

# Chapitre 4

Langages réguliers et hors-contexte

---

## Définition 4.1

Un **alphabet** est un ensemble fini et non vide de symboles.

## Exemple 4.2

$$\Sigma_1 = \{a, b, \dots, z\}$$

$$\Sigma_2 = \{a, b\}$$

$$\Sigma_3 = \{0, 1\}$$

$$\Sigma_4 = \{0, 1, \dots, 9\}$$

$$\Sigma_5 = \{\clubsuit, \diamondsuit, \heartsuit, \spadesuit\}$$

$$\Sigma_6 = \{0, 1, \dots, 9, +, -, *, /, (, )\}$$

$$\Sigma_7 = \text{ASCII}$$

$$\Sigma_8 = \text{ADN} = \{\langle \text{ADÉNINE} \rangle, \langle \text{CYTOSINE} \rangle, \langle \text{GUANINE} \rangle, \langle \text{THYMINE} \rangle\}$$

### Définition 4.3

Un **mot** est une suite finie, possiblement vide, de symboles appartenant à un alphabet.

### Exemple 4.4

*abba est un mot sur l'alphabet  $\{a, b\}$ .*

### Définition 4.5

La lettre grecque epsilon minuscule,  $\epsilon$ , dénote le **mot vide**, quel que soit l'alphabet en usage.

#### Définition 4.6

$\Sigma^*$  est l'ensemble de tous les mots pouvant être formés à partir de l'alphabet  $\Sigma$ .

#### Remarque 4.7

*On a toujours  $\varepsilon \in \Sigma^*$ , quel que soit  $\Sigma$ .*

## Exemple 4.8

$$\Sigma_2 = \{a, b\}$$

$$\Sigma_2^* = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, aab, \dots\}$$

$$\Sigma_3 = \{0, 1\}$$

$$\Sigma_3^* = \{\varepsilon, 0, 1, 00, 01, 10, 11, 000, 001, 010, \dots\}$$

$$\Sigma_6 = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, +, -, *, /, (, )\}$$

$$\Sigma_6^* \ni (25 + 50) * 56$$

$$\Sigma_6^* \ni (()25 + + + 6$$

### Définition 4.9

Soit  $w = w_1 w_2 \dots w_k \in \Sigma^*$ , alors

$$|w| = k.$$

L'entier  $k$  est la **longueur** du mot  $w$ , c'est-à-dire le nombre de symboles dans  $w$ .

### Exemple 4.10

$$|\text{allo}| = 4.$$

### Définition 4.11

Soient  $w = w_1 w_2 \dots w_k \in \Sigma^*$  et  $a \in \Sigma$ , alors

$$|w|_a = \#\{w_i \mid w_i = a\},$$

c'est-à-dire le nombre d'occurrences du symbole  $a$  dans le mot  $w$ .

### Exemple 4.12

$$|\text{yabadabadoo}|_a = 4 \quad \text{et} \quad |\text{yabadabadoo}|_b = 2.$$



### Définition 4.13

Soient

$$x = x_1x_2 \dots x_k \in \Sigma^* \quad \text{et} \quad y = y_1y_2 \dots y_l \in \Sigma^*,$$

alors

$$x \cdot y = x_1x_2 \dots x_ky_1y_2 \dots y_l$$

est la **concaténation** de  $x$  et  $y$ .

### Remarque 4.14

$w \cdot \varepsilon = w$ , quel que soit le mot  $w$ .

### Exemple 4.15

$$\text{bon} \cdot \text{jour} = \text{bonjour}$$

$$\text{bon} \cdot \varepsilon = \text{bon}$$

### Définition 4.16

Soient un mot  $w \in \Sigma^*$  et un entier positif  $i$ , alors

$$w^i = \underbrace{w \cdot w \cdots w}_{i \text{ fois}},$$

et

$$w^0 = \varepsilon.$$

De plus, on omet souvent l'opérateur de concaténation :

$$xy = x \cdot y.$$

### Définition 4.17

Soit un mot  $w \in \Sigma^*$ , alors  $w^R$  est le **renversé** de  $w$ , c'est-à-dire le mot formé avec les symboles de  $w$  pris de droite à gauche.

### Définition 4.18

Soit  $<_{\Sigma}$  une relation d'ordre sur l'alphabet  $\Sigma$ . On appelle **ordre lexicographique** la relation d'ordre suivante sur  $\Sigma^*$  :

$$\left. \begin{array}{l} |w| < |w'| \Rightarrow w <_{\Sigma^*} w' \\ |w| = |w'| \\ w = x \cdot a \cdot y \\ w' = x \cdot b \cdot z \\ a <_{\Sigma} b \end{array} \right\} \Rightarrow w <_{\Sigma^*} w'$$

### Exemple 4.19

$$\begin{aligned} \Sigma &= \{0, 1\} \\ \Sigma^* &= \{\varepsilon, 0, 1, 00, 01, 10, 11, 000, 001, 010, \dots\} \end{aligned}$$

### Définition 4.20

Un **langage** sur l'alphabet  $\Sigma$  est un sous-ensemble de  $\Sigma^*$ .

### Définition 4.21

On note  $\{\}$  ou  $\emptyset$  le **langage vide**, c'est-à-dire le langage qui ne contient aucun mot.

### Remarque 4.22

$\varepsilon \notin \emptyset$ .

### Exemple 4.23

Divers langages sur l'alphabet  $\Sigma = \{a, b\}$  :

$$L_1 = \{a, b, abba, abbb\}$$

$$\begin{aligned} L_2 &= \{w \mid |w| = 2k, k \in \mathbb{N}\} \\ &= \{\varepsilon, aa, ab, ba, bb, aaaa, aaab, aaba, aabb, \dots\} \\ &= \text{les mots qui contiennent un nombre pair de symboles} \end{aligned}$$

$$\begin{aligned} L_3 &= \{w \mid w = aw', w' \in \Sigma^*\} \\ &= \{a, aa, ab, aaa, aab, aba, abb, aaaa, \dots\} \\ &= \text{les mots qui commencent par a} \end{aligned}$$

$$\begin{aligned}
L_4 &= \{w \mid |w|_a = |w|_b\} \\
&= \{\varepsilon, ab, ba, aabb, abab, abba, baab, baba, \dots\} \\
&= \text{les mots qui contiennent autant de a que de b}
\end{aligned}$$

$$\begin{aligned}
L_5 &= \{w \mid |w|_a - |w|_b \text{ est un nombre premier}\} \\
&= \{aaab, \dots, aabbabaaaaa, \dots\}
\end{aligned}$$

## Exemple 4.24

*Divers langages sur l'alphabet  $\Sigma = \text{ASCII}$  :*

$$L_6 = \{w \mid w = (x, y, z) \text{ avec } x, y, z \in \{0, 1, \dots, 9\}^* \text{ et } xy = z\}$$

$$L_7 = \{w \mid w \text{ est une fonction syntaxiquement correcte} \\ \text{dans le langage de programmation Java} \\ \text{et cette fonction retourne toujours } \langle \text{TRUE} \rangle\}$$

$$L_8 = \{w \mid w \text{ compile sans erreur en C++}\}$$

$$L_9 = \{w \mid w \text{ est un programme C++ qui écrit "Hello world!"}\}$$



### Exemple 4.25

*Divers langages sur l'alphabet  $\Sigma = \text{ADN}$  :*

$$L_9 = \{w \mid w \text{ est la séquence d'un gène bactérien}\}$$

$$L_{10} = \{w \mid w \text{ est le génome d'un individu}\}$$

- l'union
- l'intersection
- le complément
- le produit cartésien
- l'ensemble des parties

### Définition 4.26 (Union de deux langages)

L'**union** de deux langages  $L_1$  et  $L_2$  :

$$L_1 \cup L_2 = \{w \mid w \in L_1 \text{ ou } w \in L_2\}.$$

### Exemple 4.27

*Soient*

$$\Sigma = \{a, b\}, L_1 = \{\varepsilon, a, aa\} \text{ et } L_2 = \{a, bb\},$$

*alors*

$$L_1 \cup L_2 = \{\varepsilon, a, aa, bb\}.$$

### Définition 4.28 (Intersection de deux langages)

L'**intersection** de deux langages  $L_1$  et  $L_2$  :

$$L_1 \cap L_2 = \{w \mid w \in L_1 \text{ et } w \in L_2\}.$$

### Exemple 4.29

*Soient*

$$\Sigma = \{a, b\} \text{ et } L_1 = \{\varepsilon, a, aa\} \text{ et } L_2 = \{a, bb\},$$

*alors*

$$L_1 \cap L_2 = \{a\},$$

$$L_1 \cap \emptyset = \emptyset.$$

### Définition 4.30 (Complément d'un langage)

Le **complément** d'un langage  $L$  sur l'alphabet  $\Sigma$  :

$$\overline{L} = \{w \in \Sigma^* \mid w \notin L\}.$$

### Exemple 4.31

*Soient*

$$\Sigma = \{a, b\} \text{ et } L = \{\varepsilon, a, aa\},$$

*alors*

$$\overline{L} = \{b, ab, ba, bb, aaa, aab, aba, baa, bba, \dots\}.$$

### Définition 4.32 (Produit cartésien de deux langages)

Le *produit cartésien* de deux ensembles  $\mathcal{A}$  et  $\mathcal{B}$  :

$$\mathcal{A} \times \mathcal{B} = \{(x, y) : x \in \mathcal{A} \text{ et } y \in \mathcal{B}\}.$$

### Remarque 4.33

On a que  $\#(\mathcal{A} \times \mathcal{B}) = (\#\mathcal{A})(\#\mathcal{B})$  et que  $\mathcal{A} \times \emptyset = \emptyset$ .

### Exemple 4.34

Soient

$$\Sigma = \{a, b\}, L_1 = \{\varepsilon, a, aa\}, \text{ et } L_2 = \{a, bb\},$$

alors

$$L_1 \times L_2 = \{(\varepsilon, a), (\varepsilon, bb), (a, a), (a, bb), (aa, a), (aa, bb)\}.$$

### Notation 4.35

Pour un ensemble  $\mathcal{A}$  et un entier  $k \geq 2$  :

$$\begin{aligned}\mathcal{A}^k &= \underbrace{\mathcal{A} \times \mathcal{A} \times \cdots \times \mathcal{A}}_{k \text{ fois}} \\ &= \text{l'ensemble des } k\text{-tuplets d'éléments de } \mathcal{A}.\end{aligned}$$

### Remarque 4.36

$$\#(\mathcal{A}^k) = (\#\mathcal{A})^k.$$

### Définition 4.37 (Ensemble des parties)

L'ensemble des parties d'un ensemble  $\mathcal{A}$  est défini comme :

$$\mathcal{P}(\mathcal{A}) = \{\mathcal{B} : \mathcal{B} \subseteq \mathcal{A}\}.$$

### Remarque 4.38

$$\#\mathcal{P}(\mathcal{A}) = 2^{\#\mathcal{A}}.$$

### Exemple 4.39

*Soient*

$$\Sigma = \{a, b\} \text{ et } L = \{\varepsilon, a, aa\},$$

*alors*

$$\mathcal{P}(L) = \{\emptyset, \{\varepsilon\}, \{a\}, \{aa\}, \{\varepsilon, a\}, \{\varepsilon, aa\}, \{a, aa\}, \{\varepsilon, a, aa\}\}.$$



## Définition 4.40

Une **classe** de langages est un ensemble de langages.

