

# Langage de Programmation et Compilation

Jean-Cristophe Filiâtre

13 octobre 2023

## Table des matières

<b>I</b>	<b>Aperçu de la Compilation - Assembleur x86-64</b>	<b>2</b>
<b>1</b>	<b>Introduction à la Compilation</b>	<b>2</b>
1.1	Un Compilateur . . . . .	2
1.2	Le Bon et le Mauvais Compilateur . . . . .	2
1.3	Le Travail d'un Compilateur . . . . .	3
<b>2</b>	<b>L'assembleur</b>	<b>3</b>
2.1	Arithmétique des ordinateurs . . . . .	3
2.2	Architecture . . . . .	3
2.3	L'architecture x86-64 . . . . .	3
2.4	L'assembleur x86-64 . . . . .	4
2.5	Le Défi de la Compilation . . . . .	4
<b>II</b>	<b>Syntaxe abstraite, sémantique, interprètes</b>	<b>5</b>
<b>3</b>	<b>Sémantique Formelle</b>	<b>5</b>
3.1	Syntaxe Abstraite . . . . .	5
3.2	Sémantiques Useless . . . . .	6
3.2.1	Sémantique Axiomatique - Logique de Floyd-Hoare . . . . .	6
3.2.2	Sémantique Dénotationnelle . . . . .	6
3.2.3	Sémantique Par Traduction . . . . .	6
3.3	Sémantique Opérationnelle . . . . .	6
3.3.1	Sémantique Opérationnelle à Grand Pas . . . . .	6
3.3.2	Sémantique à Petits Pas . . . . .	7
3.3.3	Equivalence des Sémantiques . . . . .	7
3.3.4	Langages Impératifs . . . . .	8
<b>4</b>	<b>Interprète</b>	<b>8</b>
<b>III</b>	<b>Analyse Lexicale</b>	<b>8</b>
<b>5</b>	<b>Blancs</b>	<b>9</b>
<b>6</b>	<b>Outils pour l'analyse lexicale</b>	<b>9</b>
6.1	Expressions Régulières et Automates Finis . . . . .	9
6.1.1	Expressions Régulières . . . . .	9
6.1.2	Automates Finis . . . . .	10
6.2	Analyseur Lexical . . . . .	10
6.2.1	Principe . . . . .	10

6.2.2	Construction . . . . .	10
<b>7</b>	<b>L'outil ocamllex</b>	<b>11</b>
7.1	Analyseur Lexical en ocamllex . . . . .	11
7.2	Efficacité . . . . .	12
7.3	D'autres utilisations ocamllex . . . . .	12

## Première partie

# Cours 1 29/09

## Table des matières

<b>1</b>	<b>Introduction à la Compilation</b>	<b>2</b>
1.1	Un Compilateur . . . . .	2
1.2	Le Bon et le Mauvais Compilateur . . . . .	2
1.3	Le Travail d'un Compilateur . . . . .	3
<b>2</b>	<b>L'assembleur</b>	<b>3</b>
2.1	Arithmétique des ordinateurs . . . . .	3
2.2	Architecture . . . . .	3
2.3	L'architecture x86-64 . . . . .	3
2.4	L'assembleur x86-64 . . . . .	4
2.5	Le Défi de la Compilation . . . . .	4

## 1 Introduction à la Compilation

Maîtriser les mécanismes de la compilation, transformation d'un langage dans un autre. Comprendre les aspects des langages de programmation.

### 1.1 Un Compilateur

Un compilateur est un traducteur d'un langage source vers un langage cible. Ici le langage cible sera l'assembleur.

Tous les langages ne sont pas compilés à l'avance, certains sont interprétés, transpilés puis interprétés, compilés à la volée, transpilés puis compilés... Un compilateur prend un programme  $P$  et le traduit en un programme  $Q$  de sorte que :  $\forall P, \exists Q, \forall x, P(x) = Q(x)$ . Un interpréteur effectue un travail simple mais le refait à chaque entrée, et donc est moins efficace.

Exemple : le langage *lilypond* va compiler un code source en fichier .pdf.

