

Compilation

CM5 - Génération de code

ISTIC, Université de Rennes 1
`Sebastien.Ferre@irisa.fr`

COMP, M1 info

Plan

- 1 Introduction
- 2 Machine intermédiaire
 - Opérateurs arithmétiques
 - Opérateurs logiques
 - Opérateurs de comparaison
 - Opérateurs de contrôle
 - Opérateurs d'adressage
- 3 Langage source
- 4 Génération de code 3-adresses
 - Expressions arithmétiques
 - Expressions booléennes
 - Instructions
 - Appel et définitions de fonctions

Plan

- 1 Introduction
- 2 Machine intermédiaire
 - Opérateurs arithmétiques
 - Opérateurs logiques
 - Opérateurs de comparaison
 - Opérateurs de contrôle
 - Opérateurs d'adressage
- 3 Langage source
- 4 Génération de code 3-adresses
 - Expressions arithmétiques
 - Expressions booléennes
 - Instructions
 - Appel et définitions de fonctions

Rappels

Analyse sémantique : 2 rôles

- Vérifications

- table des symboles
- vérification de types

- traduction vers le langage cible

= génération de code

- 2 aspects :

① impératif : préservation de la sémantique

que ça fasse ce que l'on veut

② secondaire : code généré efficace (optimisation)

que ça aille aussi vite que possible

que ça consomme le moins de ressources possible

Motivation du code intermédiaire

Les 2 aspects de la traduction motivent un découpage de la génération de code en 2 phases :

AST décoré $\xrightarrow{1}$ code intermédiaire $\xrightarrow{2}$ code cible

- ① se concentre sur la **traduction** proprement dite
 - **machine intermédiaire** offrant une **abstraction** des différentes machines
 - **processeur + mémoire + pile**
 - la programmation de cette phase doit être assez **simple** pour garantir la préservation de la sémantique
 - **approche compositionnelle**
 - on accepte que le code produit soit **inefficace**
- ② spécialisation du code vers une **machine particulière**
 - traduction relativement simple : **mot à mot**

Optimisations possible à chaque niveau

Motivation du code intermédiaire

Autre bénéfice du découpage en 2 phases :

- factorisation d'un compilateur à l'autre

.....

Nature de la machine intermédiaire

soit **machine** au sens classique

- instructions, emplacements mémoires
- mais **abstrait** des détails et contraintes
- ex : **machine 3-adresses**

soit **langage** pour lequel il existe déjà un compilateur

-

- implique 2 phases d'analyse syntaxique

Nature de la machine intermédiaire

soit **machine** au sens classique

- instructions, emplacements mémoires
- mais **abstrait** des détails et contraintes
- ex : **machine 3-adresses**

soit **langage** pour lequel il existe déjà un compilateur

-

- implique 2 phases d'analyse syntaxique

Plan

- 1 Introduction
- 2 **Machine intermédiaire**
 - Opérateurs arithmétiques
 - Opérateurs logiques
 - Opérateurs de comparaison
 - Opérateurs de contrôle
 - Opérateurs d'adressage
- 3 Langage source
- 4 Génération de code 3-adresses
 - Expressions arithmétiques
 - Expressions booléennes
 - Instructions
 - Appel et définitions de fonctions

Machine intermédiaire

- Machine 3-adresses
 - exécutant du **code 3-adresses**
 - code 3-adresses = liste d'**instructions 3-adresses**
- **instruction 3-adresses** : (op, x, y, z)
 - op : opération élémentaire, type d'instruction
 - x, y, z : emplacement mémoires ou registres
indifférenciés à ce stade
 - en général : x = résultat, y, z = opérandes
 - en fait : *au plus* 3-adresses
toutes les opérations n'utilisent pas les 3 adresses

Opérateurs arithmétiques

Opérandes et résultats de type entier

nom	instruction	notation	arité
addition	$(+, x, y, z)$	$x = y + z$	binaire
soustraction	$(-, x, y, z)$	$x = y - z$	binaire
multiplication	$(*, x, y, z)$	$x = y * z$	binaire
division	$(/, x, y, z)$	$x = y / z$	binaire
modulo	$(\%, x, y, z)$	$x = y \% z$	binaire
- unaire	$(-unaire, x, y, _)$	$x = -y$	unaire
constante ($N \in \mathbb{N}$)	$(const\ N, x, _, _)$	$x = N$	0-aire

