

0662 / Applications 2015
 Analyse lexicale et syntaxique - Applications
 Leçon 323

I. Introduction à la compilation (voir Figure 1)

Définition 1 Un compilateur est un programme qui lit le code d'un autre programme (en langage de programmation) et le traduit en langage machine (langage cible). Il peut également détecter des erreurs dans le programme source.

Remarque: La compilation s'approche par certains aspects de la conversion de fichiers. Certains outils seront réutilisés ailleurs.

1) Analyse lexicale

L'analyseur lexical regroupe les caractères en séquences appelées lexèmes, et pour chaque lexème, transmet une unité lexicale à l'analyseur syntaxique. Il supprime les éléments inutiles.

Exemple: en langage C, `let repose = 42;` (x pour n'importe quelle position x)

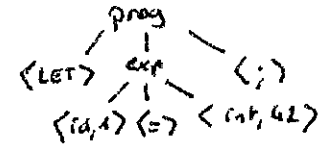
`<LET> <id, 1> <=> <int, 42> <;>`

↑
position dans la table des symboles

2) Analyse syntaxique

L'analyseur syntaxique repère la structure du programme : il construit l'arbre de dérivation correspondant.

Exemple 3



3) Analyse sémantique

L'analyseur sémantique effectue des vérifications : les types, les variables non déclarées, ...

⚠ Si on parle d'assembleur ⇒ savoir donner des exemples!
 Autres langages cibles?
 Exemple - t-on toujours un programme pour l'utiliser?
 Erreurs de conception!

II Analyse lexicale

Définition 4: - Une unité lexicale est constituée d'un nom et d'une valeur d'attribut. Le nom représente le type de l'unité : mot clé, identificateur, opérateur, entier, ...

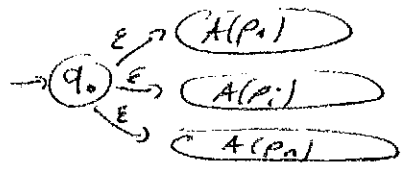
- Un motif est une expression régulière décrivant la forme que peut prendre une occurrence d'une unité lexicale.
- Un lexème est une séquence de caractères dans le programme source reconnu par un motif, donc correspondant à une unité lexicale.

Exemple 5:

unité lexicale	int	float
motif	$(- + E) \text{ chiffre}^+$	$(- + E) (\text{chiffre}^+ (\cdot + E) + \text{chiffre}) \text{chiffre}^*$ $(E + E (- + E) \text{chiffre}^+)$
lexèmes	2, -27	.4, -1.27E12, 10E-14

Construction de l'analyseur lexical:

- On associe un automate à chaque motif $P_i: A(P_i)$
- On combine ces automates avec des ϵ -transition:



L'analyse lexicale fait une lecture du code en suivant plusieurs chemins en parallèle. (l'automate est non déterministe)

Gestion des conflits:

- Détermination de l'automate → courtage (voir fig 2)
- Choix du plus long préfixe
- ordre de priorité sur les sorties : les mots clés sont prioritaires par rapport aux identificateurs, ...

requi: Certains outils pour l'écriture de programme font de l'analyse lexicale à mesure que le texte est tapé: emacs ...
 on voit alors bien l'utilisation d'automate (le → let → l'histoire)

- Certains outils font une analyse lexicale très simplifiée: sur certains calculateurs programmables, les mots clés sont accessibles via un menu.

Outil pour la construction de l'analyseur lexical: LEX/FLEX.

LEX (dans sa version plus récente FLEX) est un constructeur d'analyseur lexical, un compilateur de compilateur en somme. Il crée, à partir des motifs et des unités lexicales données en entrée, ainsi que des règles de priorité, un analyseur lexical.

III Analyse syntaxique

1) Cadre d'étude

Les langages de programmation sont généralement descriptibles par des grammaires algébriques hors contexte, G.
 L'objet de l'analyseur syntaxique est de vérifier que la suite d'unités lexicales appartenant bien au langage L(G) et de créer si tel est le cas l'arbre de dérivation correspondant, qui sera ensuite transformé en arbre de calcul.

1^{ère} Idée: Algorithme de Cocke, Younger et Kasami ou récrite la grammaire sous forme normale de Chomsky, puis on effectue une analyse utilisant la programmation dynamique. C'est peu efficace (O(n³)) et non utilisé en pratique

2^{ème} Idée: Analyse descendante: on part du symbole initial de G et on reconstruit les étapes pour arriver à la chaîne ⇒ Analyse LL

3^{ème} Idée: Analyse ascendante: on lit la chaîne et on remonte les étapes jusqu'à arriver au symbole initial de G ⇒ Analyse LR
 Les 2^{ème} et 3^{ème} idées sont utilisées en pratique.