### 1.2 Le Bon et le Mauvais Compilateur

On juge un compilateur à :

1. Sa correction
2. L'efficacité du code qu'il produit
3. Son efficacité en tant que programme

« Optimizing compilers are so difficult to get right that we dare say that no optimizing compiler is completely error-free! Thus, the most important objective in writing a compiler is that it is correct »- *Dragon Book, 2006*

## 1.3 Le Travail d'un Compilateur

Le travail d'un compilateur se compose :

- d'une phase d'analyse qui :
  1. reconnaît le programme à traduire et sa signification
  2. signale les erreurs et peut donc échouer
- d'une phase de synthèse qui :
  1. produit du langage cible
  2. utilise de nombreux langages intermédiaires
  3. n'échoue pas

Processus : source  $\rightarrow$  analyse lexicale  $\rightarrow$  suite de lexèmes (tokens)  $\rightarrow$  analyse syntaxique  $\rightarrow$  Arbre de syntaxe abstraite  $\rightarrow$  analyse sémantique  $\rightarrow$  syntaxe abstraite + table des symboles  $\rightarrow$  production de code  $\rightarrow$  langage assembleur  $\rightarrow$  assembleur  $\rightarrow$  langage machine  $\rightarrow$  éditeur de liens  $\rightarrow$  exécutable.

## 2 L'assembleur

### 2.1 Arithmétique des ordinateurs

On représente les entiers sur  $n$  bits numérotés de droite à gauche. Typiquement,  $n$  vaut 8, 16, 32 ou 64. On peut représenter des entiers non signés jusqu'à  $2^n - 1$ . On peut représenter les entiers en définissant  $b_{n-1}$  comme un bit de signe, on peut alors représenter  $[-2^{n-1}, 2^{n-1} - 1]$ . La valeur d'une suite de bits est alors :  $-b_{n-1}2^{n-1} + \sum_{k=0}^{n-2} b_k 2^k$ . On ne peut pas savoir si un entier est signé sans le contexte.

La machine fournit des opérations logiques (bit à bit), de décalage (ajout de bits 0 de poids fort, 0 de poids faible ou réplication du bit de signe pour interpréter une division), d'arithmétique (addition, soustraction, multiplication).

### 2.2 Architecture

Un ordinateur contient :

- Une unité de calcul (CPU) qui contient un petit nombre de registres et des capacités de calcul
- Une mémoire vive (RAM), composée d'un très grand nombre d'octets (8 bits), et des données et des instructions, indifférenciables sans contexte.

L'accès à la mémoire coûte cher : à 1B instructions/s, la lumière ne parcourt que 30cm entre deux instructions.

En réalité, il y a plusieurs (co)processeurs, des mémoires cache, une virtualisation de la mémoire... Principe d'exécution : un registre (%rip) contient l'adresse de l'instruction, on lit (fetch) un ou plusieurs octets dans la mémoire, on interprète ces bits (decode), on exécute l'instruction (execute), on modifie (%rip) pour l'instruction suivante. En réalité, on a des pipelines qui branchent plusieurs instructions en parallèle, et on essaie de prédire les sauts conditionnels.

Deux grandes familles d'Architectures : CISC (complex instruction set), qui permet beaucoup d'instructions différentes mais avec assez peu de registres, et RISC (Reduced Instruction Set) avec peu d'instruction effectuées très régulièrement et avec beaucoup de registres. Ici, on utilisera l'architecture *x86-64*.

### 2.3 L'architecture x86-64

Extension 64 bits d'une famille d'architectures compatibles Intel par AMD adoptée par Intel. Architecture à 16 registres, avec adressage sur 48 bits au moins et de nombreux modes d'adressage. On ne programme pas en langage machine mais en assembleur, langage symbolique avec allocation de données globales, qui est transformé en langage machine par un assembleur qui est en réalité un compilateur. On utilise l'assembleur GNU avec la syntaxe AT&T (la syntaxe Intel existe aussi).

## 2.4 L'assembleur x86-64

Pour assembler un programme assembleur, appeler `as -o file.o` puis appeler l'édition de lien avec `gcc -no-pie file.s -o exec-name`. On peut déboguer en ajoutant l'option `-g`. La machine est petite boutiste (little-endian) si elle stocke les valeurs dans la RAM en commençant par le bit de poids faible, gros boutiste (big-endian) pour le poids fort.

