

## Traduction et Sémantique Analyse syntaxique

Luigi Santocanale  
LIF, Université de Provence  
Marseille, FRANCE

12 février 2010

Rappels de la théorie des langages  
Grammaires  
Automates à pile – AP  
Grammaires, APs, AFNs

## Plan

### Une grammaire non contextuelle ...

... est un tuple  $\mathcal{G} = \langle V, \Sigma, S, P \rangle$

où :

- $V$  est un alphabet fini,
- $\Sigma \subseteq V$  est l'alphabet des symboles **terminaux**,
- $S \in V \setminus \Sigma$  est le **symbole de départ**,
- $P \subseteq (V \setminus \Sigma) \times V^*$  est un ensemble fini de **productions**.

On appelle un symbole  $X \in V \setminus \Sigma$  un **non-terminal**.

## Un exemple

- $V = \{S, T, a, b\}$ ,
- $\Sigma = \{a, b\}$ ,
- et  $P = \{(S, T), (T, TT), (T, aTb), (T, \varepsilon)\}$ ,  
ce qu'on écrit :

$$\begin{aligned} S &\rightarrow T \\ T &\rightarrow TT \\ T &\rightarrow aTb \\ T &\rightarrow \varepsilon \end{aligned}$$

ou, alternativement :

$$\begin{aligned} S &\rightarrow T \\ T &\rightarrow TT \mid aTb \mid \varepsilon. \end{aligned}$$

## Langage d'une grammaire

Défini par :

$$L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}.$$

Attention :

si  $w \in L(\mathcal{G})$ , alors chaque lettre de  $w$  est un symbole terminal.

Une **dérivation** de  $w \in \Sigma^*$  de  $S$

est une suite  $w_0, \dots, w_n$  telle que :

- $S = w_0$ ,  $w_n = V$ , et
- $w_0 \Rightarrow w_1 \Rightarrow \dots \Rightarrow w_{n-1} \Rightarrow w_n$ .

## « Dériver »

Soient  $w, v \in V^*$ .

- On écrit

$$w \Rightarrow v$$

si

- ▶  $w = w_0 X w_1$ ,
- ▶  $v = w_0 u w_1$ , et
- ▶ la grammaire contient la production

$$X \rightarrow u.$$

- On écrit

$$w \Rightarrow^* v \quad (\text{\`a lire : } w \text{ se d\`erive \`a } v)$$

s'il existe  $n \geq 0$  et  $w_0, w_1, \dots, w_n$  tels que :

- ▶  $w = w_0$ ,  $w_n = V$ , et
- ▶  $w_0 \Rightarrow w_1 \Rightarrow \dots \Rightarrow w_{n-1} \Rightarrow w_n$ .

## Un exemple

Le mot  $abab \in L(\mathcal{G})$ , où  $\mathcal{G}$  est la grammaire

$$\begin{aligned} S &\rightarrow T \\ T &\rightarrow TT \\ T &\rightarrow aTb \\ T &\rightarrow \varepsilon \end{aligned}$$

car :

$$S \Rightarrow T \Rightarrow TT \Rightarrow aTbT \Rightarrow aTbaTb \Rightarrow aTbab \Rightarrow abab.$$

En fait :

$$\begin{aligned} L(G) &= \{w \in \Sigma^* \mid w \text{ est bien parenthés\`e}\} \\ &= \{w \in \Sigma^* \mid |w|_a = |w|_b \\ &\quad \text{et } |u|_b \leq |u|_a \text{ si } u \text{ est un pr\`efixe de } w\}. \end{aligned}$$

## Grammaires sans $\epsilon$ -règles

Une  $\epsilon$ -règle est une production de la forme

$$T \rightarrow \epsilon.$$

**Lemme.** Toute grammaire  $\mathcal{G}$  peut se transformer en une grammaire  $\mathcal{G}'$  telle que toute  $\epsilon$ -règle a la forme :

$$S \rightarrow \epsilon$$

(où  $S$  est le symbole de départ).

Notre exemple :

$$S \rightarrow \epsilon \mid T$$

$$T \rightarrow TT \mid aTb \mid ab.$$

Nous considérerons des grammaires ayant cette forme.

9/23

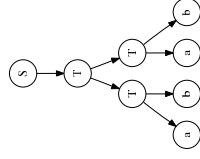
## Arbre de dérivation

On utilise les productions de  $\mathcal{G}$  pour construire un arbre.

Soit  $\mathcal{G}$  :

$$S \rightarrow \epsilon \mid T$$

$$T \rightarrow TT \mid aTb \mid ab.$$



$$S \Rightarrow T$$

$$\Rightarrow TT$$

$$\Rightarrow Tab$$

$$\Rightarrow abab$$

Le mot qui apparaît sur les feuilles est alors un mot de  $L(\mathcal{G})$ .

