

Plan

Traduction et Sémantique Introduction

Luigi Santocanale
LIF, Université de Provence
Marseille, FRANCE

12 février 2010

Préambule

Contexte, historique

LSE : un Langage Simple d'Expressions

Le langage source

Une machine virtuelle

La traduction

Des questions sémantiques

Évaluation

Typage

Plan

Des coordonnées

Préambule

Contexte, historique

LSE : un Langage Simple d'Expressions

Le langage source

Une machine virtuelle

La traduction

Des questions sémantiques

Évaluation

Typage

• Page web :

<http://www.lif.univ-mrs.fr/~lsantoca/teaching/TS/>

• Mon courriel : lsantoca@lif.univ-mrs.fr

• Formule de la note finale :

$$NF = (2 \max(E, (2E+P) / 3) + TP) / 3$$

où

- ▶ E est la note (rattrapable) de l'examen,
- ▶ P est la note (non rattrapable) du partiel,
- ▶ TP est la note (non rattrapable) du projet à élaborer en TP.
Présence à la soutenance obligatoire.

- Partiel, autour du 7ème cours.
- Projet : à la 6ème séance de TP.
(Présence à la soutenance obligatoire.)

Bibliographie :

-  **Roberto M. Amadio.**
Introduction à l'analyse syntaxique et à la compilation, 04 2009.
-  **Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman.**
Compilers : Principles, Techniques, and Tools.
Addison-Wesley, 1986.
-  **J. Levine and D. Mason, T. Brown.**
lex and yacc.
O'Reilly, 1992.

5/43

Préambule

Contexte, historique

LSE : un Langage Simple d'Expressions

- Le langage source
- Une machine virtuelle
- La traduction
- Des questions sémantiques
- Évaluation
- Typage

6/43

Compilation

Compiler :

- Littéraire :
mettre ensemble.
Fait par l'éditeur de liens.
- En Info : processus de **traduction**
d'un langage source (de haut niveau)
vers un langage cible, (de bas niveau, lisible par la machine).
- Différent de l'*interprétation*,
traduction et exécution à la volée.
Voir les langages de scripts.
- Critère de correction de la traduction :
notion de **sémantique**.

7/43

Objectifs

- Comprendre ce qui se passe dans le processus de compilation.
- Apprendre les techniques et les algorithmes de ce processus.
- Maîtriser les outils (Lex, Yacc) qui nous permettent de construire des compilateurs.
- Apprécier ces techniques et ces outils,
en dehors de leur contexte d'origine,
pour des nouvelles problématiques.

8/43

D'autres applications

- Multitude de formats, transformations entre eux ...
... voir les exercices.
- Des exemples :
 - ▶ langages de composition de textes :
tex, latex, bibtex, ...
 - ▶ langages graphiques :
dot, dot2tex, ...
 - ▶ xml, html, & co :
nsgmls, ...
- Typage pour la sûreté

9/43

Analyse lexicale

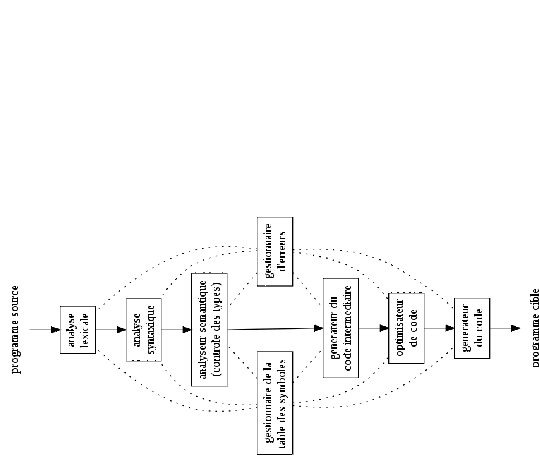
- Processus qui transforme (reconnait)
un flot de caractères vers une suite d'unités logiques
(unités lexicales).
- Exemple :

je	mange	la	pomme
pronom	verbe	article	nom

consiste de 16 caractères et 4 mots.
On peut associer à chaque mot une catégorie.
- Cadre théorique :
 - ▶ expressions et langages régulier(e)s,
 - ▶ automates à états finis.

