

→ Compilons. Aho Sethi-Ullman

I. COMPLICATION

Un compilateur est un programme qui lit le code d'un autre programme dans son langage (langage source) et le traduit dans son autre langage (langage cible). Il peut en plus effectuer des actions dans le programme source.

Dans cela, dans une première partie d'analyse, un compilateur nomme la programme source en différentes composantes (phase d'analyse lexicale) pour leur impose une structure grammaticale afin d'obtenir une représentation intermédiaire au programme (phase d'analyse syntaxique), en fin, il vérifie le sens de la représentation (phase d'analyse sémantique). En même temps, le compilateur crée une table de symboles pour stocker les informations obtenues.

II. ANALYSE LEXICALE

A. Vocabulaire

Un analyseur lexical lit donc les caractères du programme source. Les regroupes en chaîne de caractères ayant un sens, puis renvoie la composante associée, on ajoutant si besoin une entrée dans la table de symboles.

Def 1: les chaînes de caractères ayant un sens sont les lexèmes, les composantes sont les unités lexicales.

Ex 1: double = nombre * 2 ← programme source
 "double" " = " "nombre" "*" "2" ← décomposition
 < id, 2 > < = > < id, 2 > < * > < int, 2 > ← suite d'unités lexicales
 avec la première entrée de la table de lexicales
 symboles qui est "double" et la seconde "nombre".

L'analyseur lexical supprime donc au passage les caractères inutilisés par les autres comme les blancs et commentaires.

B. Principe

De manière générale, les lexèmes peuvent être reconnues par des expressions régulières. Le motif associé à l'unité lexicale Ex 3: lexème 42, motif: (i+)(\d+)+ avec diff. -{C, ..., 9}

unité lexicale: < int, 42 >
 A chaque unité lexicale et donc motif, on peut associer un automate fini déterministe. S, {q0, ..., qn}, où est l'ensemble des motifs, on obtient {q0, ..., qn}.

On peut alors construire un automate pour reconnaître tous les motifs, grâce à des é-transitions (Figure 1)

C. Résolution des conflits

Ab1: la détermination de l'automate obtenu est coûteux.

Ab2: on peut savoir choisir le plus long préfixe ← en fait si l'on dit

Ex 6: 12 doit être reconnu comme < int, 12 > et non < int, 1 > < int, 2 >

Ab 7: certains lexèmes peuvent correspondre à plusieurs unités lexicales, notamment au niveau des identifiants et mots-clés.

Il faut établir une priorité entre les unités lexicales.
 Ex 8: il doit être reconnu comme le mot id < id > et non < id, 2 > où la table de caractères contient dans la ligne 2 " id".

D. Applications et utilisation

Appl 9: notation syntaxique. Les programmes comme emacs calculent au fur et à mesure de ce qui est dans une certaine correspondance et si l'unité lexicale.

On peut alors particulièrement bien voir le système de priorités.
 id > → else → vison ...
 id, > < else > < id, >

Appl 10: certains programmes font une analyse lexicale simplifiée. Comme les compilateurs programmables par des cellules.
 mels ils doivent être choisis depuis son menu.

Appl 11: lex est un constructeur d'analyseur lexical: en lui donnant en entrée les unités lexicales et motifs associés, ainsi que les priorités entre les unités lexicales, il crée un analyseur lexical.

Ex 12: " + " { return Plus; }
 { [C - 9] } + { g } val num = calculer (h);
 extrait de programme Java permettant de reconnaître la symbole + et les entiers.

on intervient l'ordre de priorité

III. ANALYSE SYNTAXIQUE

A. Principe

Une analyse syntaxique est donc la suite d'unités lexicales renvoyées par l'analyseur lors d'un arbre dérivation qui elle peut être générée par la grammaire du langage source. Si ce n'est pas le cas, il doit renvoyer l'erreur la plus précise possible. Si c'est le cas, il doit construire l'arbre de la dérivation.

Rq 13. Les langages de programmation ont généralement une grammaire algébrique $G = (Z, V, R, St$

Notation (14) 4 : L'algorithme de Cocke-Young-Kasami utilise la programmation dynamique pour décider si la suite d'unités lexicales est reconnue par la grammaire, après que la grammaire a été mise sous forme normale de Chomsky.

Def (15) 4 : On peut se contenter de mettre la grammaire en 2NF : alors l'algorithme CYK se passe en temps $O(m^3|G|)$. En pratique, ce n'est pas assez efficace. \leftarrow voir exercice 1

Notation (15) 2 : Analyse descendante : On construit l'arbre de dérivation à partir du symbole initial de G .

Notation (17) 3 : Analyse ascendante : On construit l'arbre de dérivation à partir des feuilles : les unités lexicales.

Rq 18 : Il faut choisir en plus, pour des deux dérivées méthodes dans quel ordre lire la suite d'unités lexicales (gauche à droite ou droite à gauche). Sans la suite, de gauche à droite (L).