Commandes : Dans cette liste,  $\%(r)$  désigne l'adresse mémoire stockée dans  $r$

- `movq $a %b` permet de mettre la valeur  $a$  dans le registre  $b$
- `movq %a %b` permet de copier le registre  $a$  dans le registre  $b$
- `movq $label %b` permet de changer l'adresse de l'étiquette dans le registre  $b$
- `addq %a %b` permet d'additionner les registres  $a$  et  $b$ .
- `incq %r` permet d'incrémenter le registre  $r$ , de même pour `decq`.
- `negq %r` permet de modifier la valeur de  $r$  en sa négation
- `notq %r` permet de modifier la valeur de  $r$  en sa négation logique.
- `orq %r1 %r2` (resp. `andq` et `xorq`) permet d'affecter à  $r2$ ,  $or(r1, r2)$  (resp.  $and$ ,  $xor$ )
- `salq $n %r/salq %cl %r` décale la valeur de  $r$  de  $n$  (ou  $\%cl$ ) zéros à gauche.
- `sarq` est le décalage à droite arithmétique, `shrq` le décalage à droite logique.
- Le suffixe `q` désigne une opération sur 64 bits. `b` désigne 1 octet, `w` désigne 2 octets, `l` désigne 4 octets. Il faut préciser les deux extensions si celles-ci diffèrent.
- `jmp label` permet de jump à une étiquette.

La plupart des opérations positionnent des drapeaux selon leur résultat.

Certaines instructions : `j(suffixe)` (jump), `set(suffixe)` et `cmov(suffixe)` (move) permettent de tester des drapeaux et d'effectuer une opération selon leur valeur.

On ne sait pas combien il y a d'instructions en x86-64.

## 2.5 Le Défi de la Compilation

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

Constat : les appels de fonctions peuvent être arbitrairement imbriqués et les registres ne suffisent pas  $\Rightarrow$  on crée alors une pile car les fonctions procèdent majoritairement selon un mode LIFO.

La pile est stockée tout en haut, et croît dans le sens des adresses décroissantes, `%rsp` pointe sur le sommet de la pile. Les données dynamiques sont allouées sur le tas, en bas de la zone de données. Chaque programme a l'illusion d'avoir toute la mémoire pour lui tout seul, illusion créée par l'OS. En assembleur on a des facilités d'utilisation de la pile :

- `pushq $a` push  $a$  dans la pile
- `popq %rdi` dépile

Lorsque  $f$  (caller) appelle une fonction  $g$  (callee), on ne peut pas juste `jmp g`. On utilise `call g` puis une fois que c'est terminé `ret`.

Mézalor tout registre utilisé par  $g$  sera perdu par  $f$ . On s'accorde alors sur des **conventions d'appel**. Des arguments sont passés dans certains registres, puis sur la pile, la valeur de retour est passée dans `%rax`. Certains registres sont *callee-saved* i.e. l'appelé doit les sauvegarder pour qu'elle survive aux appels. Les autres registres sont dit *caller-saved* et ne vont pas survivre aux appels.

Il faut également qu'en entrée de fonction `%rsp + 8` doit être multiple de 16, sinon des fonctions peuvent planter.

Il y a quatre temps dans un appel :

1. Pour l'appelant, juste avant l'appel :
2. Pour l'appelé, au début de l'appel :
  - (a) Sauvegarde `%rbp` puis le positionne.
  - (b) Alloue son tableau d'activation.
  - (c) Sauvegarde les registres *callee-saved*.

3. Pour l'appelé, à la fin de l'appel :
  - (a) Placer le résultat dans `%rax`
  - (b) Restaure les registres sauvegardés
  - (c) Dépile sont tableau d'activation
  - (d) Exécute `ret`
4. Pour l'appelant, juste après l'appel :
  - (a) Dépile les éventuels arguments
  - (b) Restaure les registres *caller-saved*

## Deuxième partie

# Cours 2 : 6/10

## Table des matières

<b>3</b>	<b>Sémantique Formelle</b>	<b>5</b>
3.1	Syntaxe Abstraite . . . . .	5
3.2	Sémantiques Useless . . . . .	6
3.2.1	Sémantique Axiomatique - Logique de Floyd-Hoare . . . . .	6
3.2.2	Sémantique Dénotationnelle . . . . .	6
3.2.3	Sémantique Par Traduction . . . . .	6
3.3	Sémantique Opérationnelle . . . . .	6
3.3.1	Sémantique Opérationnelle à Grand Pas . . . . .	6
3.3.2	Sémantique à Petits Pas . . . . .	7
3.3.3	Equivalence des Sémantiques . . . . .	7
3.3.4	Langages Impératifs . . . . .	8
<b>4</b>	<b>Interprète</b>	<b>8</b>

## Introduction

La signification des programmes est définie souvent de manière informelle, en langue naturelle, e.g. le langage Java.