La plupart du temps, on n'étudie que les classes associées à certains modèles calculatoires.

Par exemple, la classe des langages qui sont reconnus par un type de machine ayant telles ressources et telles limitations...

Le terme classe de complexité est souvent utilisé dans la littérature.

### Définition 4.41

**FINI** désigne la classe de tous les **langages finis**, c'est-à-dire l'ensemble des langages qui contiennent un nombre fini de mots.

## Exemple 4.42

*Si*

$$\Sigma = \{a, b\}$$

$$L = \{abba, baba, \varepsilon\},$$

*alors*

$$L \in \text{FINI}$$

$$\emptyset \in \text{FINI}$$

$$\Sigma^* \notin \text{FINI}$$

### Définition 4.43

On désigne la classe de tous les langages par **UNIVERS**.

# Fermeture d'une classe de langages par rapport à une opération

## Définition 4.44

On dit qu'une classe est fermée par rapport à une opération, ou fermée pour une opération, si le résultat de l'opération est toujours un élément de la classe lorsque les opérandes sont pris dans la classe.

### **Théorème 4.45**

*La classe FINI est fermée pour l'union.*

#### **Démonstration :**

*Il suffit de voir que l'union de deux langages finis est un langage fini.*

### **Théorème 4.46**

*La classe FINI est fermée pour l'intersection.*

#### **Démonstration :**

*L'intersection de deux langages finis donne toujours un langage fini.*

### **Théorème 4.47**

*La classe FINI n'est pas fermée pour l'opération complément.*

#### **Démonstration :**

*On a :  $\emptyset \in \text{FINI}$ , mais  $\overline{\emptyset} = \Sigma^* \notin \text{FINI}$ .*

### **Remarque 4.48**

*Le complément d'un langage fini, quel qu'il soit, ne donne jamais un langage fini.*



# Résumé

Concept	Définition	Peut-être vide ?	Peut-être infini ?	Exemples
symbole	—	—	—	a
alphabet	ensemble de symboles	non	non	{a, b, c}
mot	suite de symboles	oui ( $\varepsilon$ )	non	abaac
langage	ensemble de mots	oui ( $\emptyset$ )	oui	{ $\varepsilon$ , ba, bba, bbba, ...}
classe	ensemble de langages	oui (inintéressant)	oui (toutes les classes intéressantes sont infinies)	FINI REG HC P

## Définition 4.49 (Automate fini déterministe)

Un **automate fini** est un quintuplet  $(Q, \Sigma, \delta, q_0, F)$  où :

- $Q$  est un ensemble fini d'**états** ;
- $\Sigma$  est un alphabet ;
- $\delta : Q \times \Sigma \rightarrow Q$  est la **fonction de transition** ;
- $q_0 \in Q$  est l'**état initial** ;
- $F \subseteq Q$  est l'ensemble **états acceptants**.

### Exemple 4.50

*Soit*

$$M = (Q, \Sigma, \delta, q_0, F),$$

*où*

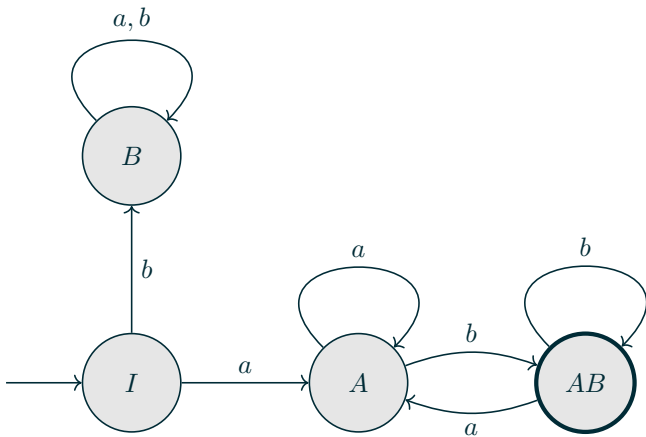
$$Q = \{I, A, B, AB\},$$

$$\Sigma = \{a, b\},$$

$$q_0 = I,$$

$$F = \{AB\},$$

*et où la fonction  $\delta$  de transition est donnée par le diagramme de la Figure ??.*



Le langage

$$\{w \mid w = \mathbf{a}x\mathbf{b} \text{ avec } x \in \{\mathbf{a}, \mathbf{b}\}^*\}.$$

est l'ensemble des mots qui amènent l'automate  $M$  de son état initial à l'état  $AB$ .

### Définition 4.51

Soit  $M = (Q, \Sigma, \delta, q_0, F)$  un automate.

Une **configuration** de l'automate  $M$  est une paire  $[w, q]$  où  $w \in \Sigma^*$  et  $q \in Q$ .

### Définition 4.52

Soit  $M = (Q, \Sigma, \delta, q_0, F)$  un automate.

Si  $a \in \Sigma$  et  $\delta(q, a) = q'$ , alors

$$[a \cdot w, q] \vdash [w, q']$$

dénote une transition de l'automate  $M$  de la configuration  $[a \cdot w, q]$  à la configuration  $[w, q']$ .

### Définition 4.53

Soit  $M = (Q, \Sigma, \delta, q_0, F)$  un automate.

Pour  $x \in \Sigma^*$  et  $w \in \Sigma^*$ ,

$$[x \cdot w, q] \stackrel{n}{\vdash} [w, q']$$

dénote  $n$  transitions successives amenant l'automate  $M$  de la configuration  $[x \cdot w, q]$  à la configuration  $[w, q']$ , et

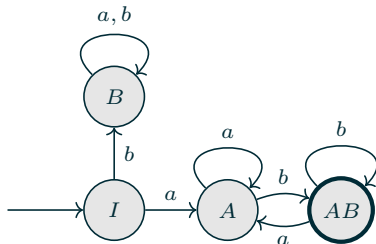
$$[x \cdot w, q] \stackrel{*}{\vdash} [w, q']$$

indique que

$$[x \cdot w, q] \stackrel{n}{\vdash} [w, q']$$

pour  $n = |x|$ .





$$\begin{array}{ll}
 [abab, I] & \vdash [bab, A] \\
 & \vdash [ab, AB] \\
 & \vdash [b, A] \\
 & \vdash [\varepsilon, AB]
 \end{array}
 \qquad
 \begin{array}{ll}
 [aab, I] & \vdash [ab, A] \\
 & \vdash [b, A] \\
 & \vdash [\varepsilon, AB]
 \end{array}$$

$$\begin{array}{ll}
 [ba, I] & \vdash [a, B] \\
 & \vdash [\varepsilon, B]
 \end{array}$$

### Définition 4.54

Soit  $M = (Q, \Sigma, \delta, q_0, F)$  un automate.

On dira que l'automate  $M$  **accepte** le mot  $w \in \Sigma^*$  s'il  $\exists q \in F$  tel que

$$[w, q_0] \stackrel{*}{\vdash} [\varepsilon, q].$$

### Définition 4.55

Soit  $M = (Q, \Sigma, \delta, q_0, F)$  un automate.

L'ensemble

$$L(M) = \{w \in \Sigma^* \mid M \text{ accepte } w\}$$

est le **langage reconnu** par l'automate  $M$ .

### Définition 4.56

Un langage  $K$  est **régulier** s'il existe un automate  $M$  tel que  $K = L(M)$ .

### Définition 4.57

**REG** désigne la classe des langages réguliers

## Exemple de langage régulier

Soit

$$M = (Q, \Sigma, \delta, q_0, F)$$

où

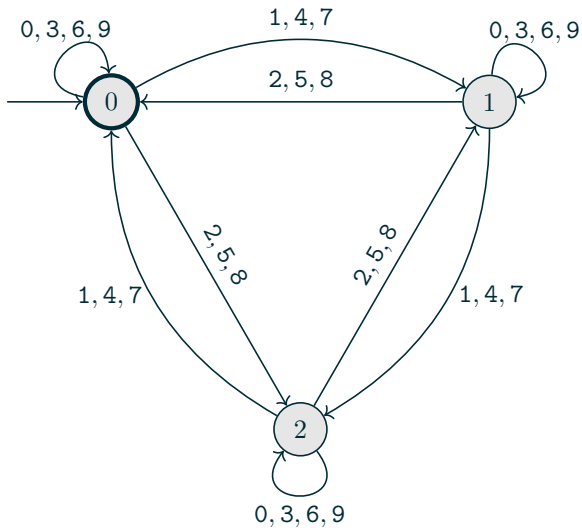
$$Q = \{\langle 0 \rangle, \langle 1 \rangle, \langle 2 \rangle\}$$

$$\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

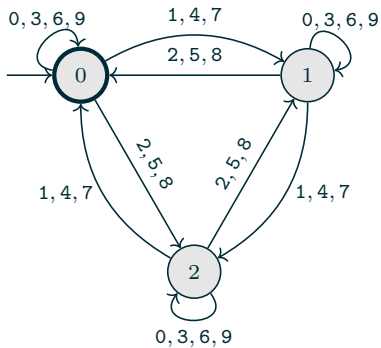
$$q_0 = \langle 0 \rangle$$

$$F = \{\langle 0 \rangle\},$$

et où la fonction de transition  $\delta$  est donnée par le diagramme de la Figure 4.1



**Fig. 4.1 :** Un exemple d'automate fini.



$$\begin{aligned}
 [68725, \langle 0 \rangle] &\vdash [8725, \langle 0 \rangle] \\
 &\vdash [725, \langle 2 \rangle] \\
 &\vdash [25, \langle 0 \rangle] \\
 &\vdash [5, \langle 2 \rangle] \\
 &\vdash [\varepsilon, \langle 1 \rangle].
 \end{aligned}$$

Donc  $68725 \notin L(M)$ .

