

NOM : PECATTE

Prénom : Timothée

Jury :

Algèbre ← Entourez l'épreuve → Analyse

Sujet choisi : 923 - Analyses lexicales et syntaxique : applications

Autre sujet :

I. Introduction à la compilation :

But: Pouvait écrire des programmes dans un langage de haut niveau.

Ex 1: TANT QUE, SI... AUTRS, appels de fonctions...

Def 2: (Interpréteur) Un interpréteur peut être vu comme une fonction qui prend en entrée un programme écrit dans un langage L et des données, et qui renvoie des données ou une erreur.

$$I_L : L \times D^* \rightarrow D^* \cup \{\text{erreur}\}$$

Rmq 3: Il faut donc concevoir un interpréteur pour chaque langage.

Def 4: (Compilateur) Un compilateur d'un langage L vers un langage L' agit comme un "traducteur" qui est fidèle vis-à-vis des interpréteurs, i.e. : $\forall d \in D^*, I_{L'}(P, d) = I_L(C_{L \rightarrow L'}, d)$

Rmq 5: Ainsi, si on possède un interpréteur pour un langage de référence, il suffit de concevoir un compilateur vers ce langage.

Def 5: Un compilateur s'organise en deux parties principales :

- Le front-end, qui vérifie que l'entrée est bien donnée et qui reconstruit le programme.
- Le back-end, qui gère le code vers le langage cible, en effectuant éventuellement des optimisations.

Qeshors: Comment passer d'un mot qu'on a vu dans un programme à une représentation sémantique de celui-ci ? Comment détecter les erreurs ? Peut-on les corriger ?

II. L'Analyse lexicale d'un programme :

II.1. Bot de l'analyse lexicale :

Def 7: (Symboles) Dans un langage de programmation, les symboles, ou unités lexicales, sont les différents types d'objets qu'une suite de caractères peut représenter.

Ex 8: Identificateur, parenthèses, mots-clés (SI, TANT QUE), chiffres, commentaires, =, <, >, :=, ...

But: Étant donné une chaîne de caractères en entrée, le rôle de l'analyse lexicale est de découper l'entrée en symboles et de garder ceux qui sont significatifs (les séparateurs sont ignorés).

Rmq 1: Au point de vue d'un compilateur, un programme et son analyse écrit sur une seule ligne sont exactement les mêmes (sauf pour l'effacement). Cela permet donc des librairies quant à l'industrialisation des programmes.

II.2. Bases théoriques :

Def 10: (Expressions régulières) Les expressions régulières sont le plus petit ensemble contenant \emptyset, ϵ, a pour $a \in \Sigma$, et stables par \cup, \cap, \cdot et \star . On associe par induction un langage à une expression régulière.

Ex 11: $L((a|b)^* @ a|b)^* \text{com}$ est l'ensemble des adresses mail ne contenant que a et b comme lettres (à part le "@" et le com)

Prop 13: L est régulier ssi L est reconnu par une expression régulière.

Rmq 13: Les expressions régulières sont plus faciles, mais les autres sont plus pratiques pour tester l'acceptation.

Rmq 14: L'automate obtenu via la preuve est non-déterministe mais on peut le rendre déterministe

Prop 15: Etant donné une expression régulière, on peut construire de manière effective un automate déterministe minimal qui reconnaît la même langage.

Rmq 16: Les automates contextuels jusqu'à présent permettent seulement de tester l'appartenance, alors que le rôle de l'analyse lexicale est de produire une suite de symboles ρ qui donne naissance le comportement de l'automate.

Def 17: Un automate généralisé est un 5-uplet (Q, q_0, F, S, δ) avec

- Q est l'ensemble des états
 - q_0 est l'état initial
 - F est l'ensemble des états final
 - S est la fonction de transition
 - δ est la fonction de production
- À l'état q , sur l'écriture a , si $\delta(q, a)$ existe, l'automate passe dans l'état $\delta(q, a)$. Sinon, si $q \in F$, on produit $\epsilon(q)$ et on passe dans l'état $\epsilon(q, \epsilon)$. Sinon, on arrête la lecture pour renvoyer un erreur.