2) Analyse descendante

Les nœuds de l'arbre de dérivation sont construits selon un parcours préfixe.

Exemple; Pour la grammaire G: $E \rightarrow E + T \mid T$ d'axiome E,
 $T \rightarrow T * F \mid F$ Terminaux {+, *, id, (,)}

Par la formule id + id, on a l'arbre de dérivation construit sur la Figure 3. a

Le problème est la chose de la règle de dérivation à appliquer

Définition 6: Une grammaire est LL(k) s'il est possible de faire un analyseur syntaxique descendant qui lit k symboles d'entrée à l'avance.

Nous allons étudier des analyseurs LL(1). Quelques outils sont nécessaires à la prédiction de la règle à choisir. On se place dans le cadre d'une grammaire $G = (A, V, P, S)$

Définition 7 - Pour $\alpha \in (A \cup V)^*$, $Primer(\alpha) = \{a \in A \mid \exists u \in A^*, \alpha \xrightarrow{*} au\}$
 - Pour $x \in V$, $suivant(x) = \{a \in A \cup \{\$ \} \mid \exists \alpha, \beta \in (A \cup V)^* \ S \xrightarrow{*} \alpha x a \beta\}$
 ou $\$$ est un marqueur de fin ($\beta = \epsilon$).

Remarque: On peut calculer primer récursivement, et suivant à l'aide de primer, en appliquant les règles suivantes jusqu'à un point fixe.

- 1) si x est un terminal, $primer(x) = \{x\}$.
- 2) si $x \rightarrow \epsilon$, $\epsilon \in primer(x)$
- 3) si $x \rightarrow y_1 \dots y_n$, et si $\epsilon \in primer(y_1) \cap \dots \cap primer(y_n)$, alors $primer(y_1) \cup \dots \cup primer(y_n)$
- 1) $\$ \in suivant(S)$
- 2) si $A \rightarrow \alpha B \beta$, $primer(B) \setminus \{\epsilon\} \subset suivant(B)$
- 3) si $A \rightarrow \alpha B$ ou $A \rightarrow \alpha B \beta$ et $\epsilon \in primer(B)$, alors $suivant(A) \subset suivant(B)$

D
V
P
1

théorème 8 une grammaire G est LL(1) si et seulement si, pour toute paire de productions distinctes $A \rightarrow \alpha$ / $A \rightarrow \beta$ de G , on a:

- 1) $\text{Premier}(\alpha) \cap \text{Premier}(\beta) = \emptyset$
- 2) Si $\epsilon \in \text{Premier}(\alpha)$, alors $\text{Premier}(\beta) \cap \text{Suivant}(A) = \emptyset$
- Si $\epsilon \in \text{Premier}(\beta)$, alors $\text{Premier}(\alpha) \cap \text{Suivant}(A) = \emptyset$

Dans le cas d'une grammaire LL(1), on peut donc construire un tableau: Table d'analyse prédictive qui contient les règles de production à appliquer à X en lisant a :

Table d'analyse prédictive (G)

<p>Pour chaque règle $X \rightarrow \alpha$</p> <p>Pour $a \in \text{Premier}(\alpha)$</p> <p>$M(a, X) := X \rightarrow \alpha$</p> <p>Si $\epsilon \in \text{Premier}(\alpha)$</p> <p>Pour chaque $b \in \text{Suivant}(X)$</p> <p>$M(b, X) := X \rightarrow \alpha$</p>

Pour une grammaire LL(1), on aura une seule règle par case et les cases sans règles sont des échecs.

Application 9 G n'est pas LL(1). on la remplace par

G' :

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

d'axe E .

D
V
P
2

On calcule premiers et suivants, et on crée la table d'analyse associée.

3) Analyse ascendante

On part des feuilles de l'arbre et on essaie de remonter jusqu'à la racine.

Exemple 10 Figure 3.5

L'analyse ascendante suit un processus de lecture / réduction; on lit des termes de l'entrée, on les réduit, on combine avec ce qu'on avait lu avant - La structure de donnée qui apparaît naturellement est une pile.

Définition 11 Une grammaire est LR(k) s'il est possible de faire un analyseur syntaxique ascendant qui lit k symboles d'entrée à l'avance.

Exemple d'analyse LR(0):

on construit l'automate à pile de l'analyse ascendante. Il est déterministe si et seulement si G est LR(0)

D
V
P
3

4) Analyse syntaxique en pratique

Comme pour l'analyse lexicale, on a vu des méthodes assez systématiques d'analyse syntaxique. Des méthodes plus générales existent, pour d'autres classes de grammaires (même ambiguës) et il y a également des constructeurs d'analyseurs syntaxiques:

Yacc par exemple. Il est utilisable sur le même principe que Lex.