Opérateurs logiques

Opérandes et résultats de type booléen

nom	instruction	notation	arité
et	<i>(and, x, y, z)</i>	$x = y \text{ and } z$	binaire
ou	<i>(or, x, y, z)</i>	$x = y \text{ or } z$	binaire
non	<i>(not, x, y, _)</i>	$x = \text{not } y$	unaire
vrai	<i>(const true, x, _, _)</i>	$x = \text{true}$	0-aire
faux	<i>(const false, x, _, _)</i>	$x = \text{false}$	0-aire

Opérateurs de comparaison

Opérandes de type entier et résultat de type booléen

nom	instruction	notation	arité
eq	(eq, x, y, z)	$x = y == z$	binaire
neq	(neq, x, y, z)	$x = y != z$	binaire
leq	(leq, x, y, z)	$x = y <= z$	binaire
lt	$(lt, x, y, _)$	$x = y < z$	binaire
geq	(geq, x, y, z)	$x = y >= z$	binaire
gt	$(gt, x, y, _)$	$x = y > z$	binaire

Opérateurs de contrôle

nom	instruction	notation	arité
saut inconditionnel	<i>(goto L, _, _, _)</i>	goto <i>L</i>	0-aire
saut si zero	<i>(ifz L, x, _, _)</i>	ifz <i>x</i> goto <i>L</i>	unaire
saut si non-zero	<i>(ifnz L, x, _, _)</i>	ifnz <i>x</i> goto <i>L</i>	unaire
entrée fonc/proc	<i>(begin L, _, _, _)</i>	begin <i>L</i>	0-aire
retour fonc	<i>(return, x, _, _)</i>	return <i>x</i>	unaire
retour proc	<i>(return, _, _, _)</i>	return	0-aire
passage argument	<i>(arg, x, _, _)</i>	arg <i>x</i>	unaire
appel fonc	<i>(call L, x, _, _)</i>	<i>x</i> = call <i>L</i>	0-aire
appel proc	<i>(call L, _, _, _)</i>	call <i>L</i>	0-aire

où *L* est une **étiquette de code**

- position dans le code
- constante car connue à la compilation

Opérateurs d'affectation et d'adressage

nom	instruction	notation	arité
copie	$(=, x, y, _)$	$x = y$	unaire
déréf. droit	$(*_D, x, y, _)$	$x = *y$	unaire
	si y contient une adresse		
déréf. gauche	$(*_G, x, y, _)$	$*x = y$	binaire
	si x contient une adresse		
“adresse de”	$(\&, x, y, _)$	$x = \&y$	unaire
	si y est un emplacement mémoire		

Diagrammes mémoire de ces opérateurs :

Plan

- 1 Introduction
- 2 Machine intermédiaire
 - Opérateurs arithmétiques
 - Opérateurs logiques
 - Opérateurs de comparaison
 - Opérateurs de contrôle
 - Opérateurs d'adressage
- 3 **Langage source**
- 4 Génération de code 3-adresses
 - Expressions arithmétiques
 - Expressions booléennes
 - Instructions
 - Appel et définitions de fonctions

Langage source

Dans le cadre de ce cours, on considère le langage **BABIL**

- langage **impératif** tel que C ou PASCAL
- contenant les principales constructions
 - expressions arithmétiques et booléennes
 - structures de contrôle : conditionnelles et boucles
 - définitions et appels de fonctions et procédures
- ordre d'évaluation des expressions non fixé (comme en C, Java)
 - sauf pour les expressions booléennes (comme en C, Java)

Syntaxe BABIL : unités lexicales

id → $[a - zA - Z][a - zA - Z0 - 9_]+$
num → $[0 - 9]^+$
bool → true | false
op → + | - | * | /
rop → == | =< | >= | < | > | !=
bop → and | or