11/43

Structure d'un compilateur



10/43

Analyse syntaxique

- Processus qui décerne les flots d'unités lexicales porteur de sens.
- Exemple :

je	mange	la	pomme
sujet	verbe	transitif	objet

verbe
- est un phrase, une phrase est constituée par un sujet et un verbe.
- En compilation :
processus qui décerne la structure algébrique (articulation logique) d'un programme à partir d'un flot d'unités lexicales.
- Une première signification à un programme.
- Cadre théorique :
 - ▶ grammaires non contextuelles, langages algébriques,
 - ▶ automates à pile.

12/43

De la linguistique à l'informatique



Noam Chomsky.

Three models for the description of language.
IRE Trans. Inform. Theory, (2 :3) :113–124, 1956.



Noam Chomsky.

On certain formal properties of grammars.
Information and Control, 2 :137–167, 1959.



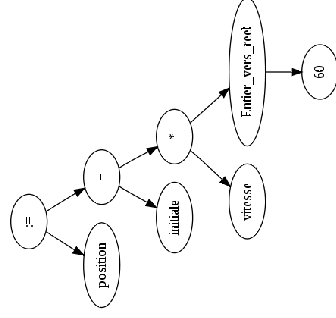
Peter Naur.

ALGOL—the international language for description of logical and numerical processes.

Nordisk Mat. Tidsskr., 8 :117–129, 144, 1960.

Transtypage

... être transformée en

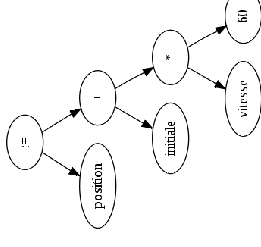


13/43

Analyse sémantique

Typiquement il s'agit de la phase de *contrôle de types*.

L'expression



peut être jugé incorrecte ou ...

14/43

Génération du code intermédiaire

- Représentation du programme dans un langage indépendant de l'architecture.
- Il s'agit, en général, d'un langage pour une machine abstraite.
- Code à trois adresses :

```
temp1 := EntiertVersReel(60)
temp2 := id3 * temp1
temp3 := id2 + temp2
id2   := temp3
```

15/43

16/43

- Traduction dans le langage cible, lisible par une architecture donnée (CISC et x86, RISC et Sun, ...).
- Assignment des variables aux registres.

Ainsi,

```
Un exemple :
/* temp1 := EntiertVersReel(60) */
temp2 := id3 * 60.0
/* x temp3 := id2 + temp2 */
id2 := id2 + temp2
```

```
position := initiale + vitesse * 60
```

devient :

```
11:position.c      ****      position=initiale + vitesse* 60;
34                .loc 1 11 0
35 001c D945F4      ifds      -12(%ebp)
36 001f D9050000    ifds      .LC0
36                0000
37 0025 DEC9       fmulp     %st, %st(1)
38 0027 D845F8     fadds    -8(%ebp)
39 002a D85DFC     fsips    -4(%ebp)
```

17/43

18/43

Plan

Préambule

Contexte, historique

LSE : un Langage Simple d'Expressions

Le langage source

Une machine virtuelle

La traduction

Des questions sémantiques

Évaluation

Typage

Le langage

- constantes numériques :
1, 03457
- constantes booléennes :
true, false
- opérateurs :
+, *, ...
- variables :
x, y, z, s46uu, ...
- déclarations :
let ... = ... in ...
- contrôle du flot :
if ... then ... else ...

19/43

20/43

Analyse lexicale

La chaîne de caractères

let x7 = 3 in (x7 + 4)

est analysée comme suit :

Chaîne (lexème) :	let	x7	=	3	in	(x7	+	4)
Unité lexicale :	let	id	eq	cnst	in	lpar	id	plus	cnst	rpar
Valeur (attribut) :		x7		3		x7			4	

21/43

Dérivation (gauche)

La chaîne d'unités lexicales appartient au langage de \mathcal{G} ,
car il existe une dérivation du symbole initial S vers la chaîne :

$S \Rightarrow B$
 $\Rightarrow \text{let id eq } E \text{ in } B$
 $\Rightarrow \text{let id eq const in } B$
 $\Rightarrow \text{let id eq const in } E$
 $\Rightarrow \text{let id eq const in lpar } E \text{ rpar}$
 $\Rightarrow \text{let id eq const in lpar } E \text{ OP } E \text{ rpar}$
 $\Rightarrow \text{let id eq const in lpar id OP } E \text{ rpar}$
 $\Rightarrow \text{let id eq const in lpar id plus } E \text{ rpar}$
 $\Rightarrow \text{let id eq const in lpar id plus cnst rpar}$