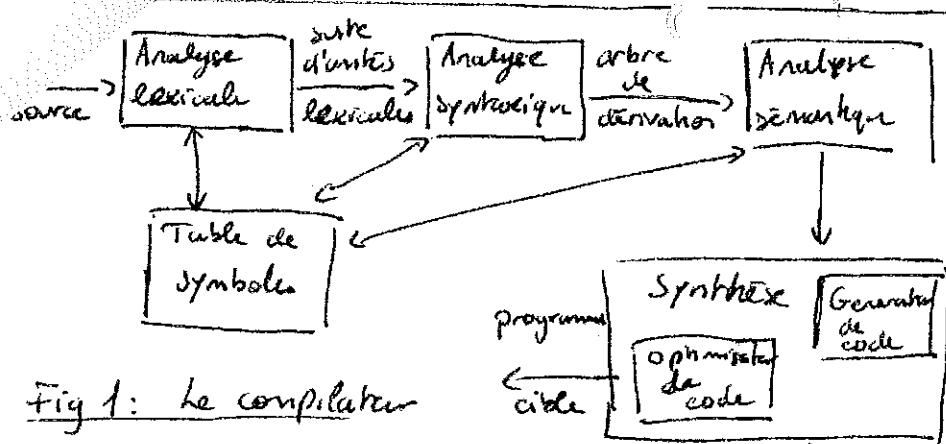


Fig 1: Le compilateur

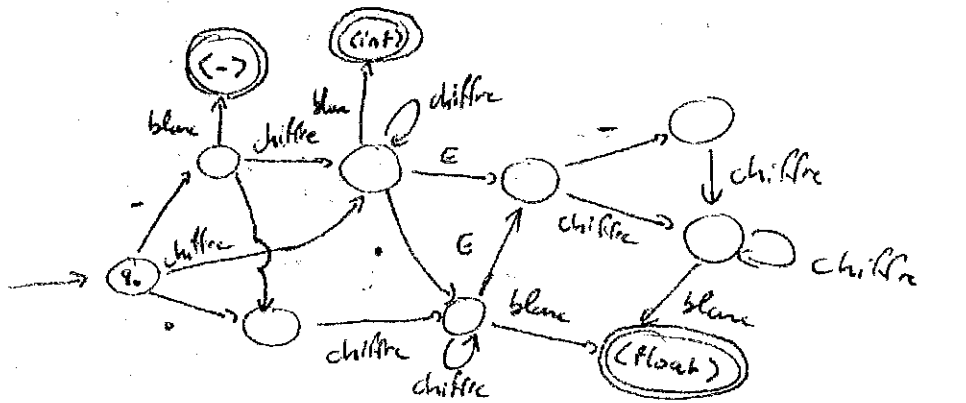
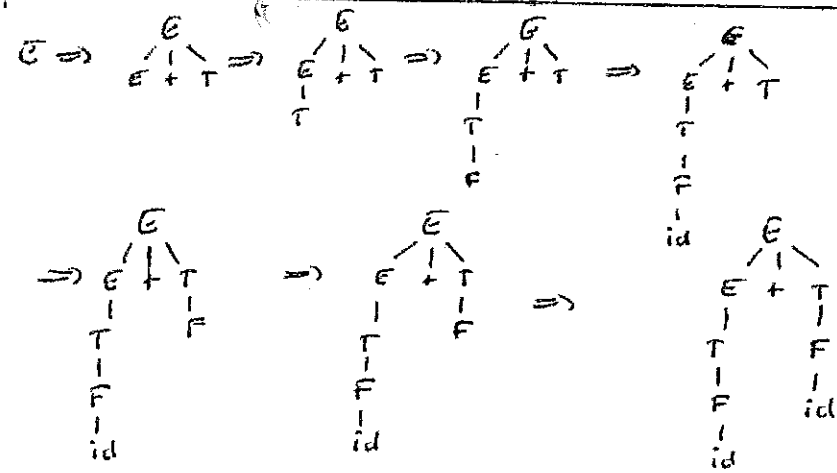
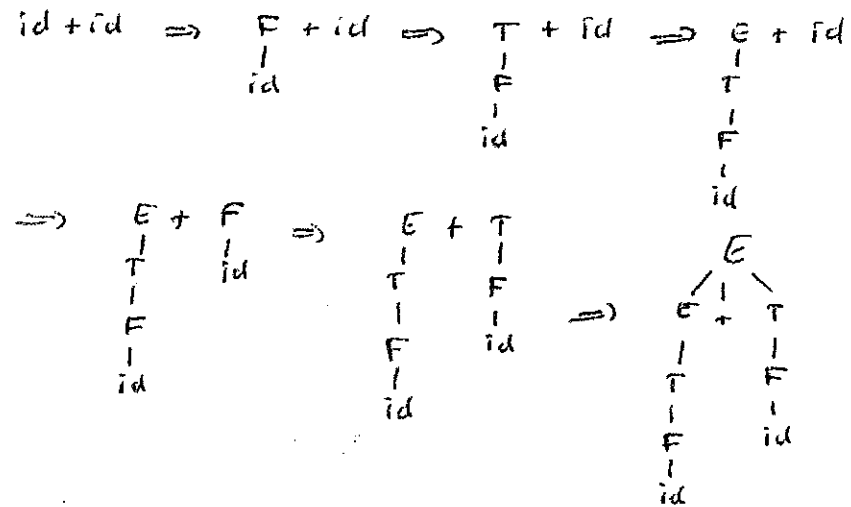


