.
- La mémoire est organisée en **mots** (32 ou 64 bits), on y accède par une adresse qui est représentée par un entier.
- Les valeurs **simples** sont stockées dans une unité de la mémoire.
- Les valeurs **complexes** sont stockées dans des cases consécutives de la mémoire (structures, tableaux), ou dans des structures chaînées (une liste est représentée par la valeur de son premier élément associée à l'adresse du reste de la liste).
- Les registres permettent d'accéder **rapidement** à des données simples.
- Le **nombre de registres** dépend de l'architecture de la machine et est limité.

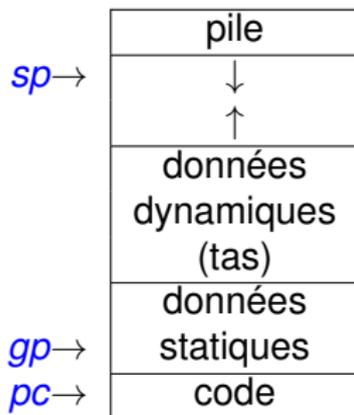
- Certaines données sont explicitement manipulées par le programme source par l'intermédiaire de **variables**.
- D'autres données sont **créées** par le compilateur :
  - valeurs intermédiaires lors de l'exécution d'un calcul arithmétique
  - variables lors de l'appel de fonctions
  - valeurs objet, valeurs fonctionnelles ...
- Certaines données ont une **durée de vie** connue à la compilation: règles de portées du langage, analyse des variables vivantes.

- Les variables **globales** du programme peuvent être allouées à des **adresses fixes**, à condition de connaître leur taille.
- La gestion des variables locales des blocs et des procédures, ou bien le stockage de valeurs intermédiaires peut se faire en allouant de la mémoire en **pile (stack)**.
- D'autres données ont par contre une durée de vie qui n'est pas connue à la compilation.  
C'est le cas lorsque l'on manipule des pointeurs (expressions dont la valeur est une adresse).  
Ces données vont être allouées en général dans une autre partie de la mémoire, organisée en **tas (heap)**.

L'espace réservé dans le tas devra être libéré s'il n'est plus utilisé.

- commandes explicites pour libérer l'espace (C, Pascal)
- l'exécution du programme utilise un programme général de récupération de mémoire appelé communément **gc** pour **garbage collector** qui est traduit en **glaneur de cellules** ou **ramasse-miettes** (Caml, Java).

# Schéma d'organisation de la mémoire



Code **indépendant** de la machine cible.

- factoriser une grande partie du travail de compilation;
- rendre le compilateur plus portable.

Le choix d'un langage intermédiaire est très important.

- assez riche pour permettre un codage aisé des opérations du langage sans créer de trop longues séquences de code
- assez limité pour que l'implantation finale ne soit pas trop coûteuse.

- Nombre de **registres** limité, indispensables pour les calculs.
- Nécessité de **sauvegarder** certaines valeurs dans une pile.
- Le passage du langage de haut niveau vers de l'assembleur passe par plusieurs **langages intermédiaires**.
- L'assembleur lui-même utilise des **pseudo-instructions** de plus haut niveau que le code machine et des étiquettes symboliques.

## 1 Introduction

- Modèle d'exécution

## 2 Génération de code utilisant une pile

- **Modèle d'exécution à pile**
- Expressions arithmétiques simples
- Variables
- Instructions conditionnelles

Nous spécifierons une fonction *code* qui prend comme argument un arbre de syntaxe abstraite pour le langage et renvoie une suite d'instructions du code de la machine à pile.

Notations concrètes pour représenter la syntaxe abstraite:

- Si  $E ::= E_1 + E_2$  est une règle de grammaire nous écrivons  $code(E_1 + E_2) = \dots code(E_1) \dots code(E_2)$  spécifie la valeur de la fonction *code* sur un arbre de syntaxe abstraite correspondant à l'addition de deux expressions.
- Si  $C_1$  et  $C_2$  représentent des suites d'instructions alors  $C_1 | C_2$  représente la suite d'instructions obtenue en concaténant  $C_2$  à la suite de  $C_1$ .
- La liste vide d'instructions est représentée par  $[]$ .

- Stocker les valeurs intermédiaires dans la pile
- Echanger les valeurs entre mémoire et registres pour les calculs
- La pile peut contenir des entiers, des flottants ou des adresses.
- Certaines données grosses en taille comme les chaînes de caractères sont allouées dans un espace supplémentaire, lorsqu'elle devra manipuler une chaîne, la pile se contentera de manipuler l'adresse de la chaîne.