11/23

## Dérivation gauche et droite

- On appelle une **dérivation gauche** si on réécrit toujours le non-terminal plus à gauche.
- On appelle une **dérivation droite** si on réécrit toujours le non-terminal plus à droite.
- Un mot  $w \in L(\mathcal{G})$  peut avoir **plusieurs dérivations** qui témoignent de son appartenance à  $L(\mathcal{G})$ .

10/23

## Dérivations et leurs arbres

- Il existe une bijection entre dérivations gauches d'un mot  $w$  et arbres de dérivation de  $w$ .
- Il existe une bijection entre dérivations droites d'un mot et ses arbres de dérivation.
- Un mot peut avoir plusieurs arbres de dérivation, « essentiellement différentes ».
- Un mot peut avoir plusieurs arbres de dérivation.

12/23

## Ambiguïté

- Un **mot**  $w \in L(\mathcal{G})$  est **ambigu** s'il possède plus qu'un arbre de dérivation.
- Une **grammaire**  $\mathcal{G}$  est **ambiguë** s'il existe un mot  $w \in L(\mathcal{G})$  qui est ambigu.
- Un **langage**  $L$  est **ambigu** si toute grammaire  $\mathcal{G}$  telle que  $L = L(\mathcal{G})$  est ambiguë.

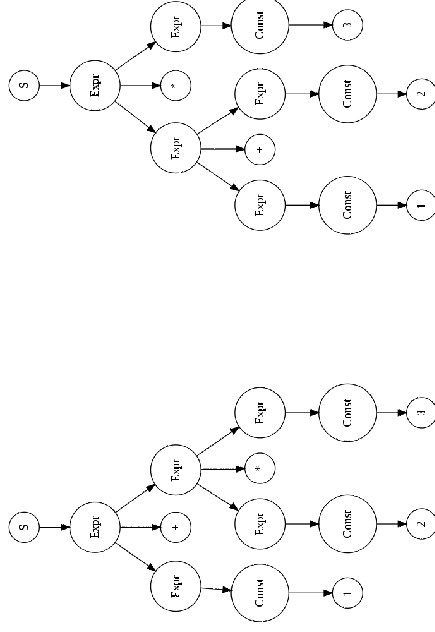
13/23

## Un mot ambigu

Le mot

$$1 + 2 * 3 \in L(\mathcal{G})$$

est ambigu, car :



15/23

## Un grammaire ambiguë

Soit  $\mathcal{G}$

$$\begin{aligned} S &\rightarrow \text{Expr} \\ \text{Expr} &\rightarrow \text{Expr} + \text{Expr} \mid \text{Expr} * \text{Expr} \mid (\text{Expr}) \mid \text{Const} \\ \text{Const} &\rightarrow \text{Chiffre} \mid \text{Chiffre Const} \\ \text{Chiffre} &\rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

Cette grammaire est ambiguë, car le mot

$$1 + 2 * 3 \in L(\mathcal{G})$$

est ambigu.

14/23

## Des langages (non) ambigus

Le langage  $L(\mathcal{G})$  est non ambigu. Soit  $\mathcal{G}'$  :

$$\begin{aligned} S &\rightarrow \text{Expr} \\ \text{Expr} &\rightarrow \text{Expr} + \text{Terme} \mid \text{Terme} \\ \text{Terme} &\rightarrow \text{Terme} * \text{Facteur} \mid \text{Facteur} \\ \text{Facteur} &\rightarrow (\text{Expr}) \mid \text{Const} \\ \text{Const} &\rightarrow \text{Chiffre} \mid \text{Chiffre Const} \\ \text{Chiffre} &\rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

alors  $L(\mathcal{G}') = L(\mathcal{G})$ .

Le langage

$$\{ a^n b^m c^m d^m \mid n, m \geq 1 \} \cup \{ a^n b^m c^m d^m \mid n, m \geq 1 \}$$

est ambigu.

16/23

## Automate à pile

... est un tuple

$$\langle Q, \Sigma, i, F, \Gamma, \Delta \rangle$$

où

- $Q$  ensemble fini d'états,
- $\Sigma$  alphabet d'entrée,
- $i \in Q$  état initial,  $F \subseteq Q$ ,
- $\Gamma$  alphabet de pile,
- $\Delta$ , la fonction de transition a le type suivant :
 
$$\Delta : Q \times (\Sigma \cup \{\varepsilon\}) \times \Gamma \rightarrow \mathcal{P}(\Gamma^* \times Q).$$

