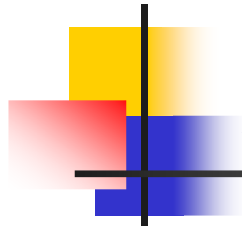




# PRINCIPES DE L'ANALYSE SYNTAXIQUE

---

*ESI – École nationale Supérieure en  
Informatique*



# Analyse syntaxique en bref I

---

- L'objectif d'une analyse syntaxique est de reconnaître si un programme donné (le programme source dont les entités lexicales ont été codées) appartient au langage engendré par une grammaire dite hors-contexte (de type 2 dans la classification de Chomsky).



# LES GRAMMAIRES NON-CONTEXTUELLES

---

- Une grammaire non-contextuelle  $G$  est défini par un quadruplet  $\langle \mathbf{N}, \mathbf{T}, \mathbf{P}, \mathbf{S}, \rangle$  où
- $\mathbf{N}$  : ensemble fini non vide de symboles appelés symboles non terminaux,
- $\mathbf{T}$  : ensemble fini de symboles appelés symboles terminaux,
- $\mathbf{S}$  :  $S \in \mathbf{N}$ , symbole de départ ou axiome de la grammaire,
- $\mathbf{P}$  : règles de réécritures ou ensembles de production,
  - Chaque production est de la forme :
  - $\mathbf{A} \rightarrow \alpha$
  - où  $A \in \mathbf{N}$  et  $\alpha \in (\mathbf{N} \cup \mathbf{T})^*$



# LES GRAMMAIRES NON-CONTEXTUELLES

- La grammaire  $G$  suivante décrit les expressions arithmétiques :
  - $G = \langle \{E\}, \{+, -, *, /, (, ), id\}, P, E \rangle$
  - $P :$   $E \rightarrow E + E$
  - $E \rightarrow E - E$
  - $E \rightarrow E * E$
  - $E \rightarrow E / E$
  - $E \rightarrow ( E )$
  - $E \rightarrow id$
- L'ensemble des productions peut être réécrit de la manière suivante :
  - $P :$   $E \rightarrow E + E \mid E - E \mid E * E \mid E / E \mid ( E ) \mid id$



# LES GRAMMAIRES NON-CONTEXTUELLES

---

- La grammaire  $G'$  suivante décrit également les expressions arithmétiques :
- $G' = \langle \{E, T, F\}, \{+, -, *, /, (, ), id\}, P, E \rangle$
- $P : E \rightarrow E + T \mid T$
- $T \rightarrow T * F \mid F$
- $F \rightarrow ( E ) \mid id$



# LES GRAMMAIRES NON-CONTEXTUELLES

---

- Quelle Grammaire choisir ?
- Sur quels critères se baser ?



# Rappels

---

- ***Définition d'une dérivation***

- Soit  $G$  une grammaire. On dit qu'une chaîne  $\omega_2$  se dérive d'une chaîne  $\omega_1$  si  $\omega_2$  s'obtient de  $\omega_1$  par l'application d'une règle de production de  $G$ . Cette dérivation est notée :
  - $\omega_1 \Rightarrow \omega_2$



# Rappels

---

## ■ Exemple :

- $A \rightarrow \beta$  est une production de la grammaire,
- $\alpha, \gamma \in (N \cup T)^*$ ,
- $\alpha A \gamma \Rightarrow \alpha \beta \gamma$





# Rappels

---

- **Série de dérivations :**

- Soient  $\alpha_1, \alpha_2, \dots, \alpha_m$  des chaînes appartenant à  $(N \cup T)^*$ ,  $m \geq 1$  et :

- $\alpha_1 \Rightarrow \alpha_2$

- $\alpha_2 \Rightarrow \alpha_3$

- $\vdots$

- $\alpha_{m-1} \Rightarrow \alpha_m$

- La suite de dérivation notée

- $\alpha_1 \Rightarrow^* \alpha_m$



# Rappels

---

- **Dérivation la plus à gauche :**
- Dans une suite de dérivation, si à chaque étape de dérivation une production est appliquée au non-terminal le plus à gauche, la suite de dérivation est dite ***la plus à gauche*** ou en terme anglais ***leftmost derivation***.



# Rappels

---

- **Dérivation la plus à droite :**
- Dans une suite de dérivation, si à chaque étape de dérivation une production est appliquée au non-terminal le plus à droite, la suite de dérivation est dite ***la plus à droite*** ou en terme anglais ***rightmost derivation***.