## 3 Sémantique Formelle

La sémantique formelle caractérise mathématiquement les calculs décrits par un programme. C'est utile pour la réalisation d'outils (interprètes, compilateurs), et nécessaire aux raisonnements sur les programmes.

### 3.1 Syntaxe Abstraite

On ne peut pas manipuler un programme en tant qu'object syntaxique, on préfère utiliser la syntaxe abstraite (se déduit lors de la compilation à l'analyse syntaxique et sémantique.). On construit un arbre de syntaxe abstraite pour comprendre.

On définit la syntaxe abstraite par une grammaire. En OCaml, on réalise la syntaxe abstraite par des types construits. Il n'y a pas de parenthèses dans la syntaxe abstraite. On appelle sucre syntaxique toute construction de la syntaxe concrète qui n'existe pas dans la syntaxe abstraite.

C'est sur la syntaxe abstraite qu'on va définir la sémantique.

## 3.2 Sémantiques Useless

### 3.2.1 Sémantique Axiomatique - Logique de Floyd-Hoare

Tony Hoare, An axiomatic basis for computer programming, 1969, article le plus cité de l'histoire de l'informatique.

On caractérise les programmes par l'intermédiaire des propriétés satisfaites par les variables. On introduit le triplet  $\{P\} i \{Q\}$  qui signifie, si  $P$  est vraie avant l'instruction  $i$ , après,  $Q$  sera vraie

### 3.2.2 Sémantique Dénotationnelle

A chaque expression  $e$ , on associe sa définition  $\|e\|$  qui est un objet représentant le calcul désigné par  $e$ . On définit cet objet récursivement.

### 3.2.3 Sémantique Par Traduction

On définit la sémantique d'un langage en le traduisant vers un langage dont la sémantique est connue.

## 3.3 Sémantique Opérationnelle

La sémantique opérationnelle décrit l'enchaînement des calculs élémentaires qui mènent de l'expression à son résultat.

Il y a deux formes de sémantique opérationnelle :

1. La sémantique naturelle (big-steps semantics) :  $e \rightarrow v$
2. La sémantique par réduction (small steps semantics)  $e \rightarrow e_1 \rightarrow e_2 \rightarrow \dots \rightarrow v$

On applique ça à Mini-ML (qui est Turing-Complet).

### 3.3.1 Sémantique Opérationnelle à Grand Pas

On cherche à définir  $e \rightarrow v$ . On définit les valeurs parmi les constantes, les primitives non appliquées, les fonctions et les paires.

Une relation peut être définie comme la plus petite relation satisfaisant un ensemble d'axiomes notés  $\overline{P}$  et des règles d'inférences (implications). On définit  $\text{Pair}(n)$  par  $\overline{\text{Pair}(0)}$  et  $\frac{\text{Pair}(n)}{\text{Pair}(n+2)}$ . La plus petite relation qui vérifie ces deux propriétés coïncide avec la propriété «  $n$  est un entier naturel pair ».

Une dérivation est un arbre dont les noeuds correspondent aux règles et les feuilles aux axiomes (les arbres croissent vers le haut). L'ensemble des dérivations possibles caractérise exactement la plus petite relation satisfaisant les règles d'inférence.

**Définition 3.3.1.** On définit les variables libres d'une expression  $e$ , noté  $fv(e)$  par récurrence sur

$$\begin{aligned} fv(x) &= \{x\} \\ fv(c) &= \emptyset \\ fv(op) &= \emptyset \\ e \text{ avec : } \quad fv(\text{fun } x \rightarrow e) &= fv(e) \setminus \{x\} \\ \quad fv(e_1 e_2) &= fv(e_1) \cup fv(e_2) \\ \quad fv((e_1, e_2)) &= fv(e_1) \cup fv(e_2) \\ \quad fv(\text{let } x = e_1 \text{ in } e_2) &= fv(e_1) \cup (fv(e_2) \setminus \{x\}) \\ \quad fv(x) &= \{x\} \end{aligned}$$

On définit la substitution de toute occurrence libre de  $x$  dans  $e$  par  $v$  définie par :

$$\begin{aligned} x[x \leftarrow v] &= v \\ y[x \leftarrow v] &= y \text{ si } y \neq x \end{aligned}$$

On fait le choix d'une stratégie d'appel par valeur, i.e. l'argument est complètement évalué avant l'appel.