Syntaxe BABIL : expressions arith/bool

$E(xpr)$	\rightarrow	num
	$ $	id
	$ $	$E \ op \ E$
	$ $	$- \ E$
	$ $	$id \ (\ LA \)$
$L(ist)A(rgs)$	\rightarrow	$E \ , \ LA \ \ E \ \ \epsilon$
$B(oolExpr)$	\rightarrow	$bool$
	$ $	id
	$ $	$E \ rop \ E$
	$ $	$B \ bop \ B$
	$ $	$\underline{not} \ B$

Syntaxe BABIL : programmes, fonctions, instructions

$P(\text{rogram})$	\rightarrow	LF
$LF(\text{unction})$	\rightarrow	$F \text{ } LF \mid F$
$F(\text{unction})$	\rightarrow	$\underline{\text{define}} \text{ } id \text{ } (\text{ } LP \text{ }) \text{ } LS \text{ } \underline{\text{end}}$
$LP(\text{arameter})$	\rightarrow	$id \text{ } , \text{ } LP \mid id \mid \epsilon$
$LS(\text{atement})$	\rightarrow	$S \text{ } ; \text{ } LS \mid S$
$S(\text{atement})$	\rightarrow	$id = E$
		$\underline{\text{if}} \text{ } B \text{ } \underline{\text{then}} \text{ } LS \text{ } \underline{\text{end}}$
		$\underline{\text{if}} \text{ } B \text{ } \underline{\text{then}} \text{ } LS \text{ } \underline{\text{else}} \text{ } LS \text{ } \underline{\text{end}}$
		$\underline{\text{while}} \text{ } B \text{ } \underline{\text{do}} \text{ } LS \text{ } \underline{\text{done}}$
		$id \text{ } (\text{ } LA \text{ })$
		$\underline{\text{return}}$
		$\underline{\text{return}} \text{ } E$

Plan

- 1 Introduction
- 2 Machine intermédiaire
 - Opérateurs arithmétiques
 - Opérateurs logiques
 - Opérateurs de comparaison
 - Opérateurs de contrôle
 - Opérateurs d'adressage
- 3 Langage source
- 4 Génération de code 3-adresses
 - Expressions arithmétiques
 - Expressions booléennes
 - Instructions
 - Appel et définitions de fonctions

Génération de code 3-adresses

- code 3-adresses = séquence d'instructions
→ impose de **linéariser** le code

- exemple : $x + 2 * z$

- 3 opérations : $+$, $*$, 2
- donc minimum 3 instructions : $+$, $*$, *const 2*
- ordre : *const 2* \rightarrow $*$ \rightarrow $+$

- code :

1	(<i>const 2</i> , $t1$, $_$, $_$)	$t1 = 2$
2	($*$, $t2$, $t1$, z)	$t2 = t1 * z$
3	($+$, $t3$, x , $t2$)	$t3 = x + t2$

- x, y : variables du programme source

a priori en mémoire

- $t1, t2, t3$: variables intermédiaires

→ inventées (allouées) par le compilateur

a priori en registre

- $t3$ contient le résultat de l'expression

Génération de code 3-adresses

- code 3-adresses = séquence d'instructions
→ impose de **linéariser** le code

- exemple : $x + 2 * z$

- 3 opérations : $+$, $*$, 2
- donc minimum 3 instructions : $+$, $*$, *const 2*
- ordre : *const 2* \rightarrow $*$ \rightarrow $+$

- code :
- | | | |
|---|---|---------------|
| 1 | (<i>const 2</i> , $t1$, $_$, $_$) | $t1 = 2$ |
| 2 | ($*$, $t2$, $t1$, z) | $t2 = t1 * z$ |
| 3 | ($+$, $t3$, x , $t2$) | $t3 = x + t2$ |

- x, y : variables du programme source

a priori en mémoire

- $t1, t2, t3$: variables intermédiaires

→ inventées (allouées) par le compilateur

a priori en registre

- $t3$ contient le résultat de l'expression

Génération de code 3-adresses

- code 3-adresses = séquence d'instructions
→ impose de **linéariser** le code

- exemple : $x + 2 * z$

- 3 opérations : $+$, $*$, 2
- donc minimum 3 instructions : $+$, $*$, *const 2*
- ordre : *const 2* \rightarrow $*$ \rightarrow $+$

- code :

