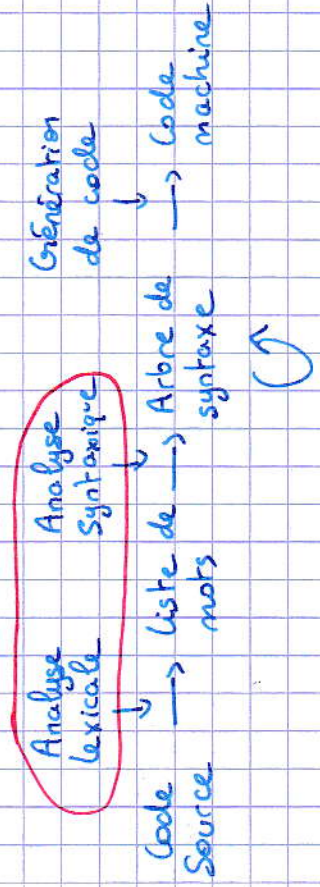


Leçon 25. Analyses lexicales et syntaxique. Applications.

Prérequis: automates, expressions régulières. OCaml.

I/Motivation: compilation

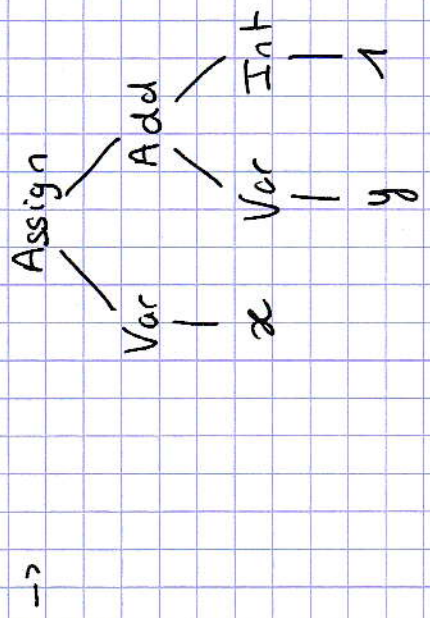


Analyses sémantiques (typage...)

Exemple:

let x = y + 1;

→ LET VAR("x") EQUAL VAR("y") PLUS INT(1)



II/Analyse lexicale

Objectif: reconnaître les mots du langage

1) Principe

Def: les lexèmes sont les mots du langage on les identifie par leur type, et éventuellement une valeur.

Ex: let → LET, x → VAR("x")

L'analyse repose sur la reconnaissance de motifs: chaque lexème est décrit par une expression régulière, dont les mots sont reconnus par un automate.

Def: règle d'analyse lexicale

Étant donné A l'alphabet des caractères, les règles sont de la forme:

$$e \rightarrow a$$

avec e une expression régulière sur A et a une action

NB: on peut simplement renvoyer le lexème, mais aussi renvoyer une autre information, modifier le texte...

Rq: une analyse lexicale gère une liste de règles, ~~une~~ pour chaque lexème. Si deux règles reconnaissent un mot, il y a ambiguïté. On considère donc ~~les~~ deux règles suivantes:

- x on choisira la règle reconnaissant le mot le plus long
- x en cas d'égalité, on suit l'ordre de priorité, prédéfini, des règles

2) Algorithme

Data: $(e_i \rightarrow a_i) \in E \rightarrow$ la liste des règles, ordonnée pour chaque e_i , A_i l'automate associé

lexeme-reconnu := -1
initialisation des A_i

Tant qu'il existe un A_i non bloqué

On fait lire le caractère suivant à tous les A_i

Si un A_i atteint l'état final

lexeme-reconnu := i

// si plusieurs A_i atteignent l'état final, on garde

// celui d'indice supérieur

// tous les automates sont bloqués

Si lexeme-reconnu == -1

erreur lexical, le texte ne correspond à aucun lexème

Si non

// lexeme-reconnu correspond au lexème le plus long

On applique a_i

on reprend à la suite de ce lexème

3) Applications

Coloration syntaxique, filtration des commentaires, reconnaissance de mots... → Développement

Exemple

ocamllex

III/ Analyse syntaxique

Objectif: reconnaître si un texte est reconnu par la grammaire
→ si oui, générer l'arbre de syntaxe abstraite

1) Contexte

Def: une grammaire est un ensemble de règles

$A \rightarrow a_1 \dots a_n$

$A \in N$: non-terminal, $a_i \in N \cup T$: non-terminal ou terminal

Ex: $S \rightarrow \text{let } \text{var} = E$

$E \rightarrow \text{var } \text{int} \mid E + E$

$(N = \{S, E\}, T = \{\text{let}, \text{var}, =, \text{int}, +\})$

Def: une dérivation est l'application successive de règles

Ex: $S \rightarrow \text{let } \text{var} = E$

$\rightarrow \text{let } \text{var} = \underline{E + E}$

$\rightarrow \text{let } \text{var} = \text{var} + \underline{E}$

$\rightarrow \text{let } \text{var} = \text{var} + \text{int}$

$\Rightarrow \text{let } x = y + 1$ est reconnu par cette grammaire.

2) Analyse descendante

Principe : en partant de la première règle (l'axiome), retrouver les règles à appliquer pour générer le texte

NB : à ne pas confondre avec l'analyse ascendante, qui part du texte et cherche à revenir à l'axiome

Problème : si on a $S \rightarrow a|b$, comment savoir quelle règle appliquer ?

Solution : on regorde les caractères suivants, c'est l'analyse $LL(k)$. Ici, on se contente de $LL(1)$

3) Analyse $LL(1)$

On utilise une pile pour mémoriser où on en est des règles appliquées.

Si il y a cET au sommet de la pile, et aET en tête du texte, et $a=c$, on avance dans la lecture (si $a \neq c$ c'est une erreur syntaxique)

Si non, on regarde dans le tableau Π qui indique pour chaque XEN , et cET, s'il y a une règle $X \rightarrow \beta$, où β commence par c.

Quand la pile est vide, soit il ne reste plus que le symbole de fin (ici #) et c'est bon, soit on a une erreur syntaxique.

Construction du tableau

$$T(X, c) = \begin{cases} \beta & \text{si } X \rightarrow \beta \text{ et } c \in \text{premier}(\beta) \\ \gamma & \text{si } X \rightarrow \gamma, \text{null}(\gamma) \text{ et } c \in \text{suivent}(X) \end{cases}$$

$$\begin{aligned} \text{avec } \text{null}(XEN) &= \text{true ssi } a \rightarrow s^*E \\ \text{premier}(XEN) &= \{a \in T \mid \exists w. a \rightarrow s^*aw\} \\ \text{suivent}(XEN) &= \{a \in T \mid \exists u, w. S \rightarrow s^*uXaw\} \end{aligned}$$

NB : la grammaire est $LL(1)$ si pour chaque case il y a au plus une règle, i.e pas d'ambiguïté

Ex : Développement 2 : analyse $LL(1)$ pour langage arithmétique

IV/ Applications

En tant que première étape de la compilation :

AST \rightarrow code machine
analyse de typage de l'AST

Pour l'interprétation : code source \rightarrow résultat

Pour la génération d'images :

LaTeX \rightarrow PDF

svg \rightarrow image

Pour l'affichage : HTML \rightarrow affichage dans un navigateur

Du vertures : Analyse ascendante, OCaml yacc
 \rightarrow Génération de code et optimisation
 \rightarrow Typage