### Proposition 4.58

*L'automate de la Figure 4.1 reconnaît le langage contenant le mot vide et les chaînes de chiffres décimaux qui forment un nombre divisible par 3.*

#### Démonstration :

L'inspection du diagramme de la fonction de transition  $\delta$  révèle que  $M$  accepte  $\varepsilon$  de même que les nombres dont la somme des chiffres est congrue à 0 modulo 3.

Or, si le développement décimal de l'entier  $n$  est  $d_i d_{i-1} \cdots d_0$ , on a :

$$n \equiv 0 \pmod{3}$$

$$\Leftrightarrow \sum_{j=0}^i d_j 10^j \equiv 0 \pmod{3}$$

$$\Leftrightarrow \sum_{j=0}^i d_j \equiv 0 \pmod{3} \quad (\text{car } 10^j \cong 1)$$

□



### **Théorème 4.59**

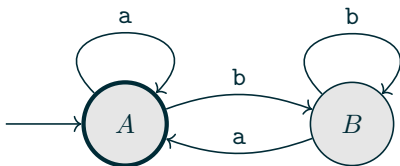
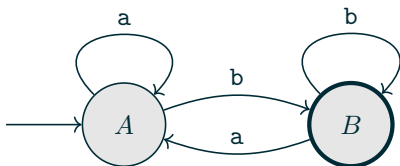
*La classe REG est fermée pour l'opération complément.*

#### **Démonstration :**

Si  $K \in \text{REG}$ , alors il est reconnu par un automate  $M$ .

Soit  $M'$  l'automate construit comme  $M$ , mais où les états acceptants et non-acceptants ont été interchangés.

Alors  $M'$  reconnaît  $\overline{K}$ . Donc  $\overline{K} \in \text{REG}$  et la classe est fermée sur cette opération. □



**Fig. 4.2 :** Exemple d'automate pour l'opération complément.

### **Théorème 4.60**

*La classe REG est fermée pour l'intersection.*

### Démonstration :

Soient  $K_1 \in \text{REG}$  et  $K_2 \in \text{REG}$  deux langages réguliers reconnus respectivement par les automates  $M_1$  et  $M_2$  :

$$\begin{aligned}M_1 &= (Q_1, \Sigma, \delta_1, q_{0,1}, F_1) \\M_2 &= (Q_2, \Sigma, \delta_2, q_{0,2}, F_2).\end{aligned}$$

Considérons l'automate  $M$  suivant :

$$M = (Q_1 \times Q_2, \Sigma, \delta, (q_{0,1}, q_{0,2}), F_1 \times F_2),$$

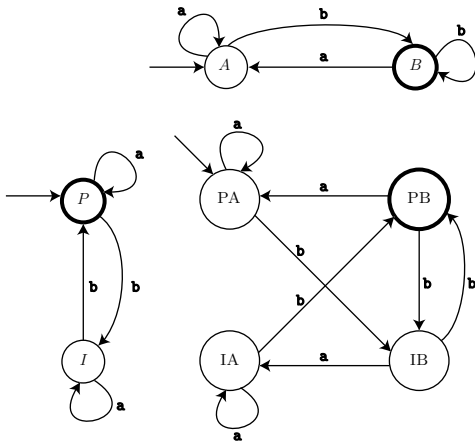
où

$$\delta((q_1, q_2), x) = (q'_1, q'_2)$$

si et seulement si

$$\delta_1(q_1, x) = q'_1 \quad \text{et} \quad \delta_2(q_2, x) = q'_2,$$

comme dans l'exemple de la Figure 4.3.



**Fig. 4.3 :** Exemple d'automate pour l'intersection.

Soit  $w \in K_1 \cap K_2$ .

Alors  $w$  amène  $M_1$  de  $q_{0,1}$  à un état final  $f_1 \in F_1$  et amène également  $M_2$  de  $q_{0,2}$  à  $f_2 \in F_2$ .

Par construction, le mot  $w$  amène l'automate  $M$  de  $(q_{0,1}, q_{0,2})$  à  $(f_1, f_2) \in F_1 \times F_2$ .

Par conséquent,  $M$  accepte  $w$  et  $w \in L(M)$ .

Inversement, si  $f_1 \notin F_1$  ou  $f_2 \notin F_2$ , alors  $(f_1, f_2) \notin F_1 \times F_2$  et  $M$  n'accepte pas  $w$ .

On a donc  $L(M) = K_1 \cap K_2$  et  $K_1 \cap K_2 \in \text{REG}$ .

□

### **Théorème 4.61**

*La classe REG est fermée sous l'union.*

#### **Démonstration :**

Soient  $K_1 \in \text{REG}$  et  $K_2 \in \text{REG}$ , deux langages réguliers.

On a :

$$K_1 \cup K_2 = \overline{\overline{K_1} \cap \overline{K_2}}.$$

Comme la classe REG est fermée pour le complément et l'intersection, (théorèmes 4.59 et 4.60) alors les langages  $\overline{K_1}$ ,  $\overline{K_2}$ ,  $\overline{K_1} \cap \overline{K_2}$ , et enfin  $\overline{\overline{K_1} \cap \overline{K_2}}$  sont réguliers.

On a donc  $K_1 \cup K_2 \in \text{REG}$ .



## **Théorème 4.62**

$$\text{FINI} \subset \text{REG}$$

### **Démonstration :**

Tout d'abord, on remarque que l'inclusion doit être stricte, étant donné que le langage  $L = \Sigma^*$ , qui contient une infinité de mots, est régulier.

Ensuite, il est facile de montrer que le langage vide est régulier.

Notons en plus qu'un automate accepte le mot vide si et seulement si son état initial est acceptant.

Nous allons maintenant considérer les langages finis, non vides, qui ne contiennent pas le mot vide.



Soient  $L = \{w_1, \dots, w_k\}$  un langage fini de  $k$  mots,  $k \geq 1$  et  $\varepsilon \notin L$ , avec

$$w_1 = x_{1,1} x_{1,2} \dots x_{1,l_1}$$

$$w_2 = x_{2,1} x_{2,2} \dots x_{2,l_2}$$

$$\vdots$$

$$w_k = x_{k,1} x_{k,2} \dots x_{k,l_k}$$

où  $l_i = |w_i|$  et  $x_{i,j} \in \Sigma$ .

Soit

$$l = \max_{1 \leq i \leq k} l_i$$

la longueur maximale des mots de  $L$ .

Considérons l'automate suivant :

$$M = (Q, \Sigma, \delta, q_0, F)$$

où

$$Q = \{\langle \text{INIT} \rangle\} \cup \Sigma \cup (\Sigma \times \Sigma) \cup \dots \cup \Sigma^l \cup \{\langle \text{PUITS} \rangle\}$$

$$q_0 = \langle \text{INIT} \rangle$$

$$F = \{\langle x_{1,1} x_{1,2} \dots x_{1,l_1} \rangle, \langle x_{2,1} x_{2,2} \dots x_{2,l_2} \rangle, \dots, \langle x_{k,1} x_{k,2} \dots x_{k,l_k} \rangle\}$$

et où la fonction  $\delta$  de transition est donnée par :

$\forall z, z_1, z_2, \dots, z_l \in \Sigma :$

$$\delta(\langle \text{INIT} \rangle, z) = \langle z \rangle$$

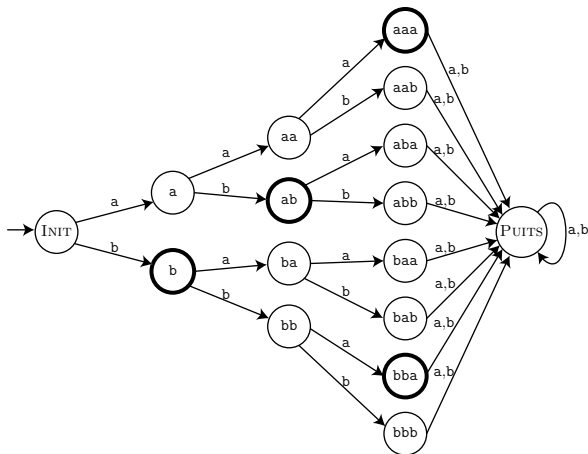
$$\delta(\langle z_1 z_2 \dots z_h \rangle, z) = \langle z_1 z_2 \dots z_h z \rangle, \quad \text{si } h < l$$

$$\delta(\langle z_1 z_2 \dots z_l \rangle, z) = \langle \text{PUITS} \rangle$$

$$\delta(\langle \text{PUITS} \rangle, z) = \langle \text{PUITS} \rangle$$

comme dans l'exemple de la Figure 4.4.

$\Sigma = \{a, b\}$     $L = \{b, ab, aaa, bba\}$     $k = 4$     $l = 3$



**Fig. 4.4 :** Exemple d'automate acceptant un langage fini.

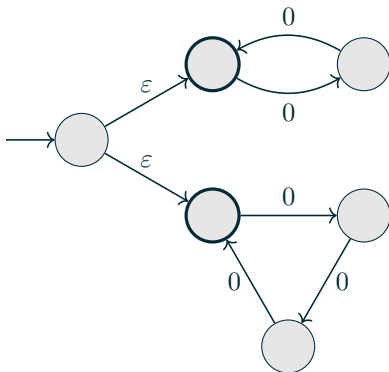
L'automate ainsi construit acceptera les mots du langage  $L$  et uniquement ces mots-là.

Donc  $L \in \text{REG}$  et  $\text{FINI} \subseteq \text{REG}$ .



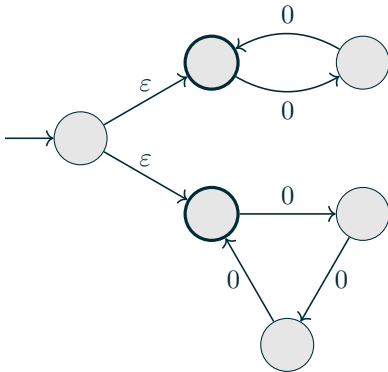
# Automates finis non-déterministes

Un automate fini *non-déterministe* (AFN) est un automate fini où, au lieu de spécifier la transition qui se produit à chaque étape, on spécifie les transitions *possibles*.



Ici,  $\Sigma = \{0\}$ ,  $\epsilon$  signifie une transition « gratuite » (qui peut s'effectuer sans lire de caractère).

Lors de l'exécution, l'automate « essaye » tous les chemins possibles, et si au moins *un* chemin possible mène à l'acceptation du mot, l'automate accepte.



Cet automate accepte les chaînes de zéros de longueur divisible par 2 ou par 3.

