

# Langages, Compilation, Automates.

## Partie 5: Grammaires, types de grammaires

Florian Bridoux

Polytech Nice Sophia

2022-2023

## 1 Grammaires

## 2 Types de grammaires

## 1 Grammaires

## 2 Types de grammaires

## Définition (Grammaires de réécriture)

Une grammaire est un 4-uplet  $(N, \Sigma, R, S)$  où :

- $N$  est un ensemble fini de **symboles non terminaux**.
- $\Sigma$  est un ensemble fini de **symboles terminaux** avec  $N \cap \Sigma = \emptyset$ .
- $R$  est un sous ensemble fini de **règles**.

Une règle de  $R$  est de la forme  $\alpha \rightarrow \beta$  avec

- $\alpha, \beta \in (N \cup \Sigma)^*$  et ,
- $\alpha$  contient **au moins un non terminal**

$\alpha$  et  $\beta$  sont appelé respectivement **partie gauche et droite de la règle**.

- $S$  est un élément de  $N$  appelé l'**axiome** de la grammaire.

Pour alléger les notations, au lieu d'écrire plusieurs règles:

$$\alpha_1 \rightarrow \beta_1$$

$$\alpha_1 \rightarrow \beta_2$$

$$\alpha_1 \rightarrow \beta_3$$

on écrit simplement:

$$\alpha_1 \rightarrow \beta_1 \mid \beta_2 \mid \beta_3$$

# Proto-mots d'une grammaire

Les **proto-mots** d'une grammaire  $G = (N, \Sigma, R, S)$  sont des mots construits sur l'alphabet  $\Sigma \cup N$ , on les définit récursivement de la façon suivante :

- $S$  est un proto-mot de  $G$
- si  $\alpha\beta\gamma$  est un proto-mot de  $G$  et  $\beta \rightarrow \delta \in R$  alors  $\alpha\delta\gamma$  est un proto-mot de  $G$ .

Un proto-mot de  $G$  ne contenant aucun symbole non-terminal est appelé un mot **engendré** par  $G$ . Le **langage engendré par  $G$** , noté  $L(G)$  est l'ensemble des mots engendrés par  $G$ .

# Dérivation

- L'opération qui consiste à générer un proto-mot  $\alpha\delta\gamma$  à partir d'un proto-mot  $\alpha\beta\gamma$  et d'une règle de production  $r$  de la forme  $\beta \rightarrow \delta$  est appelée l'opération de **dérivation**. Elle se note à l'aide d'une double flèche :

$$\alpha\beta\gamma \Rightarrow \alpha\delta\gamma$$

- On note  $\alpha \xRightarrow{k} \beta$  pour indiquer que  $\beta$  se dérive de  $\alpha$  en  $k$  étapes.
- On définit aussi les deux notations  $\xRightarrow{+}$  et  $\xRightarrow{*}$  de la façon suivante :
  - $\alpha \xRightarrow{+} \beta \equiv \alpha \xRightarrow{k} \beta$  avec  $k > 0$
  - $\alpha \xRightarrow{*} \beta \equiv \alpha \xRightarrow{k} \beta$  avec  $k \geq 0$

## Attention

Les symboles  $\Rightarrow$  et  $\rightarrow$  ne représentent pas la même chose.

- Les symboles non terminaux appartenant à  $N$  sont représentés par des lettres latines majuscules :  $A, B, C, S, E, T \dots$
- Les symboles terminaux appartenant à  $\Sigma$  sont représentés par des lettres latines minuscules :  $a, b, c, d \dots$
- Les proto-mots appartenant à  $(N \cup \Sigma)^*$  sont représentés par des lettres grecques minuscules :  $\alpha, \beta, \gamma, \varepsilon \dots$
- L'axiome est très souvent représenté par le non-terminal  $S$  et constitue la partie gauche de la première règle de production

## Remarque:

En Lex et yacc, tout est inversé: les terminaux sont en majuscule et les non terminaux sont en minuscule.