23/43

La définition du langage

...se fait par une grammaire non contextuelle

$\mathcal{G} = \langle V, \Sigma, S, P \rangle$

où

$V = \{ S, OP, E, B \} \cup$
 $\{ \text{plus, prod, \dots, id, const, lpar, rpar, if, then, else, let, eq, in} \}$

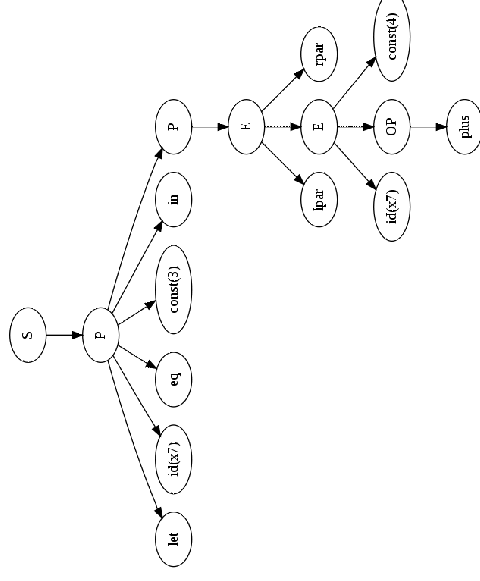
et P est

$OP \rightarrow \text{plus} \mid \text{prod} \mid \dots$
 $E \rightarrow \text{id} \mid \text{const} \mid E \text{ OP } E \mid \text{lpar } E \text{ rpar}$
 $B \rightarrow E \mid \text{if } E \text{ then } B \text{ else } B \mid \text{let id eq } E \text{ in } B \mid \text{lpar } B \text{ rpar}$
 $S \rightarrow B$

22/43

Arbre de dérivation

L'appartenance au langage de \mathcal{G} est aussi témoigné par un arbre de dérivation :



24/43

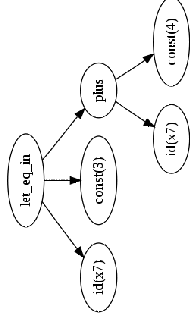
Syntaxe abstraite

A partir de l'arbre de dérivation on construit sa syntaxe abstraite.

Notre exemple donne :

`let_in_eq(⟨id, x7⟩, ⟨const, 3⟩, plus(⟨id, x7⟩, ⟨const, 4⟩))`

ou bien



25/43

26/43

Considération philosophiques sur la syntaxe abstraite

- Un programme est un terme – comme ceux de la logique.
- On peut représenter un terme comme un arbre étiqueté. Cela permet de se débarrasser de symboles non intéressants (parenthèses, virgules).
- On traite la construction `Let...eq...in...` comme une unité, un symbole de fonction dearité 3.
- La structure mathématique des arbres est inductive :
 - ▶ on peut définir des fonctions par induction sur la structure,
 - ▶ on peut définir des fonctions par induction sur la syntaxe abstraite.

Définition de la syntaxe abstraite

Il s'agit d'arbres étiquetés.

Un *E*-terme :

étiquette	infos	enfants
id	chaîne de caractères	0
const	nombre, true, false	0
plus		2 <i>E</i> -termes
prod		

Un *B*-terme : un *E*-terme, ou bien

étiquette	infos	enfants
let_eq_in		2 <i>E</i> -termes, 1 <i>B</i> -terme
if_then_else_		1 <i>E</i> -terme, 2 <i>B</i> -termes

27/43

28/43

Considération moins philosophiques sur la syntaxe abstraite

- Implantation de la syntaxe abstraite :

- ▶ en C, par les pointeurs (arbres = généralisation des listes)
 - ▶ en d'autres langages par les types inductifs.
- En OCaml :

```

type valeur = Int of int | True | False
type eterm = Id of string | Const of valeur
           | Plus of eterm * eterm
type bterm = Expr of expr
           | Ifthenelse of eterm * bterm * bterm
           | Leteqin of string * eterm * eterm
  
```

- Construire explicitement la syntaxe abstraite peut être coûteux (en espace).

Si l'on peut, on laisse cette construction implicite.