1	(<i>const 2</i> , $t1$, $_$, $_$)	$t1 = 2$
2	($*$, $t2$, $t1$, z)	$t2 = t1 * z$
3	($+$, $t3$, x , $t2$)	$t3 = x + t2$

- x, y : **variables du programme source**

a priori en mémoire

- $t1, t2, t3$: **variables intermédiaires**

→ inventées (allouées) par le compilateur

a priori en registre

- $t3$ contient le résultat de l'expression

Génération de code 3-adresses

Comme pour la vérification de types ou la production de l'AST

- le plus simple est de procéder de **façon compositionnelle**
 - définir une “valeur” (ici : code généré) pour chaque **construction** du langage, en isolation
 - **grammaire attribuée** avec comme **attribut synthétisé** principal le **code 3-adresse généré** $X.code$
- + “prises” : attributs hérités et synthétisés supplémentaires
 - permettant aux constructions englobantes d'assembler les sous-séquences de code
 - ex : $t3$ comme **emplacement du résultat de l'expression** $E.place$
 - ex : $t2$ comme **résultat de $2*z$, utilisé dans $x + 2*z$**
 - ces prises diffèrent d'un NT à l'autre

Génération de code 3-adresses

Rappel : **impératif** = préserver la sémantique

- génération compositionnelle & systématique (grammaire attribuée)
- pas d'optimisation
 - on ne réutilise pas les **variables intermédiaires**
→ fonction `nouvar()`
 - idem pour les **étiquettes de code**
→ fonction `nouvetiq()`

Génération de code 3-adresses

Remarque

La génération de code est le lieu de rencontre de **3 programmes** :

- 1 le programme **source** (construction de la grammaire)
- 2 le programme **générateur** (code du compilateur)
- 3 le programme **généré** (code 3-adresses)

Ils sont entremêlés dans la grammaire attribuée :

- 1 source : règles syntaxiques
- 2 générateur : actions/calculs associés aux règles
- 3 généré : codes 3-adresses manipulés comme valeurs par les actions/calculs

C'est sans doute la principale difficulté de la compilation !

Génération de code pour les expressions arithmétiques

.....

Grammaire attribuée pour les expressions arithmétiques

$$\begin{aligned} E &\rightarrow \text{num} \\ &\quad \begin{cases} E.place := \text{nouvar}() \\ E.code := (\text{const num.vallex}, E.place, _, _) \end{cases} \\ | \text{id} \\ &\quad \begin{cases} E.place := TS.place(\text{id.vallex}) \\ E.code := \text{vide} \end{cases} \\ | E' \text{ op } E'' \\ &\quad \begin{cases} E.place := \text{nouvar}() \\ E.code := E'.code @@ E''.code \\ \quad @@ (\text{op.vallex}, E.place, E'.place, E''.place) \end{cases} \\ | - E' \\ &\quad \begin{cases} E.place := \text{nouvar}() \\ E.code := E'.code @@ (-\text{unaire}, E.place, E'.place, _) \end{cases} \end{aligned}$$

Exemple de génération pour les expressions arithmétiques

$$E = (a + b) * (a + b) \dots\dots$$

Génération de code pour les expressions booléennes

- On peut faire comme pour les expressions arithmétiques
 - attributs *B.code* et *B.place*
- mais *B* utilisé uniquement comme condition de branchement
 - le résultat ne sert qu'à décider du branchement
- de plus, dans `(a and b) or c`
 - si *a* = *false*, inutile d'évaluer *b*
 - si *a* = *true* et *b* = *true*, inutile d'évaluer *c*

Code court-circuit pour les expressions booléennes

On va fait du **code court-circuit**

- branchement vers une étiquette *B.siVrai* si $B = true$ et vers *B.siFaux* si $B = faux$
 - plus de résultat explicite
→ l'attribut *B.place* n'est plus défini
 - *B.siVrai* et *B.siFaux* sont des attributs **hérités** qui informent *B* de "où brancher" en fonction de la valeur de *B*
- en n'évaluant que ce qui est nécessaire (aspect "court-circuit")

Génération de code court-circuit pour les expressions booléennes

.....