- $L(G)$  est défini de la façon suivante :

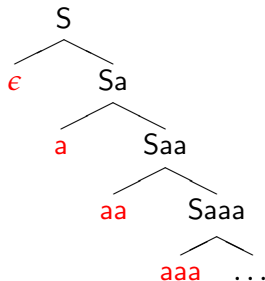
$$L(G) = \{m \in \Sigma^* \mid S \xRightarrow{+} m\}$$

- Deux grammaires  $G$  et  $G'$  sont équivalentes si  $L(G) = L(G')$ .

$$L_1 = \{\epsilon, a, aa, aaa, \dots\}$$

$$G = (\{S\}, \{a\}, \{S \rightarrow Sa \mid \epsilon\}, S)$$

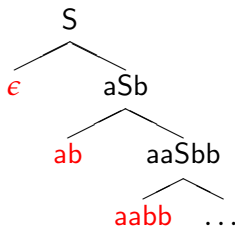
Proto-mots de  $G$ :



$$L_2 = \{\epsilon, ab, aabb, aaabbb, aaaabbbb, \dots\}$$

$$G = (\{S\}, \{a, b\}, \{S \rightarrow aSb \mid \epsilon\}, S)$$

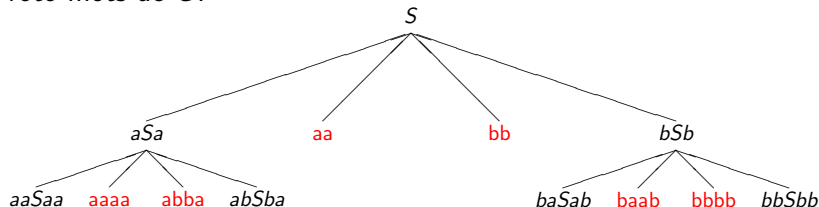
Proto-mots de  $G$ :



$$L_3 = \{aa, bb, aaaa, abba, baab, bbbb, \dots\}$$

$$G = (\{S\}, \{a, b\}, \{S \rightarrow aSa \mid bSb \mid aa \mid bb\}, S)$$

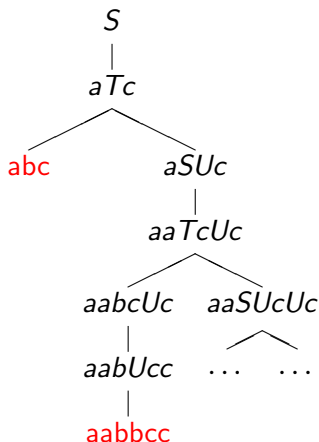
Proto-mots de  $G$ :



$$L_4 = \{\varepsilon, abc, aabbcc, aaabbbccc, \dots\}$$

$$G = (\{S, T, U\}, \{a, b, c\}, \left\{ \begin{array}{ll} S \rightarrow aTc, & T \rightarrow b \mid SU, \\ cU \rightarrow Uc, & bU \rightarrow bb \end{array} \right\}, S)$$

Proto-mots de  $G$ :



# Sens de dérivation

$G = (\{E, T, F\}, \{+, *, a\}, \{E \rightarrow T + E \mid T, T \rightarrow F * T \mid F, F \rightarrow a\}, E)$   
Les proto-mots engendrés lors d'une dérivation peuvent comporter plus d'un symbole non-terminal :

$$\begin{aligned} E &\Rightarrow T + E \\ &\Rightarrow T + T \\ &\Rightarrow F + T \\ &\Rightarrow F + F * T \\ &\Rightarrow F + a * T \\ &\Rightarrow F + a * F \\ &\Rightarrow a + a * F \\ &\Rightarrow a + a * a \end{aligned}$$

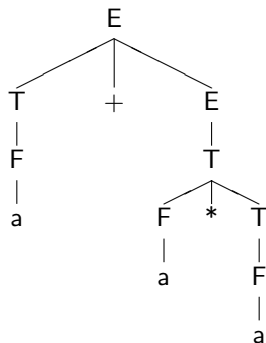
**Dérivation gauche** : on réécrit  
le non-terminal le plus à  
gauche :

$$\begin{aligned} E &\Rightarrow T + E \\ &\Rightarrow F + E \\ &\Rightarrow a + E \\ &\Rightarrow a + T \\ &\Rightarrow a + F * T \\ &\Rightarrow a + a * T \\ &\Rightarrow a + a * F \\ &\Rightarrow a + a * a \end{aligned}$$