# Rappels

---

- La chaîne aabbaa peut s'obtenir de l'axiome S en utilisant :

**i)** La dérivation la plus à gauche :

■  $S \Rightarrow aAS \Rightarrow aSbAS \Rightarrow aabAS \Rightarrow aabbaS \Rightarrow aabbaa$

**ii)** La dérivation la plus à droite :

■  $S \Rightarrow aAS \Rightarrow aAa \Rightarrow aSbAa \Rightarrow aSbbaa \Rightarrow aabbaa$



# Rappels

---

## ■ **Langages**

- Etant donné une grammaire  $G$ , on appelle langage engendré par  $G$  et on le note  $L(G)$  le langage :
- $L(G) = \{ \omega \mid \omega \in T^* \text{ et } S \Rightarrow^+ \omega \}$



# Rappels

---

- **Exemple :**

- Soit la grammaire  $G = \langle N, T, P, S, \rangle$  avec :
- $N = \{S\}$  ;  $T = \{a, b\}$  ;
- $P = \{ S \rightarrow aSb ; S \rightarrow ab \}$  ;
- **$L(G) = \{ a^n b^n \mid n \geq 1 \}$**



# Arbre syntaxique

---

- Un arbre syntaxique également appelé arbre de dérivation est une structure pour représenter une suite de dérivations.



# Arbre syntaxique

---

- Soit  $G = \langle N, T, P, S \rangle$  une grammaire context-free, un arbre de dérivation est associé à une chaîne de  $L(G)$  et défini comme suit :
  - Chaque sommet de l'arbre a une étiquette qui est un symbole de  $N \cup T \cup \{\varepsilon\}$ ,
  - L'étiquette de la racine est le symbole  $S$ ,
  - Les sommets non-feuilles de l'arbre ont des étiquettes appartenant à  $N$ ,
  - Les feuilles de l'arbre ont des étiquettes appartenant à  $T \cup \{\varepsilon\}$ ,
  - Si un sommet d'étiquette  $A \in N$  a des fils d'étiquettes  $X_1, X_2, \dots, X_N$  (de gauche à droite) alors  $A \rightarrow X_1 X_2 \dots X_N$  est une production de  $P$ .





# Arbre syntaxique

---

- **Exemple :**

- Considérer la grammaire suivante :

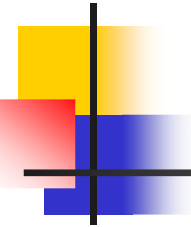
$$G = \langle \{S, A\}, \{a, b\}, P, S \rangle$$

$$P : S \rightarrow aAS \mid a$$

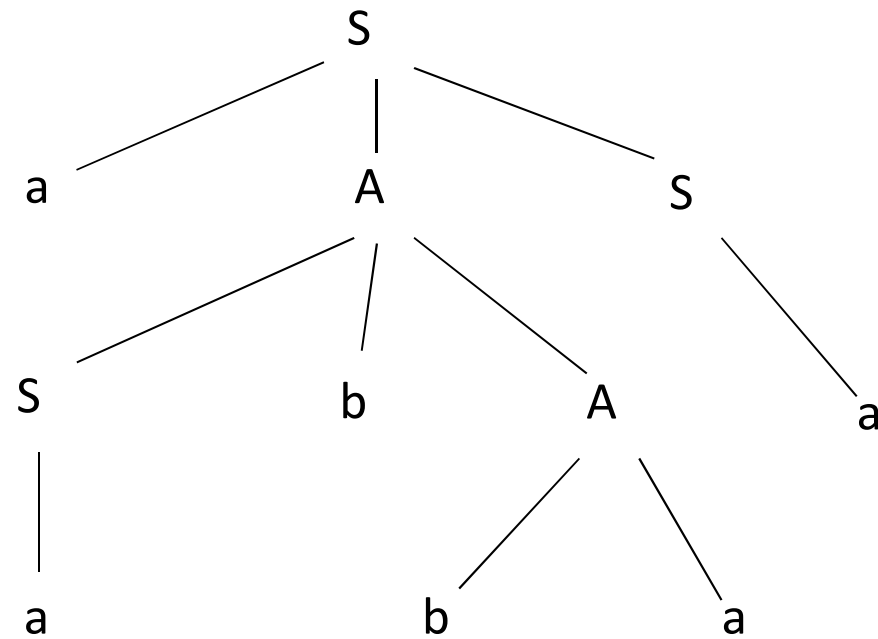
$$A \rightarrow SbA \mid SS \mid ba$$

- La suite de dérivations suivante pour obtenir la chaîne aabbaa :

- $S \Rightarrow aAS \Rightarrow aSbAS \Rightarrow aabAS \Rightarrow aabbaS \Rightarrow aabbaa$



# Arbre syntaxique





# Ambiguïté ou non ambiguïté d'une grammaire

---

- **Définition :**

- Une grammaire non-contextuelle est dite ambiguë si une chaîne  $\omega$  ( $\in T^*$ ) possède deux (ou plus) arbres syntaxiques.