Fig 2: Automate pour la reconnaissance des entiers et des flottants en langage de calculatrice (et -)



a - analyse descendante



b - analyse ascendante

Figure 3 - analyse lexicale de id+id, par la grammaire

G

Algorithme de Cook, Younger et Kasami:

Référence: Carton

Leçon: 306, 307, 310, 323

Plan : I) Mix d'une grammaire sous forme normale de Chomsky

II) Réalisation d'un algorithme en utilisant la programmation dynamique.

III) Complexités ...

Définition: Forme normale de Chomsky

Une grammaire est sous forme normale de Chomsky ssi toutes ses règles sont de la forme $T \rightarrow UV$ ou $T \rightarrow a$.

Théorème: existence de FN de Chomsky

Quelle que soit $G = (A, T, P, S_0)$ il existe G' sous forme de Chomsky telle que $L_{G'} = L_G \setminus \{\epsilon\}$.

démonstration (en deux temps)

• Pour $a \in A$, on ajoute $V_a \in T$. On définit σ une substitution telle que

$$\sigma|_V = \text{Id}, \quad \sigma(a) = V_a \text{ pour tout } a \in A$$

On pose alors

$$P' = \{V_a \rightarrow a, a \in A\} \cup \{S \rightarrow \sigma(w), S \rightarrow w \in P, w \neq \epsilon\}$$

$$T' = T \cup \{V_a, a \in A\}.$$

$G' = (A, T', P', S_0)$ a toutes ses règles de la forme

$$S \rightarrow S_1 \dots S_n \text{ ou } S \rightarrow a \text{ et } L_{G'} = L_G \setminus \{\epsilon\}$$

• Il reste à éliminer sans modifier le langage les règles

$$S \rightarrow S_1 \dots S_n \quad (n \neq 2)$$

$$(\text{c}) \quad n = 1 : S \rightarrow S_1.$$

Pour chaque règle de la forme $S \rightarrow S_1$

pour chaque règle avec S dans le membre droit

On ajoute une nouvelle règle au l'én surface

Seul membre droit par S_1

On supprime les règles $S \rightarrow S_1$.

$n > 2$ On note $P_1 \dots P_m$ les productions de la forme $S \rightarrow S_1 \dots S_n$, $n > 2$.

On remplace chacune de ces productions P_k ($1 \leq k \leq m$)

$$S \rightarrow S_1 S_2^k, \quad S_i^k \rightarrow S_i S_{i+1}^k, \quad 1 \leq i < n-1, \quad S_{n-1}^k \rightarrow S_{n-1} S_n$$

Alors G'' est sous forme normale de Chomsky et est équivalente à G' donc à G . III

II) On a alors une grammaire $G = (A, T, P, S_0)$ sous forme normale de Chomsky.

On veut résoudre le problème du mot pour G .

$$w = w_1 \dots w_n \in A^*, \quad w[i,j] = w_1 \dots w_j$$

On note $E_{ij} = \{S \in T, w[i,j] \in L_G(S)\}$

(but) $\bullet w \in L_G \iff S_0 \in E_{1,n}$

(initialisation) \bullet Si $i=j$ $w[i,j]=w_i$ donc $E_{ij} = \{S, S \rightarrow w_i \in P\}$

(hérédité) \bullet Si $i < j$, $w[i,j] \in L_G(S)$ si $\exists k \in [i, j-1]$,

$$S_1 \in E_{i,k}, \quad S_2 \in E_{k+1,j} \quad \text{et } S \rightarrow S_1 S_2 \in P.$$

L'algorithme de programmation dynamique s'en déduit.