**Dérivation droite** :  
...

$$\begin{aligned} E &\Rightarrow T + E \\ &\Rightarrow T + T \\ &\Rightarrow T + F * T \\ &\Rightarrow T + F * F \\ &\Rightarrow T + F * a \\ &\Rightarrow T + a * a \\ &\Rightarrow F + a * a \\ &\Rightarrow a + a * a \end{aligned}$$

# Arbre de dérivation



Un arbre de dérivation pour  $G = (N, \Sigma, R, S)$  est un arbre ordonné et étiqueté dont les étiquettes appartiennent à l'ensemble  $N \cup \Sigma \cup \{\varepsilon\}$ . Si un nœud de l'arbre est étiqueté par le non terminal  $A$  et ses fils sont étiquetés  $X_1, X_2, \dots, X_n$  alors la règle  $A \rightarrow X_1X_2\dots X_n$  appartient à  $P$ .



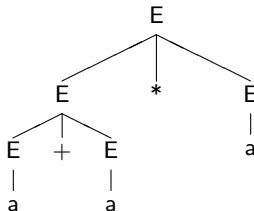
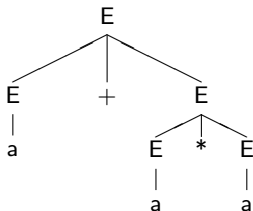
- Un arbre de dérivation indique les règles qui ont été utilisées dans une dérivation, mais pas l'ordre dans lequel elles ont été utilisées.

# Ambiguïté

Une grammaire  $G$  est **ambiguë** s'il existe au moins un mot  $m$  dans  $L(G)$  auquel correspond plus d'un arbre de dérivation.

*Exemple :*

$$E \rightarrow E + E \mid E * E \mid a$$

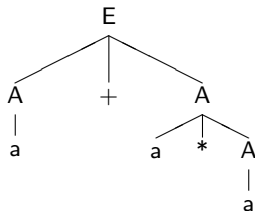


Dans le contexte de la compilation,

- Soit on essaye d'éviter les grammaires ambiguës,
- Soit on ajoute des règles pour résoudre les problèmes de conflit liés à l'ambiguïté de la grammaire.
- Pour qu'il y ait unicité de l'analyse

*Exemple :*

$$E \rightarrow A + E \mid A, A \rightarrow a * A \mid a$$



## 1 Grammaires

## 2 Types de grammaires

# Types de règles

Les grammaires peuvent être classées en fonction de la forme de leurs règles de production. On définit cinq types de règles :

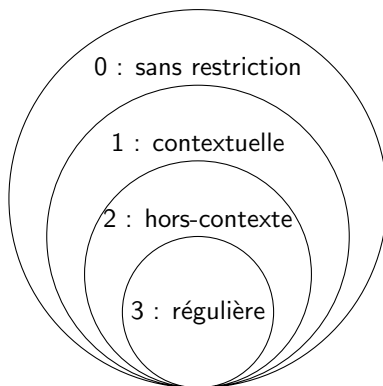
- Une règle est **régulière à gauche** si et seulement si elle est de la forme  $A \rightarrow xB$ ,  $A \rightarrow x$  ou  $A \rightarrow \epsilon$  avec  $A, B \in N$  et  $x \in \Sigma^*$ .
- Une règle est **régulière à droite** si et seulement si elle est de la forme  $A \rightarrow Bx$ ,  $A \rightarrow x$  ou  $A \rightarrow \epsilon$  avec  $A, B \in N$  et  $x \in \Sigma^*$ .
- Une règle  $A \rightarrow \alpha$  est un règle **hors-contexte** si et seulement si :  $A \in N$  et  $\alpha \in (N \cup \Sigma)^*$ .
- Une règle  $\alpha \rightarrow \beta$  est une règle **contextuelle** si et seulement si :  $\alpha = gAd$  et  $\beta = g\delta d$  avec  $g, d, \delta \in (N \cup \Sigma)^*$  et  $A \in N$ .  
*Le nom "contextuelle" provient du fait que  $A$  se réécrit  $B$  uniquement dans le contexte  $g_d$ .*
- Une règle  $\alpha \rightarrow \beta$  **sans restriction** sont les plus générale. Il faut juste que  $\alpha$  ait au moins un non terminal.