Ex 18: tuxex

Rmq 19: L'automate cherché ainsi à reconnaître la plus grande partie possible comme un symbole, ce qui autorise la dérivation de variables de la forme STTE.

II.3. Souffrance paritaire:

Rmq 20: En paritaire, ρ est labellisé d'écrire $\epsilon_1 \epsilon_2 \epsilon_3 \epsilon_4 \epsilon_5 \epsilon_6 \epsilon_7 \epsilon_8$ à chaque fois qu'on voit parité d'un caractère numérique.

Def 21: (Classe de caractère). Une classe de caractère est un identifiant qui décrit un ensemble de caractère fini.

Ex 22: $ca = a-zA-Z$, $ch = 0-9$, $oo = |$, $opar = (, \dots$

Def 23: (Séquence d'expressions régulières). Une séquence d'expressions régulières sur Σ est une suite de définition

$$\begin{aligned} A_1 &= R_1 \\ A_2 &= R_2 \\ &\vdots \\ A_n &= R_n \end{aligned}$$

et les R_i sont expr régulières sur $\Sigma \cup \{A_1, \dots, A_{i-1}\}$

Def 24: Opération "ouf" $R_1 \text{ ouf } (R_2) \sim R_1 \overline{R_2} R_2 \overline{R_1}$

Rmq 25: Les opérations, et d'autres, sont autorisées dans le générateur automatique Flux d'automate généralisé pour analyser lexicale

Def 1: Construction de l'analyse lexicale pour l'exemple suivant

$$ch = 0-9, \text{ hex} = A-F,$$

$$\begin{aligned} h_{\text{hex}} &= ch \cdot ch^* \\ h_{\text{hex}} &= h(ch | hex)^* \\ \text{non-recursif} &= \text{intrecursif} \cdot \text{intrecursif} (ch | hex) \end{aligned}$$

Def 19: Construction directe d'un automate déterministe à partir d'une expr régulière.

III. L'analyse syntaxique d'un programme:

III.1. But de l'analyse syntaxique:

Rmq 26: Après le passage de l'analyse lexicale, le programme est encore linéaire et donc peu propice à manipuler

But: Structurer le programme et détecter les erreurs éventuelles de syntaxe

Rmq 27: C'est aussi le rôle de l'analyse syntaxique de vérifier que certaines règles imposées par le langage sont vérifiées, comme par exemple la déclaration des variables avant leur usage, mais cela ne sera pas abordé dans la suite.

III.2. Grammaires hors-contexte et parsers automatiques à pile:

Def 28: Une grammaire non contextuelle est un quadruplet (N_0, V_1, P, S) où V_1 est l'ensemble des symboles non terminaux, V_1 est l'ensemble des terminaux, $P \subseteq V_1 \times (V_1 \cup V_1^*)^*$ est l'ensemble des productions, et $S \in V_1$ est le symbole de départ

Ex 39: $G_0 = (\{G, T, F\}, \{T, *, ()\}, \{ \in \rightarrow \text{ext } T, T \rightarrow T * F | F, F \rightarrow () | \text{id} \}, E)$

Def 30: (Dérivations). $\varphi \xrightarrow{G} \psi$ s'il existe des mots σ, τ, κ et un non-terminal $A \neq \epsilon$ tel que $\varphi = \sigma A \tau$ et $\psi = \sigma \kappa \tau$ et $A \rightarrow \kappa \in P$.

On définit $\xrightarrow{*}$ comme la closure transitive réflexive de \xrightarrow{G}

Def 31: (Arbre syntaxique). Un arbre syntaxique est une représentation d'une dérivation qui ne tient pas compte de l'ordre dans lequel les non-terminals sont remplacés. Il est défini par induction sur la dérivation.

Ex 32: Anorex

Def 33: (Ambiguïté). Une phrase $x \in L(G)$ est ambiguë si elle admet plusieurs arbres syntaxiques. Une grammaire est ambiguë lorsqu'elle contient au moins une phrase ambiguë.

Rmq 34: L'ambiguïté d'une grammaire n'est pas forcément générale aussi, mais induit du non-déterminisme que nous devons résoudre en faisant un choix.