Une structure de donnés : $\langle M, pc, sp \rangle$, où

- M : tableau en mémoire,
- pc : program counter, compteur d'instruction,
- sp : stack pointer, pointeur au sommet de la pile.

```

buildV      sp++; M[sp] := v; pc++;
branchj     si M[sp] = true alors pc++ sinon pc := j; sp-
Loadi      sp++; M[sp] := M[j]; pc++
add         sp++; M[sp] := M[sp] + M[sp + 1]; pc++
return     pc := 0; M[0] := M[sp]; sp := 0
    
```

29/43

30/43

Fonction de traduction

Fonction de traduction, définie par induction sur la syntaxe abstraite.

$$\begin{aligned}
 \mathcal{C}_E(x, w) &= \text{load } i(x, w) \\
 \mathcal{C}_E(v, w) &= \text{build } v \\
 \mathcal{C}_E(e_1 + e_2, w) &= \mathcal{C}_E(e_1, w) \cdot \mathcal{C}_E(e_2, w) \cdot \text{add} \\
 \mathcal{C}_B(x, w) &= (\text{load } i(x, w)) \cdot \text{return} \\
 \mathcal{C}_B(v, w) &= (\text{build } v) \cdot \text{return} \\
 \mathcal{C}_B(e_1 + e_2, w) &= \mathcal{C}_E(e_1, w) \cdot \mathcal{C}_E(e_2, w) \cdot \text{add} \cdot \text{return} \\
 \mathcal{C}_B(\text{let } x = \text{einp } w) &= \mathcal{C}_E(e, w) \cdot \mathcal{C}_B(b, w \cdot x) \\
 \mathcal{C}_B(\text{if } e \text{ then } b_1 \text{ else } b_2, w) &= \mathcal{C}_E(e, w) \cdot (\text{branch } k) \cdot \mathcal{C}_B(b_1, w) \cdot k : \mathcal{C}_B(b_2, w)
 \end{aligned}$$

31/43

Résultat :-)

Compiler

```

let x = 3 in
let y = x + x in
  let x = true in
    if x then y else x
    
```

donne

```

1 : build 3
2 : load 1
3 : load 1
4 : add
5 : build true
6 : load 3
7 : branch k
8 : load 2
9 : return
k=10 : load 3
11 : return
    
```

32/43

- Que veut dire une expression telle que

$$\text{if } x \text{ then } 0 \text{ else } 1$$
- Est ce que le programme

$$\text{if } 33 \text{ then } 1 \text{ else } 1$$
est correcte ?
- Est ce que la traduction est correcte ?

33/43

Définition de \Downarrow

$$\frac{}{v \Downarrow v}$$

$$\frac{e_1 \Downarrow n_1 \in \mathbb{N} \quad e_2 \Downarrow n_2 \in \mathbb{N}}{e_1 + e_2 \Downarrow n_1 + n_2}$$

$$\frac{e \Downarrow \text{true} \quad p_1 \Downarrow v}{\text{if } e \text{ then } p_1 \text{ else } p_2 \Downarrow v}$$

$$\frac{e \Downarrow \text{false} \quad p_2 \Downarrow v}{\text{if } e \text{ then } p_1 \text{ else } p_2 \Downarrow v}$$

$$\frac{e \Downarrow v' \quad [v'/x]p \Downarrow v}{\text{let } x = e \text{ in } p \Downarrow v}$$

35/43

Entre les expressions, on dit que

$\text{true}, \text{false}, \text{cnst}(0), \text{cnst}(1), \dots$
sont des valeurs.

On peut définir une relation

$$e \Downarrow v$$

par induction sur la syntaxe abstraite. À lire :

l'expression e s'évalue au valeur v .

34/43

Notre exemple

L'expression $\text{let } x = 3 \text{ in } (x + 4)$ s'évalue à 7 :

$$\text{let } x = 3 \text{ in } (x + 4) \Downarrow 7$$

car on peut construire un preuve de cette relation :

$$\frac{\frac{3 \Downarrow 3}{\text{let } x = 3 \text{ in } (x + 4) \Downarrow 7} \quad \frac{4 \Downarrow 4}{3 + 4 \Downarrow 7}}{\text{let } x = 3 \text{ in } (x + 4) \Downarrow 7}$$

36/43

Considérations

- Prouver que $e \Downarrow v$ « revient à calculer » la valeur de e .
- La relation \Downarrow permet de donner une signification aux expressions, **indépendamment de** ce que fait **la machine virtuelle**.
- Nous pouvons poser formellement la question de la *correction de la traduction* :
Si $e \Downarrow v$ et la machine exécute le code $\mathcal{C}_B(e, \square)$, alors la machine s'arrête et $M[0]$ contient v .
Est ce que cette proposition est vraie ?
- La relation \Downarrow est un (important) outil mathématique. Difficilement on peut s'en servir dans le cadre de l'implémentation d'un compilateur.