On a ici comme axiomes :  $\overline{c \rightarrow c}$ ,  $\overline{op \rightarrow op}$ ,  $\overline{(\text{fun } x \rightarrow e) \rightarrow (\text{fun } x \rightarrow e)}$  et comme règles d'inférences :

$$\frac{e_1 \rightarrow v_1 \quad e_2 \rightarrow v_2}{(e_1, e_2) \rightarrow (v_1, v_2)} \quad \frac{e_1 \rightarrow v_1 \quad e_2[x \leftarrow v_1] \rightarrow v}{\text{let } x = e_1 \text{ in } e_2 \rightarrow v}$$

et

$$\frac{e_1 \rightarrow (\text{fun } x \rightarrow e) \quad e_2 \rightarrow v_2 \quad e[x \leftarrow v_2] \rightarrow v}{e_1 \quad e_2 \rightarrow v}$$

On ajoute ensuite des règles pour les primitives, dépendant de la forme de chacune, e.g. :

$$\frac{e_1 \rightarrow + \quad e_2 \rightarrow (n_1, n_2) \quad n = n_1 + n_2}{e_1 \quad e_2 \rightarrow n}$$

Partant, on peut montrer qu'un programme s'évalue en une valeur en écrivant l'arbre de dérivation de celui-ci.

**Remarque 3.3.0.1.** *Il existe des expressions sans valeur :  $e = 12$  par exemple.*

On peut établir une propriété d'une relation définie par un ensemble de règles d'inférence, en raisonnant par induction sur la dérivation. Cela signifie par récurrence structurale.

**Proposition 3.3.1.** *Si  $e \rightarrow v$ , alors  $v$  est valeur. De plus si  $e$  est close alors  $v$  l'est également.*

*Démonstration.* Par induction :  $\frac{e_1 \rightarrow (\text{fun } x \rightarrow e) \quad e_2 \rightarrow v_2 \quad e[x \leftarrow v_2] \rightarrow v}{e_1 \quad e_2 \rightarrow v}$ . ■

**Proposition 3.3.2.** *Si  $e \rightarrow v$  et  $e \rightarrow v'$ , alors  $v = v'$ .*

*Démonstration.* Par induction. ■

**Remarque 3.3.0.2.** *On a donc défini une fonction plus qu'une relation.*

### 3.3.2 Sémantique à Petits Pas

La sémantique opérationnelle à petits pas remédie aux problèmes de programmes qui ne terminent pas, en introduisant une notion d'étape élémentaire de calcul  $e_1 \rightarrow e_2$

On commence par définir une relation  $\rightarrow^\varepsilon$  correspondant à une réduction en tête, au sommet de l'expression, par exemple :  $(\text{fun } x \rightarrow e) \quad v \rightarrow^\varepsilon e[x \leftarrow v]$  On se donne également des règles pour les primitives. On réduit en profondeur en introduisant la règle d'inférence :  $\frac{e_1 \rightarrow^\varepsilon e_2}{E(e_1) \rightarrow E(e_2)}$  où  $E$  est un

contexte défini par la grammaire suivante :

$$E ::= \square \begin{array}{l} | Ee \\ | vE \\ | \text{let } x = E \text{ in } e \\ | (E, e) \\ | (v, E) \end{array}$$

Un Contexte est un terme à trou où  $\square$  représente le trou.  $E(e)$  dénote le contexte  $E$  dans lequel  $\square$  a été remplacé par  $e$ . La règle d'inférence permet donc d'évaluer une sous-expression. Tels que définis, les contextes impliquent ici une évaluation en appel par valeur et de gauche à droite.

On note  $\rightarrow^*$  la cloture réflexive et transitive de  $\rightarrow$ .

**Définition 3.3.2.** *On appelle forme normale toute expression  $e$  telle qu'il n'existe pas  $e'$  telle que :  $e \rightarrow e'$*

### 3.3.3 Equivalence des Sémantiques

**Théorème 3.3.1.** *Les deux sémantiques opérationnelles sont équivalentes pour les expressions dont l'évaluation termine sur une valeur i.e. :*

$$e \rightarrow v \Leftrightarrow e \rightarrow^* v$$

*Démonstration.*

**Lemme 3.3.2** (Passage au contexte des réductions). *Supposons  $e \rightarrow e'$ , alors :*

1.  $ee_2 \rightarrow e'e_2$
  2.  $ve \rightarrow ve'$
  3.  $\text{let } x = e \text{ in } e_2 \rightarrow \text{let } x = e' \text{ in } e_2$
- $(\Rightarrow)$  On procède par induction sur la dérivation.  
 —  $(\Leftarrow)$  :