Def 35: (Dérivation gauche). Une dérivation est dite dérivation gauche si à chaque étape on remplace le non-terminal le plus à gauche.

Prop 36: Un langage est algébrique si il est reconnu par un automate à pile.

Rmq 37: Contrairement aux automates simples, on ne peut pas déterminer en un nombre \bar{n} de pile dans le cas général, et donc on aura pas de manière efficace pour des grammaires algébriques quelconque.

III.3. Grammaires LL(k) et analyse syntaxique descendante:

Rmq 38: Si on autorise l'automate à Pile les k prochains symboles au lieu de 1, on ne change pas la complexité de la reconnaissance, mais cela peut rendre l'automate déterministe.

Def 39: (PRETIER) On définit PRETIER(X) comme étant l'ensemble des caractères qui commencent une chaîne dérivée à partir de X.

Ex 40: PRETIER(T) = { (, id }

Def 41: (SUIVANT). On définit SUIVANT(X), pour X non-terminal, comme l'ensemble des caractères qui peuvent suivre immédiatement X dans une phrase.

Ex 42: SUIVANT(T) = { *, } }

Rmq 43: On peut définir de la même manière les fonctions PRETIER et SUIVANT.

Prop 44: PRETIER et SUIVANT peut être calculé de manière efficace

Def 45: (Grammaire LL(k)). Une grammaire LL(k) est une grammaire telle que pour toute règle de production $A \rightarrow \alpha | \beta$, on a : * PRETIER(α) \cap PRETIER(β) = \emptyset

* si $e \in$ PRETIER(β), alors PRETIER(α) \cap SUIVANT(β) = \emptyset

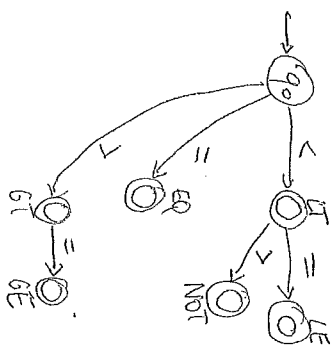
Rmq 46: On peut définir de même les grammaires LL(k) avec PRETIER et SUIVANT.

Prop 47: Si G est une grammaire LL(k), alors on peut construire un automate déterministe PDL(k) qui construit l'arbre syntaxique du mot donné en entrée ou qui renvoie erreur.

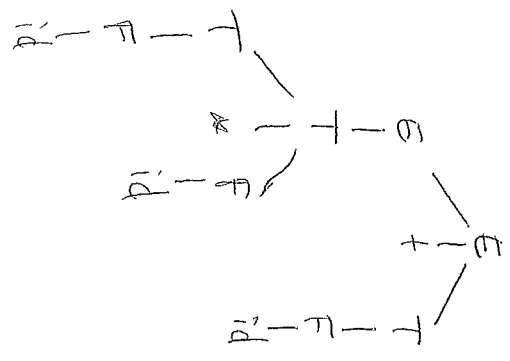
Rmq 48: Cet automate déterministe par avoir des états qui ne sont pas générés car l'automate dérive de la transition en fonction du prochain caractère.

Def 49: (Table prédictive). Une table prédictive, pour une grammaire LL(k), est un tableau dans lequel pour chaque non-terminal E, et chaque symbole \bar{a} visible possible, on a l'unique règle de production $E \rightarrow \alpha$ tq $\alpha \rightarrow \bar{a}$ ou si $\omega \in$ PRETIER(β) et "erreur" si $\omega \notin$ PRETIER(E)

Def 50: Construction de l'automate prédictif d'une grammaire LL(k) par une simple.



Exemple 18: Automate généralisé



Exemple 32: Ambiguïté de la dérivation (dans G0).

- | | | |
|---------------------------|---------------------------|---------------------|
| $S \rightarrow E$ | $E' \rightarrow +E$ | $T' \rightarrow *T$ |
| $E \rightarrow TE'$ | $T \rightarrow FT'$ | $F \rightarrow (E)$ |
| $E' \rightarrow \epsilon$ | $T' \rightarrow \epsilon$ | $F \rightarrow id$ |

Exemple de grammaire LL(1)