37/43

Un premier essai

$$\frac{n \in \mathbb{N}}{n : \text{nat}} \qquad \frac{v \in \{\text{true}, \text{false}\}}{v : \text{bool}}$$
$$\frac{e_1 : \text{nat} \quad e_2 : \text{nat}}{e_1 + e_2 : \text{nat}} \qquad \frac{e : \text{bool} \quad b_1 : \tau \quad b_2 : \tau}{\text{if } e \text{ then } b_1 \text{ else } b_2 : \tau}$$
$$\frac{e : \tau' \quad x : \tau' \dashv b : \tau}{\text{let } x = e \text{ in } b}$$

Nous devons aussi ajouter des contraintes de types sur les variables !!!

39/43

Typage

On souhaite définir une relation de la forme

$$e : \tau$$

où $\tau \in \{\text{bool}, \text{nat}\}$, à lire :

l'expression e est de type τ ,

de façon que :

1. le fait suivant soit vrai :
 $e : \tau$ implique $\exists v. v : \tau$ et $e \Downarrow v$,
2. on puisse répondre à la question si $e : \tau$ de façon efficace.

38/43

Environnement

- Environnement E :
suite finie de contraintes de type sur les variables :

$$x_1 : \tau_1, x_2 : \tau_2, \dots, x_n : \tau_n$$

tel que $x_i \neq x_j$ si $i \neq j$.

- \emptyset : environnement vide.
- Environnement $E[\tau/x]$:
même environnement que E ,
sauf pour x qui est de type τ dans $E[\tau/x]$.

40/43

Deuxième essai

$$\frac{n \in \mathbb{N}}{E \vdash n : \text{nat}}$$

$$\frac{v \in \{ \text{true}, \text{false} \}}{E \vdash v : \text{bool}}$$

$$\frac{x \text{ a type } \tau \text{ dans } E}{E \vdash x : \tau}$$

$$\frac{E \vdash e_1 : \text{nat} \quad e_2 : \text{nat}}{E \vdash e_1 + e_2 : \text{nat}}$$

$$\frac{E \vdash e : \text{bool} \quad E \vdash b_1 : \tau \quad E \vdash b_2 : \tau}{E \vdash \text{if } e \text{ then } b_1 \text{ else } b_2 : \tau}$$

$$\frac{E \vdash e : \tau' \quad E[\tau'/x] \vdash b : \tau}{E \vdash \text{let } x = e \text{ in } b}$$

41/43

Le théorème principal

Théorème.

Si on peut typer l'expression e ,
alors cette expression s'évalue à un valeur.

Plus en détails :

Proposition.

Si $\emptyset \vdash e : \tau$, alors il existe v tel que $v : \tau$ et $e \Downarrow v$.

Pour pouvoir montrer cela :

Lemme.

Si $x_1 : \tau_1, \dots, x_n : \tau_n \vdash e : \tau$ et $\emptyset \vdash v_j : \tau_j$ pour $j = 1, \dots, n$,
alors il existe v tel que $v : \tau$ et $[v_1/x_1, \dots, v_n/x_n]e \Downarrow v$.

43/43

Exemples

L'expression

if true then 0 else 1

a type *nat*, car on peut « la typer » par *nat* :

$$\frac{\frac{\text{true} \in \{ \text{true}, \text{false} \}}{\emptyset \vdash \text{true} : \text{bool}} \quad \frac{0 \in \mathbb{N}}{\emptyset \vdash 0 : \text{nat}} \quad \frac{1 \in \mathbb{N}}{\emptyset \vdash 1 : \text{nat}}}{\emptyset \vdash \text{if true then 0 else 1} : \text{nat}}$$

On ne peut pas typer l'expression

if x then 0 else 1

car il n'existe pas $\tau \in \{ \text{bool}, \text{nat} \}$ et un arbre de la forme

$$\frac{\vdots}{\emptyset \vdash \text{if } x \text{ then } 0 \text{ else } 1 : \tau}$$

42/43