# Type d'une grammaire

Une grammaire est :

- **régulière** ou de type 3 si elle est régulière à droite ou régulière à gauche. Une grammaire est régulière à gauche si **toutes** ses règles sont régulières à gauche et une grammaire est régulière à droite si **toutes** ses règles sont régulières à droite.
- **hors contexte** ou de type 2 si toutes ses règles de production sont hors contexte.
- **dépendante du contexte** ou de type 1 si toutes ses règles de production sont dépendantes du contexte.
- **sans restrictions** ou de type 0 si toutes ses règles de production sont sans restrictions.

# Hiérarchie de Chomsky-Schützenberger





# Type d'un langage

Un langage pouvant être engendré par une grammaire de type  $x$  et pas par une grammaire d'un type supérieur dans la hiérarchie, est appelé un **langage de type  $x$** .

Type du langage	Nom
3	régulier
2	hors contexte
1	dépendant du contexte
0	rékursivement énumérable

Les grammaires hors contextes sont aussi appelé algébriques.

# Exemples de langages réguliers

$$L_1 = \{m \in \{a, b\}^*\}$$

$$L_2 = \{m \in \{a, b\}^* \mid |m|_a \bmod 2 = 0\}$$

$$L_3 = \{m \in \{a, b\}^* \mid m = xaaa \text{ avec } x \in \{a, b\}^*\}$$

$$L_4 = \{m \in \{a, b\}^* \mid |m|_a \bmod 2 = 0 \text{ et } |m|_b \bmod 2 = 0\}$$

# Exemples de langages réguliers

$$L_1 = \{m \in \{a, b\}^*\}$$

$$G_1 = (\{S\}, \{a, b\}, \{S \rightarrow aS \mid bS \mid \varepsilon\}, S)$$

$$L_2 = \{m \in \{a, b\}^* \mid |m|_a \bmod 2 = 0\}$$

$$L_3 = \{m \in \{a, b\}^* \mid m = xaaa \text{ avec } x \in \{a, b\}^*\}$$

$$L_4 = \{m \in \{a, b\}^* \mid |m|_a \bmod 2 = 0 \text{ et } |m|_b \bmod 2 = 0\}$$

# Exemples de langages réguliers

$$L_1 = \{m \in \{a, b\}^*\}$$

$$G_1 = (\{S\}, \{a, b\}, \{S \rightarrow aS \mid bS \mid \varepsilon\}, S)$$

$$L_2 = \{m \in \{a, b\}^* \mid |m|_a \bmod 2 = 0\}$$

$$G_2 = (\{S, T\}, \{a, b\}, \left\{ \begin{array}{l} S \rightarrow aT \mid bS \mid \varepsilon, \\ T \rightarrow aS \mid bT \end{array} \right\}, S)$$

$$L_3 = \{m \in \{a, b\}^* \mid m = xaaa \text{ avec } x \in \{a, b\}^*\}$$

$$L_4 = \{m \in \{a, b\}^* \mid |m|_a \bmod 2 = 0 \text{ et } |m|_b \bmod 2 = 0\}$$

# Exemples de langages réguliers

$$L_1 = \{m \in \{a, b\}^*\}$$

$$G_1 = (\{S\}, \{a, b\}, \{S \rightarrow aS \mid bS \mid \varepsilon\}, S)$$

$$L_2 = \{m \in \{a, b\}^* \mid |m|_a \bmod 2 = 0\}$$

$$G_2 = (\{S, T\}, \{a, b\}, \left\{ \begin{array}{l} S \rightarrow aT \mid bS \mid \varepsilon, \\ T \rightarrow aS \mid bT \end{array} \right\}, S)$$

$$L_3 = \{m \in \{a, b\}^* \mid m = xaaa \text{ avec } x \in \{a, b\}^*\}$$

$$G_3 = (\{S, T, U\}, \{a, b\}, \left\{ \begin{array}{l} S \rightarrow aS \mid bS \mid aT, \\ T \rightarrow aU, \\ U \rightarrow a \end{array} \right\}, S)$$

$$L_4 = \{m \in \{a, b\}^* \mid |m|_a \bmod 2 = 0 \text{ et } |m|_b \bmod 2 = 0\}$$