Autres exemples :



$$L(M) = \emptyset$$



$$L(M) = \{\varepsilon\}$$



$$L(M) = \{a\}$$



### Définition 4.63 (Automate fini non-déterministe)

Un automate fini non-déterministe (AFN) est un 5-tuple  $(Q, \Sigma, \delta, q_0, F)$ , où

- $Q$  est un ensemble d'états
- $\Sigma$  est un alphabet
- $\delta : Q \times (\Sigma \cup \{\varepsilon\}) \rightarrow \mathcal{P}(Q)$  est la fonction de transition
- $q_0 \in Q$  est l'état de départ
- $F \subseteq Q$  est l'ensemble d'états acceptants.

⇒ même définition que les automates déterministes, mais  $\delta$  prend en entrée l'état et le prochain symbole **ou la chaîne vide** et produit **l'ensemble des états possibles** à la prochaine étape.

## **Théorème 4.64**

*La classe des langages reconnus par un AFN est REG*

Les AFN ne sont donc pas plus puissants que les AFD.

### Démonstration :

Clairement, les AFN reconnaissent tous les langages dans REG, car un AFD est aussi un AFN. On doit donc prouver que pour tout AFN, il existe un AFD qui accepte le même langage.

Soit  $M = (Q, \Sigma, \delta, q_0, F)$  un AFN quelconque. On construit un AFD  $M' = (Q', \Sigma, \delta', q'_0, F')$  avec  $Q' = \mathcal{P}(Q)$  avec la propriété qu'à chaque étape du calcul, l'état  $q \in \mathcal{P}(Q)$  de l'AFD contient l'ensemble des états dans lequel l'AFN pourrait se trouver à cette étape.

Comment construire l'automate formellement ?

Supposons d'abord qu'il n'y a pas de transition gratuites (donc pas de flèches  $\varepsilon$ ). Alors, pour tout état de l'AFD  $S \subseteq Q$  et tout symbole  $a \in \Sigma$ ,

$$\begin{aligned}\delta'(S, a) &= \{q \mid q \text{ est atteignable en une étape depuis } S\} \\ &= \bigcup_{q \in S} \delta(q, a).\end{aligned}$$

L'état de départ est alors :

$$q'_0 = \{q_0\},$$

et les états acceptants sont tous les états contenant un état acceptant :

$$F' = \{S \subseteq Q \mid S \text{ contient un état } q \in F\}.$$

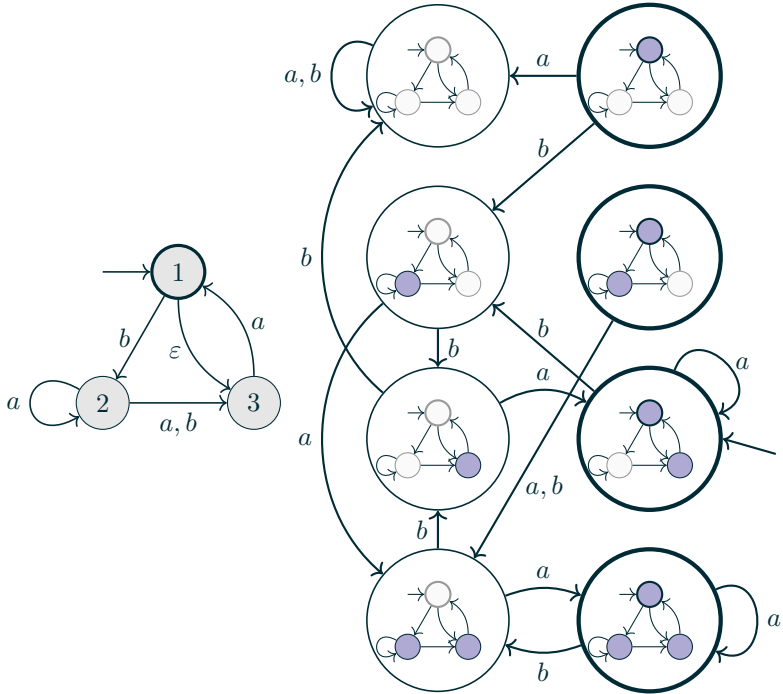
Si on a des flèches  $\varepsilon$ , alors définissons pour tout ensemble d'états  $S \subseteq Q$ , l'ensemble des états atteignables par des flèches  $\varepsilon$  :

$$E(S) = \{q \mid q \text{ peut être atteint avec 0 ou plus flèches } \varepsilon \text{ depuis } S\}.$$

On peut donc ajuster la définition précédente :

$$\begin{aligned}\delta'(S, a) &= \{q \mid q \text{ est atteignable en une étape depuis } S\} \\ &= \bigcup_{q \in S} E(\delta(q, a)) \\ q'_0 &= E(\{q_0\}).\end{aligned}$$





### **Théorème 4.65**

*La classe REG est fermée sous la concaténation*

## Démonstration :

Soient  $L_1$  et  $L_2$  deux langages réguliers sur le même alphabet. On doit démontrer que  $L_1 \circ L_2$  est également régulier.

Soient  $M_1 = (Q_1, \Sigma, \delta_1, q_{0,1}, F_1)$  et  $M_2 = (Q_2, \Sigma, \delta_2, q_{0,2}, F_2)$  deux AFN acceptant  $L_1$  et  $L_2$  respectivement. Nous allons construire un AFN  $M = (Q, \Sigma, \delta, q_0, F)$  qui accepte  $L_1 \circ L_2$  en « concaténant »  $M_1$  et  $M_2$  :

$$Q = Q_1 \cup Q_2$$

$$q_0 = q_{0,1}$$

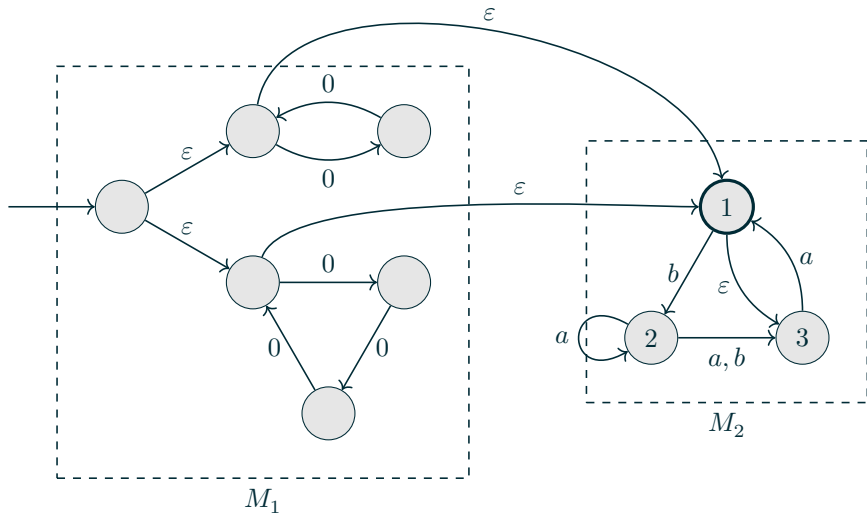
$$F = F_2$$

$$\delta(q, a) = \begin{cases} \delta_1(q, a) & \text{si } q \in Q_1 \text{ et } q \notin F_1 \\ \delta_1(q, a) & \text{si } q \in F_1 \text{ et } a \neq \varepsilon \\ \delta_2(q, a) & \text{si } q \in Q_2 \\ \delta_1(q, a) \cup \{q_{0,2}\} & \text{si } q \in F_1 \text{ et } a = \varepsilon. \end{cases}$$

Pour passer de  $M_1$  à  $M_2$ , il faut avoir lu un mot dans  $L_1$ , et ensuite pour accepter, il faut avoir lu un mot dans  $L_2$ . □



Exemple :



### **Théorème 4.66**

*La classe REG est fermée sous l'étoile*

Ici, pour tout langage  $L$ ,

$$L^i = \{w_1 w_2 \dots w_i \mid w_j \in L \ \forall i \geq 0\}$$

$$L^* = \bigcup_{i=0}^{\infty} L^i.$$

### Démonstration :

Soit  $M = (Q, \Sigma, \delta, q_0, F)$  un AFN reconnaissant  $L$ . Alors, on construit un AFN  $M' = (Q', \Sigma, \delta', q'_0, F')$  reconnaissant  $L^*$  :

$$Q' = Q \cup \{\langle init \rangle\}$$

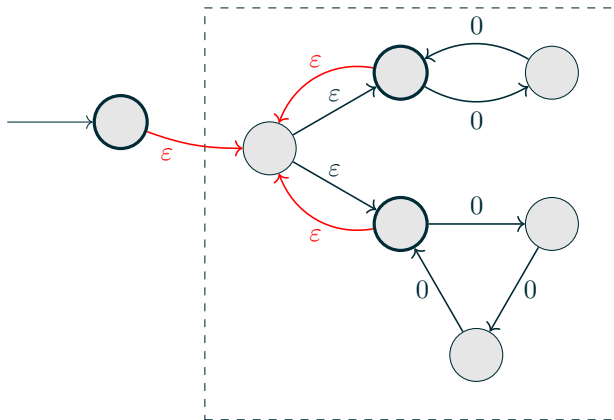
$$q'_0 = \langle init \rangle$$

$$F' = F \cup \{\langle init \rangle\}$$

$$\delta'(q, a) = \begin{cases} q_0 & \text{si } q = \langle init \rangle \text{ et } a = \varepsilon \\ \delta(q, a) \cup \{q_0\} & \text{si } q \in F \text{ et } a = \varepsilon \\ \delta(q, a) & \text{sinon.} \end{cases}$$

Donc, on a le même automate, mais on ajoute des flèches  $\varepsilon$  de tous les états acceptants à l'état de départ, et on ajoute un nouvel état de départ acceptant.

Exemple :



## Définition 4.67 (Expression régulière)

$R \subseteq (\Sigma \cup \{\varepsilon, \emptyset, \cup, \circ, *\})^*$  est une expression régulière si  $R$  est

1.  $a$  où  $a \in \Sigma$ ,
2.  $\varepsilon$ ,
3.  $\emptyset$ ,
4.  $(R_1 \cup R_2)$ , où  $R_1$  et  $R_2$  sont des expressions régulières,
5.  $(R_1 \circ R_2)$ , où  $R_1$  et  $R_2$  sont des expressions régulières,
6.  $(R_1)^*$ , où  $R_1$  est une expression régulière.

On écrira également  $R_1 \circ R_2$  comme  $R_1 R_2$ .

Ces expressions représentent :

1.  $R = a$  : le langage  $\{a\}$
2.  $R = \varepsilon$  : le langage  $\{\varepsilon\}$
3.  $R = \emptyset$  : le langage vide
4.  $(R_1 \cup R_2)$  : l'union des langages représentés par  $R_1$  et  $R_2$
5.  $(R_1 \circ R_2)$  : la concaténation des langages représentés par  $R_1$  et  $R_2$
6.  $(R_1)^*$  : l'étoile du langage représenté par  $R_1$ .