# Macros instructions pour la machine à pile

- un registre `$sp` (stack pointer) pointe sur la première cellule libre de la pile.
- L'adressage se fait par octet : un mot de 32 bits tient sur 4 octets.

Fonction de sauvegarde et de chargement sur la pile:

```
let pushr r = sub $sp, $sp, 4 |  
             sw r, 0($sp)
```

```
let pop r = lw r, 0($sp) |  
           add $sp, $sp, 4
```

## 1 Introduction

- Modèle d'exécution

## 2 Génération de code utilisant une pile

- Modèle d'exécution à pile
- **Expressions arithmétiques simples**
- Variables
- Instructions conditionnelles

# Schéma de compilation des expressions

Le résultat de l'expression est stocké dans le registre `$a0`.  
Les valeurs intermédiaires sont stockées dans la pile.

Le code correspondant est donné par :

```
code(entier) = li $a0 entier  
code( $E_1 + E_2$ ) = code( $E_1$ ) | pushr $a0 |  
                    code( $E_2$ ) | popr $t0 |  
                    add $a0 $t0 $a0
```

- Ce code est très inefficace...
- On peut l'améliorer en utilisant les registres temporaires `$ti` ( $i = 0 \dots 9$ ) et `$si` ( $i = 0 \dots 7$ )
- Peut se faire naïvement (en passant en argument une pile des registres disponibles) ou par des analyses fines qui seront vues plus tard.

La fonction de génération de code est récursive sur l'ast des expressions arithmétiques.

## Invariant

Exécution  $code(E)$  à partir d'un état où  $sp = n$ :

- la valeur de l'expression  $E$  est calculée dans le registre  $\$a0$
- les valeurs de  $P[m]$  pour  $m < n$  n'ont pas été modifiées

## 1 Introduction

- Modèle d'exécution

## 2 Génération de code utilisant une pile

- Modèle d'exécution à pile
- Expressions arithmétiques simples
- **Variables**
- Instructions conditionnelles

Les instructions **lw** et **sw** permettent de déplacer des valeurs entre une adresse mémoire et un registre.

- **sw**: affecter une valeur calculée dans un registre à un espace réservé pour les variables globales;
- **lw**: charger dans un registre une valeur contenue dans un variable globale.

**Notation** Si  $\$r$  est un registre dont la valeur est une adresse  $a$  dans la pile et  $n$  est une constante entière alors  $n(\$r)$  désigne la valeur à l'adresse  $a + n$ .

# Code pour les variables globales

$$E := \text{ident}$$

Une variable globale  $x$  pourra être stockée dans la zone de données associée à une étiquette symbolique  $\text{et}(x)$ .

```
.data
etiq_x:
    .word 0
.text
```

$$\text{code}(\text{id}) = \text{lw } \$a0 \text{ et}(\text{id})$$

On peut aussi repérer une variable (globale ou locale)  $x$  par son décalage  $\text{adr}(x)$  par rapport à un pointeur global  $\$gp$ .

$$\text{code}(\text{id}) = \text{lw } k(\$gp) \quad \text{si } k = \text{adr}(\text{id})$$

La place nécessaire aux variables est réservée avant l'exécution en retirant à  $\$sp$ , 4 fois le nombre de variables à stocker.

Un langage dont les seules instructions (non-terminal  $I$ ) sont des suites d'affectations.

## Règles de grammaire

$$\begin{aligned} I &::= \epsilon \\ I &::= I_1 A; \\ A &::= id := E \end{aligned}$$

## Génération de code correspondante

$$\begin{aligned} \text{code}(\epsilon) &= [] \\ \text{code}(I_1 A;) &= \text{code}(I_1) \mid \text{code}(A) \\ \text{code}(id := E) &= \text{code}(E) \mid \text{sw } \$a0 \text{ et}(id) && \text{si id global} \\ &\quad \text{code}(E) \mid \text{sw } \$a0 \text{ k}(\$gp) && \text{si k=adr(id)} \end{aligned}$$

## 1 Introduction

- Modèle d'exécution

## 2 Génération de code utilisant une pile

- Modèle d'exécution à pile
- Expressions arithmétiques simples
- Variables
- Instructions conditionnelles

# Instructions conditionnelles

- Les expressions conditionnelles et les boucles utilisent des commandes pour effectuer des sauts dans la zone d'instructions.
- Les instructions seront désignées par des adresses symboliques insérées dans le code:  
`label` : introduit un adresse symbolique `label` correspondant au numéro de l'instruction.
- branchement inconditionnel: `j label`
- branchement conditionnel si `$r = 0`: `beqz $r label`

Convention pour la représentation des booléens :

- `false` est représenté par 0
- `true` est représenté par 1 ou par n'importe quel entier positif.