# Exemples de langages réguliers

$$L_1 = \{m \in \{a, b\}^*\}$$

$$G_1 = (\{S\}, \{a, b\}, \{S \rightarrow aS \mid bS \mid \varepsilon\}, S)$$

$$L_2 = \{m \in \{a, b\}^* \mid |m|_a \bmod 2 = 0\}$$

$$G_2 = (\{S, T\}, \{a, b\}, \left\{ \begin{array}{l} S \rightarrow aT \mid bS \mid \varepsilon, \\ T \rightarrow aS \mid bT \end{array} \right\}, S)$$

$$L_3 = \{m \in \{a, b\}^* \mid m = xaaa \text{ avec } x \in \{a, b\}^*\}$$

$$G_3 = (\{S, T, U\}, \{a, b\}, \left\{ \begin{array}{l} S \rightarrow aS \mid bS \mid aT, \\ T \rightarrow aU, \\ U \rightarrow a \end{array} \right\}, S)$$

$$L_4 = \{m \in \{a, b\}^* \mid |m|_a \bmod 2 = 0 \text{ et } |m|_b \bmod 2 = 0\}$$

$$G_4 = (\{S, T, U, V\}, \{a, b\}, \left\{ \begin{array}{l} S \rightarrow aT \mid bU \mid \varepsilon, \quad T \rightarrow aS \mid bV, \\ V \rightarrow aU \mid bT, \quad U \rightarrow aV \mid bS \end{array} \right\}, S)$$

# Exemples de langages hors-contexte

$$L_1 = \{a^n b^n \mid n \geq 0\}$$

$$L_2 = \{mm^{-1} \mid m \in \{a, b\}^*\} \quad (\text{langage miroir - palindromes paires})$$

# Exemples de langages hors-contexte

$$L_1 = \{a^n b^n \mid n \geq 0\}$$

$$G_1 = (\{S\}, \{a, b\}, \{S \rightarrow aSb \mid \varepsilon\}, S)$$

$$L_2 = \{mm^{-1} \mid m \in \{a, b\}^*\} \quad (\text{langage miroir - palindromes paires})$$



# Exemples de langages hors-contexte

$$L_1 = \{a^n b^n \mid n \geq 0\}$$

$$G_1 = (\{S\}, \{a, b\}, \{S \rightarrow aSb \mid \varepsilon\}, S)$$

$$L_2 = \{mm^{-1} \mid m \in \{a, b\}^*\} \quad (\text{langage miroir - palindromes paires})$$

$$G_2 = (\{S\}, \{a, b\}, \{S \rightarrow aSa \mid bSb \mid \varepsilon\}, S)$$

# Exemples de langages contextuels

$$L_1 = \{a^n b^n c^n \mid n \geq 0\}$$

# Exemples de langages contextuels

$$\begin{aligned} L_1 &= \{a^n b^n c^n \mid n \geq 0\} \\ G_1 &= (\{S, B, C\}, \{a, b, c\}, \\ &\quad \left\{ \begin{array}{ll} S \rightarrow aSBC \mid \varepsilon, & CB \rightarrow BC, \\ aB \rightarrow ab, & bB \rightarrow bb, \\ bC \rightarrow bc, & cC \rightarrow cc \end{array} \right\}, S) \end{aligned}$$

- Une **grammaire** d'un langage  $L$  permet de générer tous les mots appartenant à  $L$ .
- Un **reconnaisseur** pour un langage  $L$  est un programme qui prend en entrée un mot  $m$  et répond **oui** si  $m$  appartient à  $L$  et **non** sinon.
- Pour chaque classe de grammaire, il existe une classe de reconnaisseurs qui définit la même classe de langages.

Type de grammaire	Type de reconnaisseur
régulière	Automate fini
hors contexte	Automate à pile
contextuelle	Automate linéairement borné
sans restriction	Machine de Turing