Grammaire attribuée pour les expressions booléennes

$$\begin{aligned} B &\rightarrow id \\ &\quad \left\{ \begin{array}{l} B.code := (ifnz B.siVrai, TS.place(id.vallex), _, _) \\ \quad @@ \quad (goto B.siFaux, _, _, _) \end{array} \right. \\ &| bool \\ &\quad \left\{ \begin{array}{l} B.code := si \text{ bool.vallex} \\ \quad \text{alors } (goto B.siVrai, _, _, _) \\ \quad \text{sinon } (goto B.siFaux, _, _, _) \end{array} \right. \\ &| E' \text{ rop } E'' \\ &\quad \left\{ \begin{array}{l} resultat := \text{nouvar}() \\ B.code := E'.code @@ E''.code \\ \quad @@ \quad (rop.vallex, resultat, E'.place, E''.place) \\ \quad @@ \quad (ifnz B.siVrai, resultat, _, _) \\ \quad @@ \quad (goto B.siFaux, _, _, _) \end{array} \right. \end{aligned}$$

Grammaire attribuée pour les expressions booléennes

$$\begin{array}{lcl}
 B & \rightarrow & \text{not } B' \\
 & & \left\{ \begin{array}{lcl} B'.siVrai, B'.siFaux & := & B.siFaux, B.siVrai \\ B.code & := & B'.code \end{array} \right. \\
 & | & B' \text{ and } B'' \\
 & & \left\{ \begin{array}{lcl} B'.siVrai, B'.siFaux & := & \text{nouvetiq}(), B.siFaux \\ B''.siVrai, B''.siFaux & := & B.siVrai, B.siFaux \\ B.code & := & B'.code \\ & @@ & (\text{label } B'.siVrai, _, _, _) \\ & @@ & B''.code \end{array} \right. \\
 & | & B' \text{ or } B'' \\
 & & \left\{ \begin{array}{lcl} B'.siVrai, B'.siFaux & := & B.siVrai, \text{nouvetiq}() \\ B''.siVrai, B''.siFaux & := & B.siVrai, B.siFaux \\ B.code & := & B'.code \\ & @@ & (\text{label } B'.siFaux, _, _, _) \\ & @@ & B''.code \end{array} \right.
 \end{array}$$

Exemple de génération pour les expressions booléennes

$B = (a \text{ and } b) \text{ or } c \dots\dots$

Génération de code pour les instructions

.....

Grammaire attribuée pour les instructions

$$\begin{aligned}
 S &\rightarrow id = E \\
 &\quad \left\{ \begin{array}{l} S.code := E.code @@ (=, TS.place(id.vallex), E.place, _) \\ \text{if } B \text{ then } LS' \text{ end} \\ \left\{ \begin{array}{l} B.siVrai := alors := nouveliq() \\ B.siFaux := fin := nouveliq() \\ S.code := B.code \\ @@ (label\ alors, _, _, _) @@ LS'.code \\ @@ (label\ fin, _, _, _) \end{array} \right. \\ \text{if } B \text{ then } LS' \text{ else } LS'' \text{ end} \\ \left\{ \begin{array}{l} B.siVrai := alors := nouveliq() \\ B.siFaux := sinon := nouveliq() \\ fin := nouveliq() \\ S.code := B.code \\ @@ (label\ alors, _, _, _) @@ LS'.code \\ @@ (goto\ fin, _, _, _) \\ @@ (label\ sinon, _, _, _) @@ LS''.code \\ @@ (label\ fin, _, _, _) \end{array} \right. \end{array}
 \end{aligned}$$

$$S \rightarrow \underline{\text{while } B \text{ do } LS' \text{ done}}$$
$$\left\{ \begin{array}{ll} \text{boucle} & := \text{nouveliq}() \\ B.\text{siVrai} & := \text{corps} := \text{nouveliq}() \\ B.\text{siFaux} & := \text{sortie} := \text{nouveliq}() \\ S.\text{code} & := (\text{label } \text{boucle}, _, _, _) @@ B.\text{code} \\ & @@ (\text{label } \text{corps}, _, _, _) @@ LS'.\text{code} \\ & @@ (\text{goto } \text{boucle}, _, _, _) \\ & @@ (\text{label } \text{sortie}, _, _, _) \end{array} \right.$$
$$LS \rightarrow \begin{array}{l} S ; LS' \\ \left\{ \begin{array}{l} LS.\text{code} := S.\text{code} @@ LS'.\text{code} \\ S \\ \left\{ \begin{array}{l} LS.\text{code} := S.\text{code} \end{array} \right. \end{array} \right.$$