c'est-à-dire

$$\Delta(q, a, X) \subseteq \Gamma^* \times Q$$

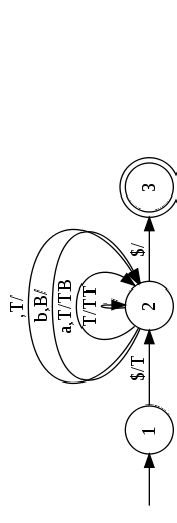
17/23

## Reconnaissance par AP

- Configuration :  $(q, w, \gamma)$  où
  - ▶  $q \in Q$ ,
  - ▶  $w \in \Sigma^*$  (la chaîne à lire),
  - ▶ et  $\gamma \in \Gamma^*$  (la pile).
- Transition :
  - $(q, w, \gamma) \vdash (q', w', \gamma')$  ssi
    - $w = aw', \gamma = T\alpha,$
    - $(\beta, q') \in \Delta(q, a, T)$  et  $\gamma' = \beta\alpha$
  - ou
    - $w = u, \gamma = T\alpha,$
    - $(\beta, q') \in \Delta(q, \varepsilon, T)$  et  $\gamma' = \beta\alpha.$

19/23

## Représentation graphique



- $Q = \{1, 2, 3\}, i = 1, F = \{3\},$
  - $\Sigma = \{a, b\},$
  - $\Gamma = \{\$, T, B\},$
  - $\Delta$
- | $\Delta$ | $(\varepsilon, \$)$             | $(a, T)$      | $(b, T)$               | ...         |
|----------|---------------------------------|---------------|------------------------|-------------|
| 1        | $\emptyset$                     | $\emptyset$   | $\emptyset$            | $\emptyset$ |
| 2        | $\{(TT, 2), (\varepsilon, 2)\}$ | $\{(TB, 2)\}$ | $\{(\varepsilon, 3)\}$ |             |
| 3        | $\emptyset$                     | $\emptyset$   | $\emptyset$            |             |

18/23

## Langage d'un AP

Le mot  $w \in L(\mathcal{A})$  ssi  
il existe  $n \geq 0$  et  $(q_0, w_0, \gamma_0), \dots, (q_n, w_n, \gamma_n)$  tels que :

- $(i, w, \$) = (q_0, w_0, \gamma_0),$
- $(q_n, w_n, \gamma_n) = (q_n, \varepsilon, \gamma_n)$  et  $q_n \in F,$
- $(q_i, w_i, \gamma_i) \vdash (q_{i+1}, w_{i+1}, \gamma_{i+1})$  pour  $i = 0, \dots, n-1.$

20/23

Exemple :  $abaabb \in L(\mathcal{A})$  car

$(1, abaabb, \$) \vdash (2, abaabb, T\$)$   
 $\vdash (2, a**ba**abb, T T\$)$   
 $\vdash (2, **ba**aabb, T B T\$)$   
 $\vdash (2, **ba**aabb, **B** T\$)$   
 $\vdash (2, **a**abb, T \$)$   
 $\vdash (2, **a**bb, T B\$)$   
 $\vdash (2, **bb**, T B B\$)$   
 $\vdash (2, **bb**, **B** B\$)$   
 $\vdash (2, **b**, **B** \$)$   
 $\vdash (2, , \$)$   
 $\vdash (3, , )$

21/23

## Le problème $w \in L(\mathcal{G})$

- Algorithme en temps  $O(n^3)$  :

 T. Kasami.

An efficient recognition algorithm for context-free languages.

Technical Report AFCRL-65-758, Air Force Cambridge research Lab., Bedford, MA, 1965.

 Daniel H. Younger.

Recognition and parsing of context-free languages in time  $n^3$ .

*Information and Control*, 10(2) :189–208, 1967.

- Pas satisfaisant au fin de l'analyse syntaxique.

Prochain objectifs :

- étudier des algorithmes – qui marchent en temps  $O(n)$  – adaptés à des sous-classes de grammaires.

23/23

**Théorème.**

Pour tout grammaire  $\mathcal{G}$  il existe un AP  $\mathcal{A}$  tel que  $L(\mathcal{A}) = L(\mathcal{G})$ .

Vice-versa :

Pour tout AP  $\mathcal{A}$  il existe une grammaire  $\mathcal{G}$  tel que  $L(\mathcal{G}) = L(\mathcal{A})$ .

**Théorème.**

Pour tout AFN  $\mathcal{A}$  il existe une grammaire  $\mathcal{G}$  tel que  $L(\mathcal{G}) = L(\mathcal{A})$ .

Considérez le langage

$$L = \{ a^n b^n \mid n \geq 0 \}.$$

**Théorème.**

Il existe un AP  $\mathcal{A}$  telle que  $L = L(\mathcal{A})$ .

Il n'existe aucun AFN  $\mathcal{A}$  tel que  $L = L(\mathcal{A})$ .

22/23