# Ambiguïté ou non ambiguïté d'une grammaire

---

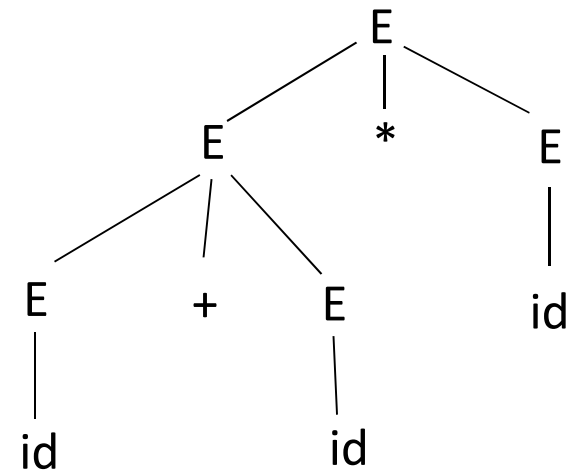
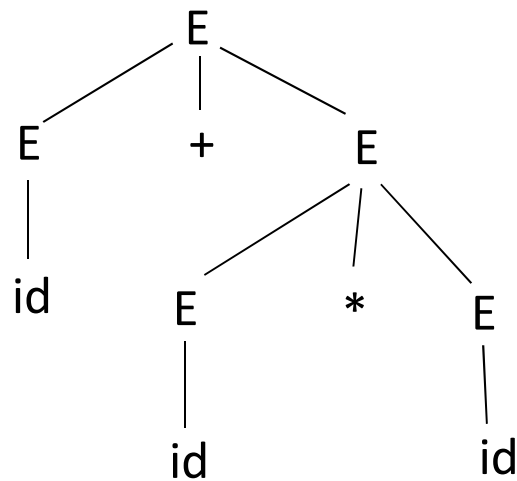
- **Exemple :**

- Soit  $G = \langle \{E\}, \{+, -, *, /, (, ), id\}, P, E \rangle$
- $P : E \rightarrow E + E \mid E - E \mid E * E \mid E / E \mid ( E ) \mid id$



# Ambiguïté ou non ambiguïté d'une grammaire

---





# Ambiguïté ou non ambiguïté d'une grammaire

---

## ■ Remarques :

- Obtenir deux suites de dérivations pour une chaîne  $\omega \in T^*$  n'implique pas forcément l'existence de deux arbres syntaxiques.
- Dans un cas général, la non-ambiguïté d'une grammaire est un problème indécidable. On dispose seulement de conditions suffisantes mais non nécessaires pour prouver qu'une grammaire est non ambiguë.
- Dans le cas des grammaires où nous disposons d'une connaissance sur le langage engendré (e.g. expressions arithmétiques) nous pouvons donner une grammaire non ambiguë qui génère le même langage.



# Ambiguïté ou non ambiguïté d'une grammaire

- La grammaire  $G' = \langle \{E, T, F\}, \{+, -, *, /, (, ), id\}, P', E \rangle$  dont les productions sont données ci-après, génère les expressions arithmétiques et elle est non ambiguë :
  - $E \rightarrow E + T \mid E - T \mid T$
  - $T \rightarrow T * F \mid T / F \mid F$
  - $F \rightarrow (E) \mid i$
- Il sera prouvé ultérieurement que cette grammaire est non ambiguë.



# Transformations de Grammaires non-contextuelles

---

- Pourquoi transformer ?
- Comment transformer ?
- Est-il possible de toujours transformer ?





# Transformations de Grammaires non-contextuelles

---

- Si un langage non-contextuel  $L$  est non vide, il peut être généré par une grammaire non-contextuelle ayant les propriétés suivantes :
  - Chaque non-terminal et chaque terminal apparaît dans la dérivation d'un mot de  $L$ .
  - Il n'y a pas de production unitaire (de la forme  $A \rightarrow B$  où  $(A,B) \in N^2$ ).
  - Si  $\varepsilon \notin L$  alors les  $\varepsilon$ -productions (de la forme  $A \rightarrow \varepsilon$  où  $A \in N$ ) peuvent être supprimées.