Ici l'étoile est défini comme :

$$R^* := \{w_1 w_2 w_3 \cdots w_n \mid n \geq 0, \forall i, w_i \in R\}.$$

Voici quelques exemples d'expressions régulières :

1.  $a \cup b$  : le langage fini  $\{a, b\}$
2.  $0^*10^*$  : le langage qui ne contient qu'un seul 1
3.  $(00)^* \cup (000)^*$  : le langage qui contient un nombre de zéro divisible par 2 ou par 3
4.  $((0 \cup 1)(0 \cup 1))^*$  : le langage qui contient un nombre pair de bits
5.  $0^*\emptyset$  : le langage vide

Quelle est la classe des langages reconnus pouvant être exprimés par une expression régulière ?

## **Théorème 4.68**

*$L \in \text{REG}$  si et seulement si  $L$  peut être exprimé par une expression régulière.*

On prouvera les deux directions de ce théorème via deux lemmes séparés.



#### **Lemme 4.69**

*Si  $L$  peut être exprimé par une expression régulière, alors  $L \in \text{REG}$ .*

**Démonstration :**

Par induction structurelle sur l'expression régulière. On peut voir l'expression régulière comme une structure définie récursivement avec

$$\mathcal{S}_0 = \{\emptyset, \varepsilon\} \cup \Sigma,$$

et avec trois règles de construction pour l'union, la concaténation, et l'étoile.

### Base de l'induction :

- $R = \emptyset$  : le langage vide est clairement régulier.
- $R = \varepsilon$  : le mot vide est clairement régulier.
- $R = \{a\}$  où  $a \in \Sigma$  : aussi clairement régulier.

**Pas d'induction :** Supposons que  $R_1$  et  $R_2$  sont réguliers. Alors :

- $R_1 \cup R_2$  est régulier, par la fermeture sous l'union.
- $R_1 \circ R_2$  est régulier, par la fermeture sous la concaténation.
- $R_1^*$  est régulier, par la fermeture sous l'étoile.



### **Lemme 4.70**

*Si  $L \in \text{REG}$ , alors il existe une expression régulière  $R$  qui exprime  $L$ .*

**Démonstration :**

L'idée consiste à prendre un AFN reconnaissant  $L$ , et de le transformer étape par étape en une expression régulière. On va donc généraliser les AFN pour permettre les expressions régulières sur les flèches. Ensuite, on enlèvera les états un par un dans le but d'arriver à un automate ayant seulement deux états et l'expression régulière qu'on veut sur la seule flèche restante.

### Définition 4.71 (AFNER)

Un automate fini non-déterministe à expression régulière (AFNER) est un tuple  $(Q, \Sigma, \delta, q_0, q_F)$  où

- $Q$  est l'ensemble des états
- $\Sigma$  est l'alphabet
- $\delta : Q \times Q \rightarrow \mathcal{R}(\Sigma)$  est la fonction de transition, où  $\mathcal{R}(\Sigma)$  représente l'ensemble des expressions régulières sur  $\Sigma$ .
- $q_0$  est l'état de départ
- $q_F$  est l'état final ( $q_0 \neq q_F$ )

De plus nous requerrons que  $q_0$  n'ait aucune flèche entrante ( $\delta(q, q_0) = \emptyset$  pour tout  $q$ ), et que  $q_F$  n'ait aucune flèche sortante ( $\delta(q_F, q) = \emptyset$  pour tout  $q$ ).

Par convention, une flèche étiquetée  $\emptyset$  sera représentée par une absence de flèche.

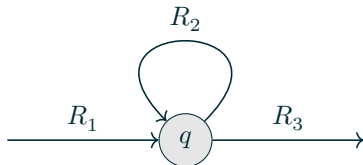
Pour transformer un AFN en AFNER, on doit

1. Créer un nouvel état de départ avec une flèche  $\varepsilon$  vers l'ancien état de départ
2. Créer un nouvel état acceptant unique avec des flèches  $\varepsilon$  depuis les anciens états acceptants
3. Convertir les flèches de type  $a, b$  en  $a \cup b$

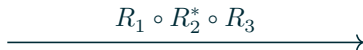


L'algorithme d'élimination :

1. On transforme l'AFN en AFNER.
2. On choisit un état  $q$  à éliminer (sauf  $q_0$  et  $q_F$ ), et pour toute paire (flèche entrante, flèche sortante), on remplace :

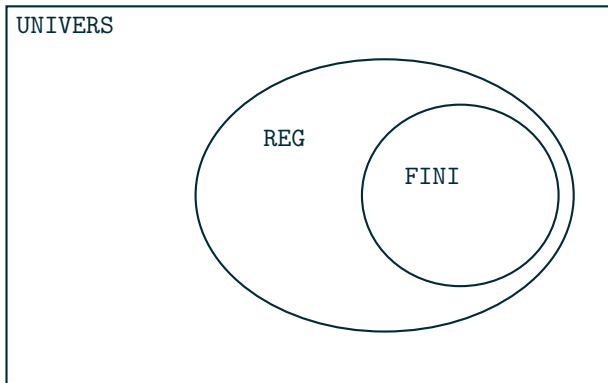


par :



3. Si des flèches parallèles sont créées, on les remplace par une seule flèche avec l'union des flèches parallèles.
4. On répète jusqu'à ce qu'il ne reste que  $q_0$  et  $q_F$ .
5. Output : l'expression régulière sur la flèche de  $q_0$  à  $q_F$ .





# Pour prouver qu'un langage $L$ est régulier

On peut...

- construire un automate qui reconnaît  $L$ ,
- utiliser les propriétés de fermeture de la classe REG,
- montrer que  $L$  est fini.

### **Théorème 4.72 (Lemme du pompiste)**

*Soit  $L \in \text{REG}$ . Il existe un entier  $p \geq 1$  tel que pour tout  $w \in L$  avec  $|w| \geq p$ , on peut écrire  $w = x \cdot y \cdot z$  où*

- $|y| > 0$ ,
- $|x \cdot y| \leq p$ ,
- $\forall i \geq 0 : x \cdot y^i \cdot z \in L$ .

**Démonstration :**

Soit  $M = (Q, \Sigma, \delta, q_0, F)$  un AFD qui reconnaît le langage  $L$ .

Choisissons  $p = \#Q$ .

Soit  $w \in L$  tel que  $|w| \geq p$ .

L'automate  $M$  va passer par au moins  $p + 1$  états avant d'accepter  $w$ .

Par le principe du pigeonnier, il existe un état par lequel  $M$  va passer au moins deux fois, parmi les  $p + 1$  premiers états visités.

Soit  $r \in Q$  le premier état par lequel  $M$  passe deux fois.

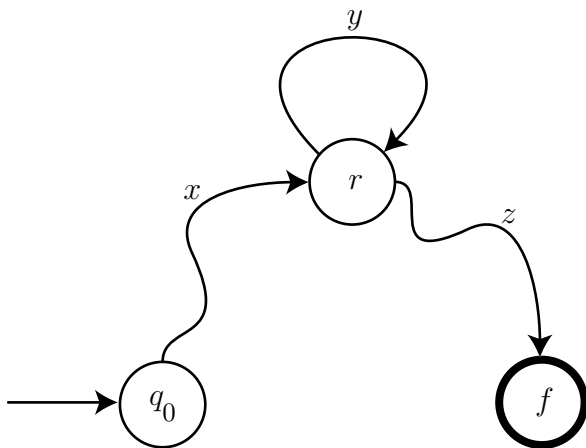
Soit  $x$  la partie de  $w$  qui est lue entre l'état initial  $q_0$  et la première visite de  $r$ .

Soit  $y$  la partie de  $w$  lue entre les deux premières visites de l'état  $r$ .

Soit  $z$  la partie restante de  $w$  après  $x \cdot y$ . On a donc :

$$\begin{aligned}w &= x \cdot y \cdot z, \\|y| &> 0, \\|x \cdot y| &\leq p,\end{aligned}$$

comme illustré dans la Figure 4.5.



**Fig. 4.5** : Illustration du lemme du pompiste.

Il est maintenant clair que les mots suivants seront acceptés par  $M$  :

$$x \cdot z$$

$$x \cdot y \cdot z$$

$$x \cdot y \cdot y \cdot z$$

$$x \cdot y \cdot y \cdot y \cdot z$$

$$\vdots$$

Ces mots appartiennent donc tous à  $L$ .



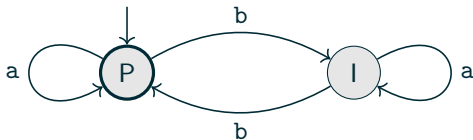


### Exemple 4.73

On sait que le langage

$$L = \{w \in \{a, b\}^* \mid |w|_b \text{ est pair}\}$$

est régulier parce qu'il existe un automate qui le reconnaît :



Cet automate a 2 états.

Par la preuve du lemme du pompiste, on peut prendre  $p = 2$ .

Considérons le mot  $w = abb$ .

On a bien  $w \in L$  et  $|w| \geq p$ .

Choisissons la décomposition  $w = x \cdot y \cdot z$  où

$$x = \varepsilon,$$

$$y = a,$$

$$z = bb.$$

On a bien :

- $|y| > 0,$
- $|x \cdot y| \leq p,$
- $bb, abb, aabb, aaabb, aaaabb, \dots \in L.$



# Pour prouver qu'un langage $L$ n'est pas régulier

On peut...

- utiliser les propriétés de fermeture de la classe REG,
- appliquer le lemme du pompiste.

### **Théorème 4.74**

*Le langage*

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

*n'est pas régulier.*

**Démonstration :**

Supposons au contraire que  $L \in \text{REG}$ , afin d'arriver à une contradiction.

Soit  $p \geq 1$  tel que donné par le lemme du pompiste.

Considérons le mot  $w = a^p b^p$ .

On a bien  $w \in L$  et  $|w| = 2p \geq p$ .

Soient  $w = x \cdot y \cdot z$ ,  $|y| > 0$ ,  $|x \cdot y| \leq p$ , tels que donnés par le lemme du pompiste.

Il est clair que  $y$  ne contient que des  $a$ .

Considérons  $i = 2$  dans l'expression  $x \cdot y^i \cdot z$ .

Ce mot contient plus de a que de b, et donc  $x \cdot y^2 \cdot z \notin L$ , en contradiction avec l'énoncé du lemme.

Donc  $L \notin \text{REG}$ .