**Lemme 3.3.3** (Evaluation des Valeurs). *On a  $v \rightarrow v$ .*

**Lemme 3.3.4.** *Si  $e \rightarrow e'$  et  $e' \rightarrow v$  alors  $e \rightarrow v$ .*

*Démonstration.* On commence par les réductions de tête, puis on procède par induction aux applications de contexte. ■

On a alors, par récurrence sur le nombre de pas, l'implication souhaitée. ■

### 3.3.4 Langages Impératifs

Pour un langage impératif les sémantiques ci-dessus sont insuffisantes. On associe alors typiquement un état  $S$  à l'expression évaluée. L'état peut être décomposé en plusieurs éléments pour modéliser par exemple une pile (des variables locales), des tas...

## 4 Interprète

On peut programmer un interprète en suivant les règles de la sémantique naturelle. On se donne un type pour la syntaxe abstraite des expressions et on définit une fonction correspondant à la relation  $\rightarrow$

Un interprète renvoie la (ou les) valeur(s) d'une expression, souvent récursivement.

On peut éviter l'opération de substitution, en interprétant l'expression  $e$  à l'aide d'un environnement donnant la valeur courante de chaque variable (un dictionnaire). Ceci pose problème car le résultat de  $\text{let } x = 1 \text{ in fun } y \rightarrow +(x, y)$  est une fonction qui doit « mémoriser » que  $x = 1$ .

On utilise alors le module Map pour les environnements (c'est une *fermeture*). On représente alors la valeur d'une fonction avec son environnement.

Pour un interprète de la sémantique à petits pas, il vaut mieux utiliser un *zipper* que de recalculer le contexte tout le temps.

## Troisième partie

# Cours 3 : 13/10

## Table des matières

<b>5</b>	<b>Blancs</b>	<b>9</b>
<b>6</b>	<b>Outils pour l'analyse lexicale</b>	<b>9</b>
6.1	Expressions Régulières et Automates Finis . . . . .	9
6.1.1	Expressions Régulières . . . . .	9
6.1.2	Automates Finis . . . . .	10
6.2	Analyseur Lexical . . . . .	10
6.2.1	Principe . . . . .	10
6.2.2	Construction . . . . .	10



<b>7</b>	<b>L'outil ocamllex</b>	<b>11</b>
7.1	Analyseur Lexical en ocamllex . . . . .	11
7.2	Efficacité . . . . .	12
7.3	D'autres utilisations ocamllex . . . . .	12

## Introduction

L'objectif est de partir du code source, une suite de caractère, pour obtenir une suite de lexèmes plus compréhensible et simple à analyser syntaxiquement

## 5 Blancs

Les blancs (espace, retour chariot, tabulation) jouent un rôle dans l'analyse lexicale, car ils permettent de séparer deux lexèmes. De nombreux blancs sont inutiles e.g.  $x + 1$ , et seront ignorés. Les conventions diffèrent selon les langages et certains des caractères blancs peuvent être significatifs, par exemple l'indentation en python ou en Haskell, ou les retours chariots transformés en points-virgules comme en python. Les commentaires jouent le rôle de blancs.

## 6 Outils pour l'analyse lexicale

On va utiliser des expressions régulières pour décrire les lexèmes et des automates finis pour les reconnaître. On exploite en particulier la capacité à construire un automate partant d'une ou plusieurs expressions régulières.

### 6.1 Expressions Régulières et Automates Finis <sup>1</sup>

#### 6.1.1 Expressions Régulières

**Définition 6.1.1** (Syntaxe). *On définit la syntaxe des expressions régulières :*

$$r = \begin{array}{|l} \emptyset \\ \varepsilon \\ a \in \Sigma \\ r \cdot r \\ r + r \\ r^* \end{array}$$

**Définition 6.1.2** (Sémantique). *On définit alors la sémantique basée sur cette syntaxe par les langages rationnels :*

$$\begin{aligned} L(\emptyset) &= \emptyset \\ L(a) &= \{a\} \\ L(r_1 r_2) &= L(r_1) \cdot L(r_2) \\ L(r_1 + r_2) &= L(r_1) \cup L(r_2) \\ L(r^*) &= \bigcup_{n \in \mathbb{N}} L^n(r) \end{aligned}$$

Pour les constantes flottantes de CamL on a par exemple, si  $d = 0|1| \dots |9$  :

$$d \, d^* \, (.d^* \mid (\varepsilon \mid .d^*)(e \mid E)(\varepsilon \mid + \mid -) \, d \, d^*)$$

On peut alors écrire un algorithme pour savoir si une suite de caractère appartient à une regexp.