CYK (w, G) (G est sous forme de Chomsky)

$\forall (i,j) \in [1,n]^2$,

initialiser $E_{ij} \leftarrow \emptyset$

$\forall i \in [1,n], \forall S \rightarrow a \in P$

Si $w_i = a$

$E_{ii} = \{S\} \cup E_{ii}$

Pour d de 1 à $n-1$

Pour i de 1 à $n-d$

Si $j = i+d$

Appliquer la récurrence pour trouver E_{ij}

Retourner $S_0 \in E_{1,n}$

III) Complexités

- Pour la mise sous forme de Chomsky, on peut grandement augmenter le nombre de productions, mais si on prend en compte la taille des productions, c'est en fait seulement linéaire. $\Rightarrow O(n)$
- CYK comporte 3 boucles imbriquées de taille n (une boucle cachée dans "appliquer la récurrence" et une boucle de taille $O(|G|)$ cachée également. $\Rightarrow O(n^3 |G|)$

Construction d'un analyseur syntaxique

Référence : Dragon (Aho, Lam, Sethi, Ullman "Compilateurs")

Leçon 907, 910, 923

Pour rappel:

$$G = (\{+, *, (,), id\}, \{E, T, F\}, P, \epsilon)$$

$$P: \begin{aligned} E &\rightarrow E+T \mid T \\ T &\rightarrow T*F \mid F \\ F &\rightarrow (E) \mid id \end{aligned}$$

Cette grammaire n'est pas LL(1) (transitions à gauche)

On la remplace par

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \epsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \epsilon \\ F &\rightarrow (E) \mid id \end{aligned}$$

On va calculer les ensembles Prémix et suivant (cf Def 7)

	+	*	()	id	:	F	T	E'	E
Prémix	+	*	()	id	(, id				
Suivants), +, *,), + \$)	+, ε	(,), \$), \$), φ

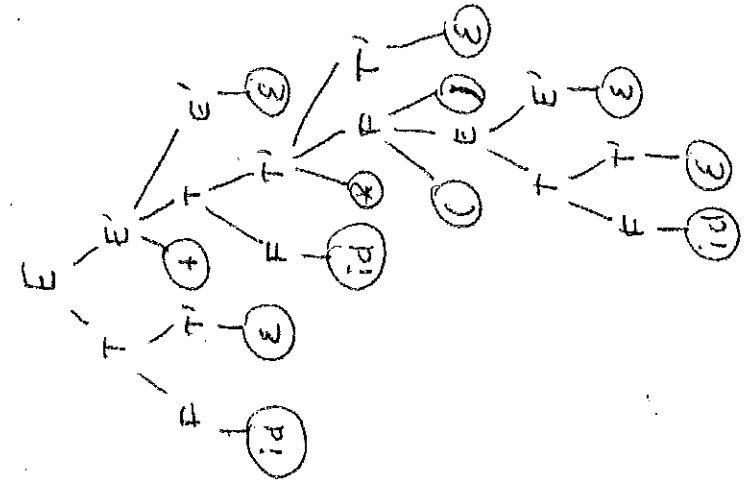
Table d'analyse:

+	*	()	id	\$
E		E → TE'		F → id	
E'	E' → TE'		E' → ε		
T		T → FT'			
T'	T' → ε	T' → FT'	T' → ε		
F		F → (E)			

Analyse d'une expression: $id + id * (id)$

Précédent	variable	regle
id	E	$E \rightarrow TE'$
	T, E'	$T \rightarrow FT'$
	F, T', E'	$F \rightarrow id$
+	T', E'	$T' \rightarrow \epsilon$
	E'	$E' \rightarrow +TE'$
id	T, E'	$T \rightarrow FT'$
	FT', E'	$F \rightarrow id$
*	T', E'	$T' \rightarrow *FT'$
(FT', E'	$F \rightarrow (E)$
id	E, T', E'	$E \rightarrow TE'$
	TE', T', E'	$T \rightarrow FT'$
)	FT', E', T', E'	$F \rightarrow id$
	T', E', T', E'	$T' \rightarrow \epsilon$
	E', T', E'	$E' \rightarrow \epsilon$
	T', E'	$T' \rightarrow \epsilon$
	E'	$E' \rightarrow \epsilon$

Construction de l'arbre: on applique les règles dans l'ordre sur la feuille la plus à gauche



- Un langage opératoire LL(1) respectivement utilisée? LI SP.

- $LL(1) \Rightarrow ?$ LI SP.