### **Théorème 4.75**

*Le langage*

$$L = \{w \in \{\mathbf{a}, \mathbf{b}\}^* \mid |w|_{\mathbf{a}} = |w|_{\mathbf{b}}\}$$

*n'est pas régulier.*

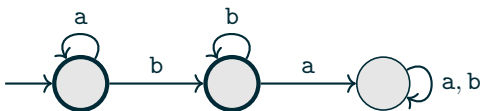
### Démonstration :

Supposons au contraire que  $L \in \text{REG}$ , afin d'arriver à une contradiction.

Considérons le langage

$$R = \{a^k b^l \mid k \geq 0, l \geq 0\}.$$

Il est facile de voir que  $R \in \text{REG}$ , car il est reconnu par l'automate suivant :

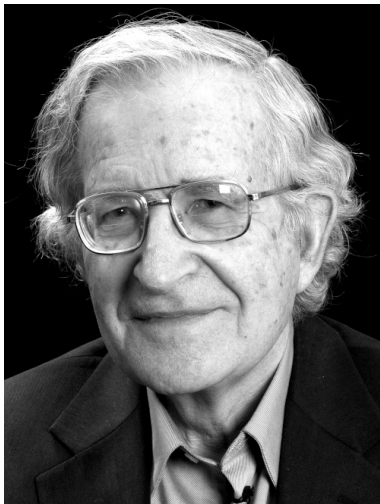




On a donc  $L \cap R \in \text{REG}$ , par la fermeture de la classe REG par rapport à l'opération d'intersection (théorème 4.60).

Mais  $L \cap R = \{a^n b^n \mid n \geq 0\}$  n'est pas régulier par le théorème 4.74.

On a une contradiction et le théorème est démontré. □



Exemple de grammaire hors-contexte :

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$

Les grammaires hors-contexte ont été définies par Noam Chomsky dans les années 50 :

#### Définition 4.76 (Grammaire hors-contexte)

Une **grammaire hors contexte (GHC)** est un quadruplet  $(V, \Sigma, R, S)$  où :

- $V$  est un ensemble fini de **variables** ;
- $\Sigma$  est un alphabet, c'est-à-dire un ensemble non vide et fini de symboles appelés **terminaux**,  $V \cap \Sigma = \emptyset$  ;
- $R$  est un ensemble fini de **règles** de la forme  $v \rightarrow z$  où  $v \in V$  et  $z \in (V \cup \Sigma)^*$  ;
- $S \in V$  est la variable **initiale**.

### Exemple 4.77

$$G = (V, \Sigma, R, S)$$

où :

$$V = \{S, X\}$$

$$\Sigma = \{a, b, c\}$$

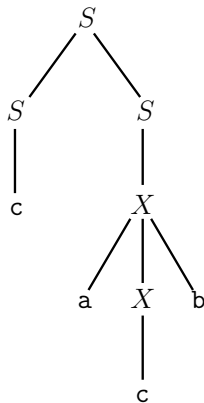
$$R = \{S \rightarrow SS \mid X \mid c, \quad X \rightarrow aXb \mid c\}$$

$$S \rightarrow SS \mid X \mid c$$

$$X \rightarrow aXb \mid c$$

Dérivation du mot *cacb* :

$$\begin{aligned} S &\Rightarrow SS \\ &\Rightarrow cS \\ &\Rightarrow cX \\ &\Rightarrow caXb \\ &\Rightarrow cacb \end{aligned}$$



**Fig. 4.6 :** Arbre de dérivation

Exemple linguistique :

$$\langle \text{PHRASE} \rangle \rightarrow \langle \text{GROUPE NOMINAL} \rangle \langle \text{GROUPE VERBAL} \rangle$$
$$\langle \text{GROUPE NOMINAL} \rangle \rightarrow \langle \text{NOM} \rangle$$
$$\langle \text{GROUPE VERBAL} \rangle \rightarrow \langle \text{VERBE} \rangle \mid \langle \text{VERBE} \rangle \langle \text{COMPLEMENT} \rangle$$
$$\langle \text{COMPLEMENT} \rangle \rightarrow \langle \text{NOM} \rangle \mid \text{que } \langle \text{PHRASE} \rangle$$
$$\langle \text{NOM} \rangle \rightarrow \text{la table} \mid \text{Paulette} \mid \text{Jean} - \text{Pierre}$$
$$\langle \text{VERBE} \rangle \rightarrow \text{voit} \mid \text{dit}$$

On peut dériver des phrases :

la table voit que Paulette dit Jean-Pierre

Grammaire HC pour les programmes RÉPÉTER :

$$\begin{aligned} S &\rightarrow \varepsilon \\ &| \langle \text{INCRÉMENTATION} \rangle S \\ &| \langle \text{AFFECTATION} \rangle S \\ &| \langle \text{RÉPÉTER} \rangle S \end{aligned}$$

$$\langle \text{INCRÉMENTATION} \rangle \rightarrow \text{inc}(V)$$

$$\langle \text{AFFECTATION} \rangle \rightarrow V \leftarrow V$$

$$\langle \text{RÉPÉTER} \rangle \rightarrow \text{répéter } V \text{ fois } [S]$$

$$V \rightarrow r_N$$

$$N \rightarrow C | CN$$

$$C \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$$



### Définition 4.78

Soit  $G = (V, \Sigma, R, S)$  une GHC et soient  $u, v, w \in (V \cup \Sigma)^*$ .

On dit que  $uAv$  **donne** ou **produit** ou **engendre**  $uwv$ , noté

$$uAv \Rightarrow uwv,$$

si  $A \rightarrow w \in R$ .

On note  $u \xRightarrow{n} v$  si  $u$  donne  $v$  en  $n$  étapes :

$$u \Rightarrow u_1 \Rightarrow u_2 \Rightarrow \cdots \Rightarrow u_n = v.$$

Aussi,  $u \xRightarrow{*} v$  s'il existe un  $n \geq 0$  tel que  $u \xRightarrow{n} v$ .

### Définition 4.79

Soit  $G = (V, \Sigma, R, S)$  une GHC.

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

est le langage engendré par  $G$ .

### Définition 4.80

Un langage  $L$  est **hors contexte** s'il existe une GHC  $G$  telle que  $L = L(G)$ .

### Définition 4.81

La classe des langages hors contexte est notée **HC**.

## Exemple 4.82

Soit  $G = (\{S\}, \{a, b\}, R, S)$  une GHC où

$$R = \{S \rightarrow aSb \mid \varepsilon\}.$$

Alors

$$\begin{aligned} L(G) &= \{\varepsilon, ab, aabb, aaabbb, \dots\} \\ &= \{a^n b^n \mid n \geq 0\}. \end{aligned}$$

Dérivation du mot  $aaabbb$  :

$$\begin{aligned} S &\Rightarrow aSb \\ &\Rightarrow aaSbb \\ &\Rightarrow aaaSbbb \\ &\Rightarrow aaabbb \end{aligned}$$

### Exemple 4.83

Soit  $G = (V, \Sigma, R, S)$  une GHC où

$$V = \{S, M, N\}$$

$$\Sigma = \{0, 1, \dots, 9, +, *, (, )\}$$

$$R : \begin{cases} S \rightarrow N \mid (S+S) \mid (S*S) \\ N \rightarrow 0 \mid 1M \mid 2M \mid \dots \mid 9M \\ M \rightarrow \varepsilon \mid 0M \mid 1M \mid \dots \mid 9M \end{cases}$$

La variable  $N$  donne les nombres :

$$N \Rightarrow 2M \Rightarrow 23M \Rightarrow 23.$$

La variable  $S$  donne les expressions arithmétiques :

$$S \Rightarrow (S+S) \Rightarrow (S+(S*S)) \xRightarrow{*} (N+(N*N)) \xRightarrow{*} (123 + (245 * 19))$$

### Proposition 4.84

Soit  $G = (\{S\}, \{a, b\}, R, S)$  une GHC où

$$R : S \rightarrow aSb \mid bSa \mid SS \mid \varepsilon,$$

et soit

$$L = \{w \mid |w|_a = |w|_b\},$$

alors

$$L(G) = L.$$

**Démonstration :**

Il est facile de voir que  $L(G) \subseteq L$  car toutes les règles de  $G$  produisent un nombre égal de  $a$  et de  $b$ .

Il reste à montrer que  $L \subseteq L(G)$ .

Soit  $w \in L$ .

La preuve fonctionne par induction sur la longueur  $n = |w|$ .



**Base de l'induction.** Soit  $n = 0$ . Alors :

$$\begin{aligned} S &\Rightarrow \varepsilon \\ &= w, \end{aligned}$$

donc  $w \in L(G)$ .

**Hypothèse d'induction.** Pour  $n > 0$ , supposons que

$$\forall z \in L, |z| < n : S \stackrel{*}{\Rightarrow} z.$$

**Pas d'induction.** Remarquons que tous les mots de  $L$  ont un nombre pair de symboles.

Considérons 4 cas pour  $w$  :

1.  $w = azb$ . Dans ce cas, on a la dérivation :

$$\begin{aligned} S &\Rightarrow aSb \\ &\stackrel{*}{\Rightarrow} azb \quad (\text{par hypothèse d'induction}) \\ &= w. \end{aligned}$$

2.  $w = \mathbf{bza}$  :

$$\begin{aligned} S &\Rightarrow \mathbf{bSa} \\ &\stackrel{*}{\Rightarrow} \mathbf{bza} \quad (\text{par hypothèse d'induction}) \\ &= w. \end{aligned}$$

3.  $w = aza$ . Dans ce cas,  $z$  contient deux  $b$  de plus que de  $a$ , et on peut donc le décomposer en  $z = z_1 \cdot z_2$ , où  $z_1$  et  $z_2$  contiennent chacun un  $b$  de plus que de  $a$ . D'où :

$$\begin{aligned} S &\Rightarrow SS \\ &\stackrel{*}{\Rightarrow} az_1 S \quad (\text{car } az_1 \in L \text{ et par hypothèse d'induction}) \\ &\stackrel{*}{\Rightarrow} az_1 z_2 a \quad (\text{car } z_2 a \in L \text{ et par hypothèse d'induction}) \\ &= aza \\ &= w. \end{aligned}$$

4.  $w = bz b$ . Ce cas se traite comme le cas 3, *mutatis mutandis*.

Dans tous les cas, on a  $w \in L(G)$ .



Construire une GHC pour un langage donné demande souvent de la créativité.

Cependant, il y a un certain nombre de techniques qui aident beaucoup.  
Entre autres :

- décomposer le langage en l'union de deux langages plus simples ;
- utiliser des règles qui produisent des listes ;

comme dans l'exemple qui suit.