**Définition 6.1.3** (Dérivée de Brzozowski). *On pose, pour  $r$  une regexp et  $c \in \Sigma$  :  $\delta(r, c) = \{w \mid cw \in L(r)\}$*

---

1. Voir Cours de LFCC

### 6.1.2 Automates Finis

**Définition 6.1.4** (Syntaxe). *Un automate fini est un quintuplet  $(Q, \Sigma, I, F, T)$  où :*

- $Q$  est un ensemble fini d'états
- $\Sigma$  est un ensemble fini appelé alphabets
- $I \subseteq Q$  est un ensemble d'états initiaux
- $F \subseteq Q$  est un ensemble d'états finaux
- $T \subseteq Q \times \Sigma \times Q$  est un ensemble de transitions

**Théorème 6.1.1** (De Kleene). *Les expressions régulières et les automates finis définissent les mêmes langages.*

*Démonstration.* Voir Cours de LFCC ■

## 6.2 Analyseur Lexical

### 6.2.1 Principe

Un analyseur lexical est un automate fini pour la réunion de toutes les expressions régulières définissant les lexèmes. Le fonctionnement de l'analyseur lexical est différent de la simple reconnaissance d'un mot par un automate car :

- Il faut décomposer un mot (le source) en une suite de mots reconnus
- Il peut y avoir des ambiguïtés
- Il faut construire les lexèmes (les états finaux contiennent des actions)

### Ambiguïtés

Le mot `funx` est reconnu par l'expression régulière des identificateurs mais contient un préfixe reconnu par une autre expression régulière (`fun`) :  $\Rightarrow$  On choisit de reconnaître le lexème le plus long possible.

Le mot `fun` est reconnu par la regexp du mot clef `fun` mais aussi par celle des identificateurs :  $\Rightarrow$  On classe les lexèmes par ordre de priorité.

### Retour en arrière

Un analyseur va échouer sur l'entrée `abc` avec les trois regexp `a`, `ab`, `bc`. L'analyseur lexical doit donc mémoriser le dernier état final rencontré, le cas échéant.

Lorsqu'il n'y a plus de transition possible dans l'automate, de deux choses l'une :

- Soit aucun état final mémorisé : échec de l'analyse lexicale
- Soit on a lu le préfixe `wv` de l'entrée, avec `w` le lexème reconnu par le dernier état final rencontré : on renvoie `w` et l'analyse lexicale redémarre avec `v` préfixé au reste de l'entrée

En pratique, on va renvoyer dans l'analyseur lexical une fonction de calcul du prochain lexème, puisque l'analyse lexicale est faite pour l'analyse syntaxique

### 6.2.2 Construction

#### L'automate de Thompson

On construit par induction un automate compatible.

## L'automate de Berry-Sethi

On met en correspondance les lettres d'un mot reconnu et celles apparaissant dans la regexp : On distingue les différentes lettres de la regexp puis on construit un automate dont les états sont des ensembles de lettres. Pour construire les transitions de  $s_1$  à  $s_2$ , on détermine les lettres qui peuvent apparaître après une autre dans un mot reconnu : follow. Pour calculer follow, on a besoin de savoir calculer les premières et dernières lettres d'un mot reconnu (first et last). On a alors besoin d'une dernière notion : null, est-ce que  $\varepsilon$  appartient au langage reconnu. On obtient :