$I ::= \text{if } E \text{ then } I \mid \text{if } E \text{ then } I_1 \text{ else } I_2$   
 $I ::= \text{while } E \text{ do } I \text{ done}$

*code*(if  $E$  then  $I$ )

= *code*( $E$ ) | **beqz** \$a0 new-suiv | *code*( $I$ ) | new-suiv:

*code*(if  $E$  then  $I_1$  else  $I_2$ )

= *code*( $E$ ) | **beqz** \$a0 new-faux | *code*( $I_1$ ) | **j** new-suiv  
| new-faux: *code*( $I_2$ ) | new-suiv:

*code*(while  $E$  do  $I$  done)

= new-loop: *code*( $E$ ) | **beqz** \$a0 new-suiv | *code*( $I$ )  
| **j** new-loop | new-suiv:

Nouvelles étiquettes new-faux, new-loop, new-suiv.

# Compilation d'expressions booléennes

$B ::= B_1 \text{ or } B_2 \mid B_1 \text{ and } B_2 \mid \text{not } B_1$

$B ::= \text{true} \mid \text{false}$

$B ::= E_1 \text{ relop } E_2$

Compiler les expressions booléennes comme des expressions arithmétiques. Empiler les valeurs 0 ou 1 en utilisant les opérations arithmétiques pour simuler les opérations booléennes.

`code(true)` = `li $a0 1`

`code(false)` = `li $a0 0`

`code( $B_1$  and  $B_2$ )` = `code( $B_1$ )` | `pushr $a0`  
| `code( $B_2$ )` | `popr $t0` | `mul $a0, $t0, $a0`

`code(not  $B_1$ )` = `code( $B_1$ )` | `neg $a0` | `add $a0, $a0, 1`

Utilisation de conditions.

$B_1$  **and**  $B_2$  = **if**  $B_1$  **then**  $B_2$  **else faux**  
 $B_1$  **or**  $B_2$  = **if**  $B_1$  **then vrai** **else**  $B_2$   
**not**  $B_1$  = **if**  $B_1$  **then faux** **else vrai**

Dans un langage avec effets de bord, le fait de calculer ou non les sous-expressions d'une expression booléenne peut avoir des comportements très différents.

- Les expressions booléennes servent souvent à des opérations de contrôle.
- On peut se donner a priori deux étiquettes  $E\text{-vrai}$  et  $E\text{-faux}$  et compiler l'expression booléenne sans calculer de valeur:
  - le résultat de l'exécution de  $E$  aboutit à  $pc = E\text{-vrai}$  lorsque la valeur de  $E$  est **vrai** et à  $E\text{-faux}$  sinon.
- On définit donc une fonction *code-bool* qui prend comme argument l'expression booléenne et les deux étiquettes et qui renvoie le code de contrôle.

```
code-bool(true,  $e_v$ ,  $e_f$ ) = j  $e_v$   
code-bool(false,  $e_v$ ,  $e_f$ ) = j  $e_f$   
code-bool(not  $B_1$ ,  $e_v$ ,  $e_f$ ) = code-bool( $B_1$ ,  $e_f$ ,  $e_v$ )  
code-bool( $B_1$  or  $B_2$ ,  $e_v$ ,  $e_f$ ) =  
    code-bool( $B_1$ ,  $e_v$ ,  $\text{new-e}$ ) |  $\text{new-e}$  : code-bool( $B_2$ ,  $e_v$ ,  $e_f$ )  
code-bool( $B_1$  and  $B_2$ ,  $e_v$ ,  $e_f$ ) =  
    code-bool( $B_1$ ,  $\text{new-e}$ ,  $e_f$ ) |  $\text{new-e}$  : code-bool( $B_2$ ,  $e_v$ ,  $e_f$ )  
code-bool( $E_1$  relop  $E_2$ ,  $e_v$ ,  $e_f$ ) =  
    code( $E_1$ ) | pushr  $\$a0$  | code( $E_2$ ) | popr  $\$t0$  |  
    code(relop)  $\$a0$ ,  $\$t0$ ,  $\$a0$  | beqz  $\$a0$   $e_f$  | j  $e_v$ 
```

$\text{new-e}$  : nouvelle étiquette

*code*(**relop**) l'instruction de comparaison pour l'opérateur **relop**: **sle**, **slt**, ...

- On évite l'empilement de valeurs intermédiaires

```
code(if E then I1 else I2) =  
code-bool(E, new-vrai, new-faux)  
| new-vrai: code(I1) | j new-suiv  
| new-faux: code(I2) | new-suiv:
```