### Exemple 4.85

Soit

$$L = \{a^i b^j a^k b^l \mid (i = k \text{ ou } i = l), j \geq 0, k \geq 1, l \geq 1\},$$

c'est-à-dire :

$$\begin{aligned} L &= \{a^i b^j a^i b^l \mid i \geq 1, j \geq 0, l \geq 1\} \\ &\cup \{a^i b^j a^k b^i \mid i \geq 1, j \geq 0, k \geq 1\}. \end{aligned}$$



Une GHC pour  $L$  est donnée par :

$$G = (\{S, S_1, S'_1, S_2, S'_2, A, B\}, \{a, b\}, R, S)$$

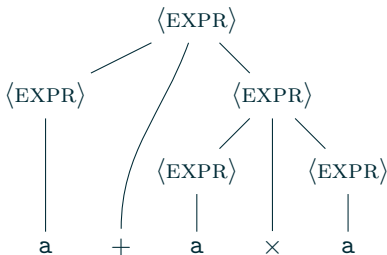
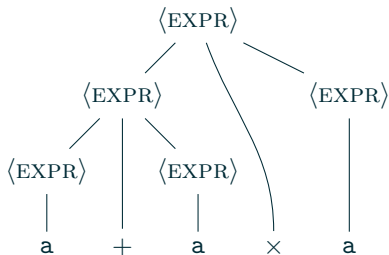
avec

$$R : \begin{cases} S & \rightarrow S_1 \mid S_2 \\ S_1 & \rightarrow S'_1 b B \\ S'_1 & \rightarrow a S'_1 a \mid a B a \\ B & \rightarrow b B \mid \varepsilon \\ S_2 & \rightarrow a S'_2 b \\ S'_2 & \rightarrow a S'_2 b \mid B a A \\ A & \rightarrow a A \mid \varepsilon \end{cases}$$

# Ambiguïté dans les GHC

Il est possible d'avoir plusieurs dérivations d'une même mot :

$$\begin{aligned}\langle \text{EXPR} \rangle &\rightarrow \langle \text{EXPR} \rangle + \langle \text{EXPR} \rangle \\ &\quad | \langle \text{EXPR} \rangle \times \langle \text{EXPR} \rangle \\ &\quad | (\langle \text{EXPR} \rangle) \\ &\quad | a\end{aligned}$$



## Définition 4.86

La GHC  $G = (V, \Sigma, R, S)$  est sous **forme normale de Chomsky (FNC)** si toutes ses règles sont de la forme

$$A \rightarrow BC \text{ ou}$$

$$A \rightarrow x,$$

avec

$$A, B, C \in V,$$

$$B \neq S,$$

$$C \neq S,$$

$$x \in \Sigma,$$

en permettant aussi

$$S \rightarrow \varepsilon.$$

### **Théorème 4.87**

*Il existe un algorithme qui prend en entrée une GHC  $G$  et qui retourne une grammaire en FNC  $G'$  telle que  $L(G) = L(G')$ .*

## Algorithme 4.88

Prendre en entrée :  $G = (V, \Sigma, R, S)$ , une GHC.

1. (INIT) Ajouter  $S_0$  à  $V$ , la nouvelle variable initiale, et ajouter la règle  $S_0 \rightarrow S$  à  $R$ .
2. (TERM) Éliminer les règles avec des terminaux non-isolés :
  - 2.1 Pour tout terminal  $a$  apparaissant à droite d'une règle avec d'autres symboles, on ajoute une variable  $N_a$  et la règle  $N_a \rightarrow a$ .
  - 2.2 On remplace toute règle de forme  $A \rightarrow X_1 \dots a \dots X_n$  par  $A \rightarrow X_1 \dots N_a \dots X_n$ .
3. (BIN) Éliminer de  $R$  toutes les règles de la forme  $A \rightarrow u_1 u_2 \dots u_k$  où  $k \geq 3$  et où  $u_1, u_2, \dots, u_k \in V \cup \Sigma$  :
  - 3.1 Enlever  $A \rightarrow u_1 u_2 \dots u_k$ .
  - 3.2 Ajouter  $A_1, A_2, \dots, A_{k-2}$  à  $V$ .
  - 3.3 Ajouter les règles  $A \rightarrow u_1 A_1, A_1 \rightarrow u_2 A_2, \dots, A_{k-2} \rightarrow u_{k-1} u_k$ .

4. (VIDE) Éliminer de  $R$  toutes les règles de la forme  $A \rightarrow \varepsilon$ , sauf possiblement la règle  $S_0 \rightarrow \varepsilon$  :
  - 4.1 Déterminer toutes les variables « annulables »  $A \in V$  telles que  $A \xRightarrow{*} \varepsilon$ .
  - 4.2 Pour toute règle de type  $A \rightarrow BC$ , ajouter  $A \rightarrow C$  si  $B$  est annulable et  $A \rightarrow B$  si  $C$  est annulable.
  - 4.3 Si  $S_0$  est annulable, ajouter  $S_0 \rightarrow \varepsilon$  ; éliminer toutes les autres règles de type  $A \rightarrow \varepsilon$ .
5. (UNAIRE) Éliminer les règles de type  $A \rightarrow B$  :
  - 5.1 On enlève  $A \rightarrow B$
  - 5.2 Pour chaque règle de type  $B \rightarrow w$  (où  $w \in (\Sigma \cup V)^*$ ), on ajoute  $A \rightarrow w$ , sauf si cette a déjà été enlevée, ou est égale à  $A \rightarrow A$ .

Produire en sortie la grammaire ainsi modifiée.



### Exemple 4.89

$$G = (\{S, A, B\}, \{a, b\}, R, S)$$

où

$$R : \begin{cases} S & \rightarrow ASA|aB \\ A & \rightarrow B|S \\ B & \rightarrow b|\varepsilon. \end{cases}$$

# 1. (INIT) Ajouter une nouvelle variable initiale

Avant :

$$S \rightarrow ASA \mid \mathbf{a}B$$

$$A \rightarrow B \mid S$$

$$B \rightarrow \mathbf{b} \mid \varepsilon$$

Après :

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid \mathbf{a}B$$

$$A \rightarrow B \mid S$$

$$B \rightarrow \mathbf{b} \mid \varepsilon$$



## 2. (TERM) Éliminer les règles avec des non-terminaux non-isolés :

Avant :

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid \mathbf{a}B$$

$$A \rightarrow B \mid S$$

$$B \rightarrow \mathbf{b} \mid \varepsilon$$

Après :

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid N_a B$$

$$A \rightarrow B \mid S$$

$$B \rightarrow \mathbf{b} \mid \varepsilon$$

$$N_a \rightarrow \mathbf{a}$$

### 3. (BIN) Éliminer les règles à plus de deux symboles à droite :

Avant :

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid N_a B$$

$$A \rightarrow B \mid S$$

$$B \rightarrow \mathbf{b} \mid \varepsilon$$

$$N_a \rightarrow \mathbf{a}$$

Après :

$$S_0 \rightarrow S$$

$$S \rightarrow AM \mid N_a B$$

$$M \rightarrow SA$$

$$A \rightarrow B \mid S$$

$$B \rightarrow \mathbf{b} \mid \varepsilon$$

$$N_a \rightarrow \mathbf{a}$$

#### 4. (VIDE) Éliminer les règles $A \rightarrow \varepsilon$ :

Avant :

$$S_0 \rightarrow S$$

$$S \rightarrow AM \mid N_a B$$

$$M \rightarrow SA$$

Variables annulables :  $A, B$

$$A \rightarrow B \mid S$$

$$B \rightarrow \mathbf{b} \mid \varepsilon$$

$$N_a \rightarrow \mathbf{a}$$

Après :

$$S_0 \rightarrow S$$

$$S \rightarrow AM \mid M \mid N_a B \mid N_a$$

$$M \rightarrow SA \mid S$$

$$A \rightarrow B \mid S$$

$$B \rightarrow \mathbf{b}$$

$$N_a \rightarrow \mathbf{a}$$

5. (UNAIRE) Éliminer les règles  $A \rightarrow B$  :

Avant :

$$S_0 \rightarrow S$$

$$S \rightarrow AM \mid M \mid N_a B \mid N_a$$

$$M \rightarrow SA \mid S$$

$$A \rightarrow B \mid S$$

$$B \rightarrow \mathbf{b}$$

$$N_a \rightarrow \mathbf{a}$$

Après :

$$S_0 \rightarrow AM \mid SA \mid N_a B \mid \mathbf{a}$$

$$M \rightarrow SA \mid AM \mid N_a B \mid \mathbf{a}$$

$$B \rightarrow \mathbf{b}$$

$$S \rightarrow AM \mid SA \mid N_a B \mid \mathbf{a}$$

$$A \rightarrow \mathbf{b} \mid AM \mid SA \mid N_a B \mid \mathbf{a}$$

$$N_a \rightarrow \mathbf{a}$$

### **Théorème 4.90**

*L'algorithme 4.88 produit une grammaire de taille  $O(n^2)$ , où  $n$  est la taille de la grammaire initiale (somme du nombre de symboles à droite des règles).*

## Démonstration :

1. Ajout de  $S_0$  : on ajoute une seule règle.
2. Élimination des règles avec non-terminaux isolés : dans le pire cas, on ajoute  $|\Sigma|$  règles, ce qui est plus petit que  $n$  si tous les éléments de  $\Sigma$  sont utilisés  $\Rightarrow O(n)$  jusqu'à maintenant.
3. Élimination des règles à plus de deux symboles : dans le pire cas, on double  $n$ , donc toujours linéaire jusqu'ici.
4. Élimination des  $\varepsilon$  : dans le pire cas, toutes les règles de forme  $A \rightarrow BC$  deviennent  $A \rightarrow BC \mid B \mid C$ , donc encore là on double  $n$  au pire.

4. Élimination des règles  $A \rightarrow B$  : au final, tout ce qui est à droite d'une règle doit déjà avoir été présent avant, possiblement avec une autre variable à gauche. Il y a au plus  $n$  variables, et les règles ayant une même variable  $A$  à gauche ont au total au plus  $n$  variables à droite. La grammaire finale est donc de taille au plus  $O(n^2)$ .

Les remarques 4.91 et 4.92 sont cruciales pour la compréhension de l'algorithme 4.94.

#### **Remarque 4.91**

*Si  $G = (V, \Sigma, R, S)$  est sous FNC, et si  $\varepsilon \in L(G)$ , alors la seule façon de produire  $\varepsilon$  est*

$$S \Rightarrow \varepsilon.$$



### Remarque 4.92

*Si  $G = (V, \Sigma, R, S)$  est sous FNC, si  $x \in \Sigma$  et  $x \in L(G)$ , alors la seule façon de produire  $x$  est*

$$A \Rightarrow x$$

*où  $A \in V$ .*

## **Théorème 4.93**

*Il existe un algorithme qui prend en entrée un mot  $w$  et une GHC  $G$  en FNC et qui détermine si  $w \in L(G)$ . Pour une grammaire  $G$  fixée, l'algorithme fonctionne en temps  $O(n^3)$ , où  $n$  est la taille du mot donné.*

### **Aperçu de la démonstration :**

Nous donnons l'algorithme et argumentons pour son temps d'exécution.