$\text{null}(\emptyset)$	=	false
$\text{null}(\varepsilon)$	=	true
$\text{null}(a)$	=	false
$\text{null}(r_1 r_2)$	=	$\text{null}(r_1) \wedge \text{null}(r_2)$
$\text{null}(r_1 \mid r_2)$	=	$\text{null}(r_1) \vee \text{null}(r_2)$
$\text{null}(r\star)$	=	true
On en déduit :		
$\text{first}(\emptyset)$	=	$\emptyset$
$\text{first}(\varepsilon)$	=	$\emptyset$
$\text{first}(a)$	=	$\{a\}$
$\text{first}(r_1 r_2)$	=	$\text{first}(r_1) \cup \text{first}(r_2)$ si $\text{null}(r_1)$
	=	$\text{first}(r_1)$ sinon
$\text{first}(r_1 \mid r_2)$	=	$\text{first}(r_1) \cup \text{first}(r_2)$
$\text{first}(r\star)$	=	$\text{first}(r)$
On définit last de même		
$\text{follow}(c, \emptyset)$	=	$\emptyset$
$\text{follow}(c, \varepsilon)$	=	$\emptyset$
$\text{follow}(c, a)$	=	$\emptyset$
$\text{follow}(c, r_1 r_2)$	=	$\text{follow}(c, r_1) \cup \text{follow}(c, r_2) \cup \text{first}(r_2)$ si $c \in \text{last}(r_1)$
	=	$\text{follow}(c, r_1) \cup \text{follow}(c, r_2)$ sinon
$\text{follow}(c, r_1 \mid r_2)$	=	$\text{follow}(c, r_1) \cup \text{follow}(c, r_2)$
$\text{follow}(c, r\star)$	=	$\text{follow}(c, r) \cup \text{first}(r)$ si $c \in \text{last}(r)$
	=	$\text{follow}(c, r)$ sinon

On construit alors l'automate reconnaissant  $r$  en ajoutant  $\#$  à la fin de  $r$  :

1. L'état initial est l'ensemble  $\text{first}(r\#)$ .
2. Tant qu'il existe un état  $s$  dont on doit calculer les transitions, pour chaque  $c$  de l'alphabet, on pose  $s'$  l'état  $\bigcup_{c_i \in s} \text{follow}(c_i, r\#)$
3. Les états acceptants sont ceux contenant  $\#$

## 7 L'outil ocamllex

### 7.1 Analyseur Lexical en ocamllex

Un fichier ocamllex porte le suffixe .mll et a la forme suivante :

- Code OCaml arbitraire
- $\text{rule } f_1 = \text{parse} \mid \text{regexp1} \{ \text{action 1} \}$
- ...
- Code OCaml arbitraire

A la compilation par ocamllex `file.mll`, on construit un fichier OCaml contenant des fonctions de types `Lexing.lexbuf  $\rightarrow$  type` où le type de sortie dépend de l'action dans la fonction.

Les regexp en ocamllex s'écrivent sous la forme :

<code>_</code>	n'importe quel caractère
<code>'a'</code>	le caractère 'a'
<code>"foobar"</code>	la chaîne "foobar" (en particulier $\varepsilon = ""$ )
<code>[caractères]</code>	ensemble de caractères
<code>[^caractères]</code>	complémentaire d'un ensemble de caractères
<code><math>r_1 \mid r_2</math></code>	alternative
<code><math>r_1 r_2</math></code>	concaténation
<code><math>r^*</math></code>	étoile
<code><math>r^+</math></code>	une ou plusieurs occurrences de $r$ , i.e. $rr^*$
<code><math>r^?</math></code>	au plus une occurrence de $r$
<code>eof</code>	fin du fichier

Pour remplacer la reconnaissance du lexème le plus long par celui le plus court, remplacer `parse` par `shortest`.

A longueur égale, c'est la règle qui apparaît en premier qui l'emporte.

On peut nommer la chaîne reconnue, ou des sous-chaînes reconnues par des sous-expressions régulières, à l'aide de la construction `as`.

On peut dans une action, rappeler récursivement l'analyseur lexical ou l'un des autres analyseurs simultanément définis. Le tampon d'analyse lexical doit être passé en argument, il est contenu dans une variable appelée `lexbuf`. Ceci est utile pour séparer les blancs, ou reconnaître les commentaires, même imbriqués.

Par défaut `ocamllex` construit l'automate dans une table interprétée à l'exécution, mais l'option `-ml` construit du code CamL pur.

## 7.2 Efficacité

Pour des raisons de stockage, et même en utilisant une table, l'automate peut prendre beaucoup de place. Il est donc préférable d'utiliser une seule expression régulière pour les identificateurs et les mots-clefs, puis de les séparer ensuite grâce à une table des mots-clefs.

On peut de même ne stocker que les caractères en minuscule pour être insensible à la casse.

## 7.3 D'autres utilisations `ocamllex`

On peut utiliser `ocamllex` pour :

1. Réunir plusieurs lignes vides consécutives en une seule
2. Compter les occurrences d'un mot dans un texte
3. Un petit traducteur OCaml vers HTML pour embellir le source mis en ligne (mettre les mots-clefs en vert, les commentaires en rouge, numéroter les lignes ...), le tout en moins de 100 lignes de code.