B. Analyse syntaxique descendante

Ex 19 Pour la grammaire $G = E \rightarrow E+T \mid T; T \rightarrow T * F \mid F; F \rightarrow (E) \mid id$, quel algorithme est le plus efficace ? $\{+, *, id, (,)\}$, la famille $\{id\}$ a pour arbre de dérivation les figures à :

Ph 20. Dans problèmes se posent alors : Quelle règle choisir ? Quels termes de la suite choisir pour la part à la règle ?

Algo 21 : Pour dans les non terminaux A on définit : pour un dérivé

Prop. A (1)

Choisir une règle $A \rightarrow X_1 \dots X_k$.

Pour $x = A$ et B

Si X_i est un non terminal, alors $Der(X_i)$

Si $x = (P]$, alors $P \in P \rightarrow A$

Si $x = E$, alors

Ph 22 : On peut avoir mes choix une règle $A \rightarrow X_1 \dots X_k$ dans ce cas, il faut revenir en arrière et en choisir une autre.

Def 23 : Une grammaire est dite LL(k), $k \geq 0$ si il est possible de faire une analyse syntaxique ^{descendant} qui lit k symboles en entrée et n'a jamais à faire le retour en arrière.

Def 24 : Premier (a) $= \{a \in \Sigma, \exists u \in \Sigma^*, u \neq \epsilon \text{ such } u \in E, u \rightarrow \epsilon\}$ et IT est défini pour $\alpha \in (E \cup V)^*$ et est calculé récursivement. Pour $A \in V$, $Succant(A) = \{a \in \Sigma, \exists u \in (E \cup V)^*, S \rightarrow^* uAu \text{ et } u \in IT\}$ et IT est calculé récursivement de droite à gauche. $S \rightarrow^* uAu$ IT est calculé récursivement de droite à gauche. $A \rightarrow \alpha \mid \beta$ de G , on a :

Prop 25. Une grammaire est LL(1) si pour toute A non terminale ayant 2 règles dont si est la partie gauche. $A \rightarrow \alpha \mid \beta$ de G , on a :

(a) $Premier(\alpha) \cap Premier(\beta) = \emptyset$

(b) $\epsilon \in Premier(\alpha) \Rightarrow Premier(\beta) \cap Succant(A) = \emptyset$

(c) $\epsilon \in Premier(\beta) \Rightarrow Premier(\alpha) \cap Succant(A) = \emptyset$

On peut alors vérifier les implications dans un tableau.

Def 26 : La table d'analyse syntaxique d'une grammaire LL(1) est le tableau à double entrée $M[A, a]$ où $A \in V$ et $a \in \Sigma \cup \epsilon$. $M[A, a]$ contient les α tq $A \rightarrow \alpha$ et $a \in Premier(\alpha)$ (2) ou $\epsilon \in Premier(\alpha)$ et $a \in Succant(A)$ (1)).

Prop 27. Une grammaire est LL(1) si chaque case ne contient que un plus une production, les cases vides correspondant aux erreurs

Def (28) 2 : $G = E \rightarrow TE'$ d'axiome E est LL(1)

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

C. Analyse syntaxique ascendante

Ex 22: On analyse syntaxique une ascendantement correspond à un préfixe de lecture (si l'on ajoute des lettres avant d'appliquer une règle) et de réduction (si l'on a les et réduit jusqu'à

puissent correspond exactement aux numbre droit d'une règle, qui on remonte dans de droite à gauche).

On obtient ainsi une analyse LR. Redonne L de gauche à droite, R. On applique les règles de droite à gauche.

Pb 29: A nouveau, deux problèmes se posent:
 Dat - on lit ce à redonne ?
 Si on veut de redonne, quelle point de dat - on redonne ?

Def 30: Une grammaire est dite LR(LR), $k \geq 0$ si il est possible de faire un analyseur syntaxique ascendant qui lit le symboles en avance, et n'ajoute à peine de retour en arrière (pas d'erreur).

La situation qui apparaît directement est d'automate à pile, où la pile contient s'état actuel des lectures, réductions.

Def 31: Un item est un triplet $[A \rightarrow \alpha \cdot \beta]$ où $A \rightarrow \alpha \beta$ est une règle de la grammaire et $\alpha, \beta \in (\Sigma \cup \{ \epsilon \})^*$.

Def 32: On peut alors construire d'automate fini déterministe des items, ainsi que d'automate à pile LR(LR) pour la grammaire. $G: S \rightarrow SS^+ d$ comme S est déterministe $\{ \alpha, \beta \}$.
 $S \rightarrow \alpha \beta$

Def 33: Une grammaire est LR(0) si d'automate à pile associé est déterministe.

Pour une grammaire LR(0), on peut à nouveau construire un tableur regroupant les informations.

Def 34: On définit la grammaire avec une nouvelle règle $S' \rightarrow S$ avec $S' \in \Sigma \cup \{ \epsilon \}$ d'origine de la grammaire.