Exemple de génération pour les instructions

```
S = i = 1; f = 1; while i < n and f < 100 do i =  
i + 1; f = f * i done .....
```


Appels de fonctions/procédures

- 2 constructions
 - $E \rightarrow id (LA)$: appel de fonction (expression)
 - $S \rightarrow id (LA)$: appel de procédure (instruction)
- il existe 2 sémantiques du **passage de paramètres**
 - **par nom** : nom d'une mémoire, emplacement
 - permet à la fonction de modifier le contenu de cette mémoire
 - ex : **paramètre** **var** de PASCAL
 - ex : **références** **sur objets** en JAVA
 - **par valeur** : contenu d'une mémoire
 - la valeur du paramètre est recopiée dans une zone mémoire réservée aux arguments
 - sur la pile ou dans des registres
 - ne permet pas de modifier l'argument
 - mais l'argument peut être un pointeur...
 - c'est la sémantique de C... et de BABIL
- c'est l'appelant qui a la charge d'installer les paramètres avant d'appeler la fonction
 - opérateur (*arg. x. . .*)

Appels de fonctions/procédures

- 2 constructions
 - $E \rightarrow id (LA)$: appel de fonction (expression)
 - $S \rightarrow id (LA)$: appel de procédure (instruction)
- il existe 2 sémantiques du **passage de paramètres**
 - **par nom** : nom d'une mémoire, emplacement
 - permet à la fonction de modifier le contenu de cette mémoire
 - ex : **paramètre var de PASCAL**
 - ex : **références sur objets en JAVA**
 - **par valeur** : contenu d'une mémoire
 - la valeur du paramètre est recopiée dans une zone mémoire réservée aux arguments
 - sur la pile ou dans des registres**
 - ne permet pas de modifier l'argument
 - mais l'argument peut être un pointeur...**
 - c'est la sémantique de C... et de BABIL
- c'est l'appelant qui a la charge d'installer les paramètres avant d'appeler la fonction
 - opérateur (*arg. x. . .*)

Génération de code pour les appels de fonctions

.....

Grammaire attribuée pour les appels de fonctions

$$E \rightarrow id (LA)$$

$$\begin{cases} E.place &:= \text{nouvar}() \\ E.code &:= LA.code @@ (\text{call } TS.label(id.vallex), E.place, \end{cases}$$

$$S \rightarrow id (LA)$$

$$\begin{cases} S.code &:= LA.code @@ (\text{call } TS.label(id.vallex), _, _, _) \end{cases}$$

$$LA \rightarrow E , LA'$$

$$\begin{cases} LA.code &:= E.code @@ (arg, E.place, _, _) @@ LA'.code \end{cases}$$

$$|$$

$$E$$

$$\begin{cases} LA.code &:= E.code @@ (arg, E.place, _, _) \end{cases}$$

$$|$$

$$\epsilon$$

$$\begin{cases} LA.code &:= vide \end{cases}$$

Génération de code pour les définitions de fonctions

.....

Grammaire attribuée pour les définitions de fonctions

$S \rightarrow \underline{return}$
 $\left\{ \begin{array}{l} S.code := (return, _, _, _) \\ \underline{return} E \\ S.code := E.code @@ (return, E.place, _, _) \end{array} \right.$

$F \rightarrow id (LP) LS$
 $\left\{ \begin{array}{l} F.code := (beginfunc TS.label(id.vallex), _, _, _) \\ @@ LS.code \end{array} \right.$

$LF \rightarrow F LF'$
 $\left\{ \begin{array}{l} LF.code := F.code @@ LF'.code \\ F \\ LF.code := F.code \end{array} \right.$

$P \rightarrow LF$
 $\left\{ \begin{array}{l} P.code := LF.code \end{array} \right.$

Exemple de génération pour les fonctions

Exemple : `define fact(n) if n==0 then return 1
else return n * fact(n-1) end.....`