# Principe des analyses descendantes et ascendantes

---

- L'analyse doit être "déterministe" i.e. lorsque l'analyseur syntaxique s'arrête, soit le programme est syntaxiquement correct ou il existe une erreur dans ce programme.
- Il existe deux grandes classes de méthodes pour effectuer l'analyse syntaxique : les analyses descendantes et les analyses ascendantes.



# Principe des analyses descendantes et ascendantes

---

## ■ **Analyse descendante :**

- On part de l'axiome de la grammaire pour retrouver le programme source en effectuant des dérivations successives.
- Si on se place dans l'arbre syntaxique représentant le programme source, cette stratégie revient à partir de la racine (l'axiome) et à descendre vers les feuilles représentant les symboles terminaux.



# Principe des analyses descendantes et ascendantes

---

- **Analyse ascendante :**

- On part du programme source pour retrouver l'axiome de la grammaire en effectuant des dérivations successives inverses.
- On va cette fois partir des feuilles de l'arbre syntaxique pour remonter vers la racine.



# Principe des analyses descendantes et ascendantes

---

## ■ Exemple

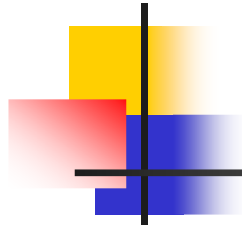
- Soit la grammaire  $G$  dont les productions sont données ci-après et la chaîne  $\omega=abba$  à analyser syntaxiquement :
  - $S \rightarrow aA \mid bB$
  - $A \rightarrow aBS \mid bS$
  - $B \rightarrow bB \mid a$



# Principe des analyses descendantes et ascendantes

---

- Le processus d'analyse descendante de la chaîne  $\omega = abba$  peut être illustré par :
  - $S \Rightarrow aA \Rightarrow abS \Rightarrow abbB \Rightarrow abba$
- Le processus d'analyse ascendante de la chaîne  $\omega = abba$  peut être illustré par :
  - $abba \Leftarrow^R abbB \Leftarrow^R abS \Leftarrow^R aA \Leftarrow^R S$
  - où  $\Leftarrow^R$  désigne une dérivation inverse



# Backus–Naur Form (BNF)

---

- La forme de Backus–Naur est une méta-syntaxe utilisée pour décrire les grammaires non contextuelles (context-free). Cette forme a été développée par John Backus et Peter Naur.
- BNF est très utilisée pour la description des grammaires des langages de programmation. Beaucoup de générateurs de compilateurs ou d'analyseurs syntaxiques s'inspirent de BNF pour l'introduction des grammaires cibles.



# Backus–Naur Form (BNF)

---

- **Définition :**
  - Une spécification BNF est un ensemble de règles de dérivation définie par :
    - **$\langle \text{symbole} \rangle ::= \text{expression}_1 \mid \text{expression}_2 \mid \dots \mid \text{expression}_N$**
- où
- $\langle \text{symbole} \rangle$  est un non-terminal, et  $\text{expression}_i$  est une suite de symboles de la grammaire (terminaux et non-terminaux et éventuellement le mot vide)
  - Les symboles qui n'apparaissent jamais en membre gauche de production sont des terminaux. Les symboles non terminaux sont toujours entre  $\langle \rangle$





# Backus–Naur Form (BNF)

---

- **Exemple :**
- $\langle \text{Expression} \rangle ::= \langle \text{Expression} \rangle "+" \langle \text{Terme} \rangle$   
|  $\langle \text{Expression} \rangle "-" \langle \text{Terme} \rangle$   
|  $\langle \text{Terme} \rangle$
- $\langle \text{Terme} \rangle ::= \langle \text{Terme} \rangle "*" \langle \text{Facteur} \rangle$   
|  $\langle \text{Terme} \rangle "/" \langle \text{Facteur} \rangle$   
|  $\langle \text{Facteur} \rangle$
- $\langle \text{Facteur} \rangle ::= "(" \langle \text{Expression} \rangle ")"$   
|  $"0" \mid "1" \mid "2" \mid "3" \mid "4" \mid "5" \mid "6" \mid "7" \mid "8" \mid "9"$