Def 35: À un ensemble d'items, on lui associe sa fermeture.

calculer récursivement: Pour tout item $[A \rightarrow \alpha \cdot \beta]$ avec $A, B \in N, \alpha, \beta \in (\Sigma \cup \{ \epsilon \})^*$, on ajoute l'ensemble $\{ [B \rightarrow \alpha \cdot \gamma], B \rightarrow \alpha \gamma \text{ une règle} \}$ construit récursivement les autres états si q est un état on obtient les états $q \cdot a, a \in (\Sigma \cup \{ \epsilon \})$ où $q \cdot a$ est la fermeture de $\{ [A \rightarrow \alpha \cdot a \cdot \beta], [A \rightarrow \alpha \cdot a \cdot \beta] \} \cap q$.

Def 36: Le premier état est la fermeture de $\{ [S' \rightarrow \cdot S] \}$. On obtient les états $q \cdot a, a \in (\Sigma \cup \{ \epsilon \})$ où $q \cdot a$ est la fermeture de $\{ [A \rightarrow \alpha \cdot a \cdot \beta], [A \rightarrow \alpha \cdot a \cdot \beta] \} \cap q$.

Def 37: La table d'analyse syntaxique LR(0) est alors la table $M[q, a]$ où q est un état de $\Sigma \cup \{ \epsilon \}$. Elle contient alors:
 (i) $q \cdot a \in \Sigma$: lecture (q') si on doit faire une étape de lecture, écrire q' sur la pile, et aller à l'état q' .
 (ii) $q \cdot a \in \Sigma$: réduction (r) si on doit faire une étape de réduction par la règle r : on doit alors dépiler tout ce qui est à droite de la règle ($A \rightarrow \alpha$) dans le pile, puis aller en $N[q', A]$ où q' est le nouvel élément en haut de pile.

Prop 38: Une grammaire est LR(0) si lorsque on a une réduction dans la case $N[q, a]$ pour un $a \in \Sigma$, alors $N[q, b], b \in \Sigma$ contient la même réduction. De plus chaque case ne contient que au plus une action.

D. Applications et utilisation
Appl 39: les navigateurs internet utilisent un analyseur syntaxique pour analyser le HTML des pages
Ex 40: Fac est un constructeur d'analyseur syntaxique: on lui donne en entrée la grammaire du langage source, il crée son analyseur syntaxique.

Ex 41: Le langage $E \rightarrow E + E$ se réduirait par un programme yacc, par exemple en:
 expr : expr '+' expr { \$\$ = mtd('+', \$1, \$3); }

Ex 42: Pour la grammaire de l'exemple 19, la formule $id + id$ a pour arbre de dérivation la Figure 3)

Figure 1 Assemblage de l'automate pour l'analyse lexicale:



Figure 2 Analyse syntaxique descendante de id + id:

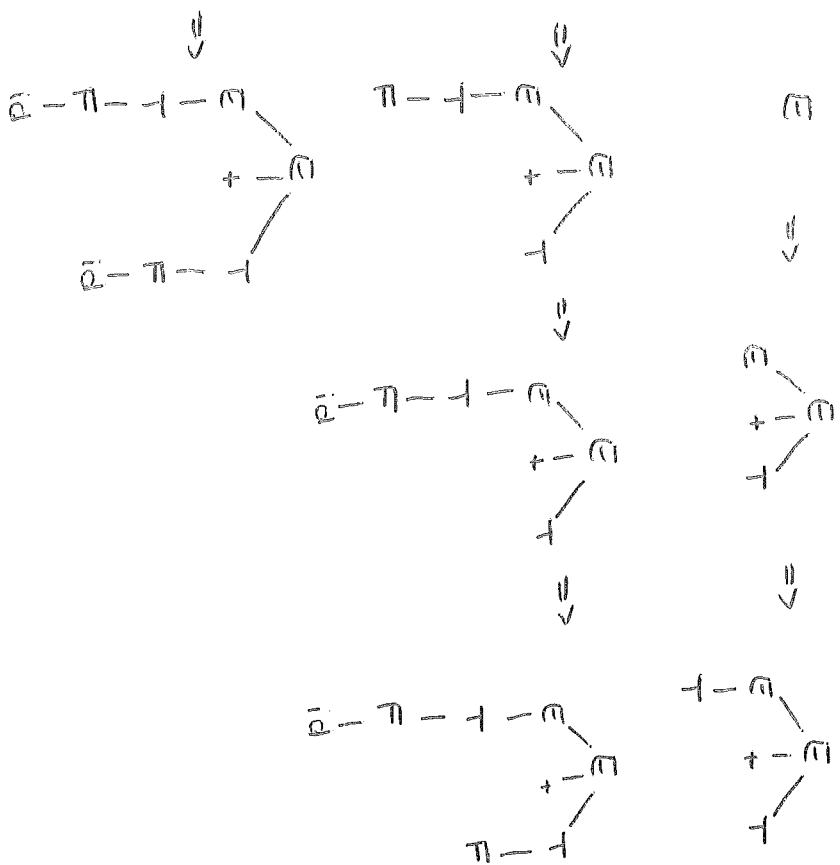


Figure 3 Analyse syntaxique ascendante de id + id