### Algorithme 4.94 (Cocke, Younger, Kasami (CYK))

*Prendre en entrée :  $G = (V, \Sigma, R, S)$ , une GHC, et  $w \in \Sigma^*$ . Soit  $w = w_1 w_2 \dots w_n$ , avec  $w_1, \dots, w_n \in \Sigma$ .*

1. *Si  $w = \varepsilon$ , alors :*

*si  $(S_0 \rightarrow \varepsilon) \in R'$ , alors retourner  $\langle \text{VRAI} \rangle$ ,  
sinon retourner  $\langle \text{FAUX} \rangle$ .*

2. *Si  $w \neq \varepsilon$  :*

*Pour  $1 \leq i \leq j \leq n$ , soit  $f(i, j)$  l'ensemble des variables dans  $V'$  qui engendrent le sous-mot  $w_i \dots w_j$ .*

*Calculer  $f(1, n)$  en appelant la sous-routine décrite à l'étape 3.*

*Si  $S_0 \in f(1, n)$ , alors retourner  $\langle \text{VRAI} \rangle$ , sinon retourner  $\langle \text{FAUX} \rangle$ .*

3. Calcul de  $f(i, j)$ , pour  $1 \leq i \leq j \leq n$  donnés.

Conserver les résultats obtenus au fur et à mesure. Si  $f(i, j)$  a déjà été calculé, alors retourner le résultat tout de suite, sans faire de calculs.

Si  $i = j$ , alors  $f(i, i) \leftarrow \{A \in V' \mid (A \rightarrow w_i) \in R'\}$ .

Si  $i < j$ , alors :

$f(i, j) \leftarrow \emptyset$ ,

pour  $k \leftarrow i$  à  $j - 1$  faire :

pour toute règle  $A \rightarrow BC$  faire :

si  $B \in f(i, k)$  et  $C \in f(k + 1, j)$ , alors :

$F \leftarrow F \cup \{A\}$

Pour une grammaire  $G$  fixée, et si la grammaire  $G'$  est calculée une fois pour toutes, le temps d'exécution de l'algorithme 4.94 sur un mot de longueur  $n$  est dominé par le nombre d'appels à la fonction  $f$  de l'étape 3.

La fonction récursive  $f$  est équivalente à un algorithme de programmation dynamique qui consiste à remplir une table de taille  $O(n^2)$  où chaque élément prend un temps  $O(n)$  à calculer. Le temps d'exécution total est dans  $O(n^3)$ . □

### Théorème 4.95

$$\text{REG} \subseteq \text{HC}$$

#### Démonstration :

Soit le langage  $L \in \text{REG}$  et soit  $M = (Q, \Sigma, \delta, q_0, F)$  un automate qui le reconnaît.

Alors on peut voir que la GHC  $G = (V, \Sigma, R, S)$  où

$$V = Q,$$

$$R = \{(A \rightarrow xB) \mid A \in Q, x \in \Sigma, B \in Q, \delta(A, x) = B\} \\ \cup \{(A \rightarrow \varepsilon) \mid A \in F\},$$

$$S = q_0,$$

engendre  $L$ .



### Corollaire 4.96

*La classe REG est strictement incluse dans la classe HC.*

#### Démonstration :

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

Nous avons déjà vu que  $L \in \text{HC}$  et  $L \notin \text{REG}$ .



## Théorème 4.97 (Lemme du pompiste pour les langages HC)

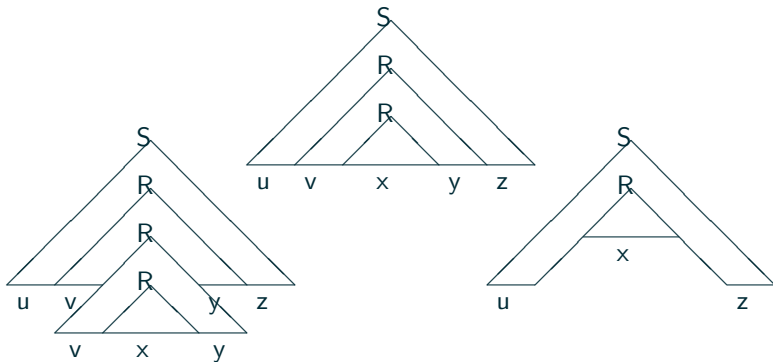
*Soit  $L \in \text{HC}$ . Il existe un entier  $p \geq 1$  tel que pour tout  $w \in L$  avec  $|w| \geq p$ , on peut écrire  $w = uvxyz$ , où*

- $|vy| > 0$ ,
- $|vxy| \leq p$ ,
- $\forall i \geq 0 : uv^i xy^i z \in L$ .



## Aperçu de la démonstration :

Si un mot  $w \in L$  est suffisamment long, alors il aura un arbre de dérivation très haut. Dans cet arbre, il y aura un chemin, du sommet jusqu'à la base, avec une répétition de variable :



### Démonstration du Théorème 4.97 :

Soit  $G$  une GHC en FNC telle que  $L(G) = L$ , et choisissons  $p = 2^{|V|}$ .

Considérons un arbre de dérivation pour  $w$ , et soit  $t_0, \dots, t_h$  le chemin le plus long de la racine à une feuille dans cet arbre. Pour tout sommet  $t$ , dénotons par  $s(t)$  la variable associée et par  $w(t)$  le mot engendré par le sous-arbre avec  $t$  comme racine.

Supposons d'abord que  $h \geq |V| + 1$ , et regardons les  $|V| + 2$  derniers sommets. Il y a au plus  $|V| + 1$  variables parmi eux et donc il y aura un  $i$  et un  $j$  tels que  $i < j$  et  $s(t_i) = s(t_j) = R \in V$  par le principe du pigeonnier.

Alors, définissons

$$x := w(t_j)$$

$$vxy := w(t_i)$$

$$uvxyz := w(t_0)$$

On a donc que

$$R \stackrel{*}{\Rightarrow} x$$

$$R \stackrel{*}{\Rightarrow} vRy$$

$$S \stackrel{*}{\Rightarrow} uRz$$

On peut maintenant prouver les trois propriétés promises par le lemme :

1. Clairement  $|vy| > 0$ , car  $G$  est en FNC et qu'il est impossible de dériver  $R \xRightarrow{*} R$  sauf trivialement en 0 étapes.
2.  $|vxy| \leq p$ , car on a choisi  $t_i$  et  $t_j$  parmi les  $|V| + 2$  derniers sommets. Le sous-arbre engendré par  $t_i$  a donc une hauteur d'au plus  $|V| + 1$  et donc  $|w(t_i)| = |vxy| \leq 2^{|V|} = p$ .
3. On peut dériver

$$S \xRightarrow{*} uRz \xRightarrow{*} uxz$$

ainsi que

$$S \xRightarrow{*} uRz \xRightarrow{*} uvRyz \xRightarrow{*} uv^2Ry^2z \xRightarrow{*} \dots \xRightarrow{*} uv^iRy^iz \xRightarrow{*} uv^ixy^iz$$

donc  $uv^ixy^iz \in L$  pour tout  $i \geq 0$ .

Comment s'assurer que  $h \geq |V| + 1$  ?

À chaque niveau, on double, sauf au dernier niveau où on a des transitions de type  $A \Rightarrow a$ . Donc,  $|w| \leq 2^{h-1}$ .

Donc, puisqu'on a choisi  $p = 2^{|V|}$ , on a que

$$p = 2^{|V|} \leq |w| \leq 2^{h-1},$$

et donc

$$h \geq |V| + 1.$$

Ceci conclut donc la démonstration.



### Exemple 4.98

Soit  $L = \{a^n b^n \mid n \geq 0\}$ . Nous avons vu que ce langage est HC, et une grammaire en FNC à 5 variables permet de l'engendrer.

Donc, le lemme du pompiste HC devrait fonctionner pour  $p = 2^5 = 32$ . Et effectivement, pour tout  $w = a^k b^k$  pour  $k \geq 16$ , on peut écrire

$$w = uvxyz$$

$$x = ab$$

$$v = a$$

$$y = b$$

$$u = a^{k-2}$$

$$z = b^{k-2}$$

et toutes les propriétés du lemme sont satisfaites.

### Théorème 4.99

*Le langage*

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

*n'est pas hors contexte.*

#### Démonstration :

Supposons au contraire que  $L \in \text{HC}$ , afin d'arriver à une contradiction.

Soit  $p \geq 1$  tel que donné par le lemme du pompiste pour les langages hors contexte.

Considérons le mot  $w = a^p b^p c^p$ .

On a bien  $w \in L$  et  $|w| = 3p \geq p$ .

Soient  $w = uvxyz$ ,  $|vy| > 0$  et  $|vxy| \leq p$ , tels que donnés par le lemme.

Considérons les deux cas suivants :

1. La partie  $v$ , ou  $y$ , contient plus d'une sorte de symboles parmi  $a$ ,  $b$  et  $c$ .

Choisissons  $i = 2$  dans l'expression  $uv^i xy^i z$ .

Ce mot n'est pas de la forme  $a^n b^n c^n$ , et donc  $uv^i xy^i z \notin L$ , en contradiction avec l'énoncé du lemme.



2. La partie  $v$  ne contient qu'une seule sorte de symbole, de même que  $y$ .

Choisissons  $i = 0$  dans l'expression  $uv^i xy^i z$ , c'est-à-dire  $uxz$ .

Ce mot ne contient pas un nombre égal des trois symboles et n'est pas de la forme  $a^n b^n c^n$ , et donc  $uv^i xy^i z \notin L$ , de nouveau en contradiction avec l'énoncé.

Donc  $L \notin \text{HC}$ .



## Théorème 4.100

*La classe HC est fermée pour l'union.*

### Démonstration :

Soient

$$L_1 = L(G_1) \text{ où } G_1 = (V_1, \Sigma, R_1, S_1)$$

et

$$L_2 = L(G_2) \text{ où } G_2 = (V_2, \Sigma, R_2, S_2),$$

en supposant que  $V_1$  et  $V_2$  soient disjoints, tout comme  $R_1$  et  $R_2$ .

Alors la GHC suivante :

$$G = (V_1 \cup V_2 \cup \{S\}, \Sigma, R_1 \cup R_2 \cup \{S \rightarrow S_1 \mid S_2\}, S)$$

engendre le langage  $L_1 \cup L_2$ .



### **Théorème 4.101**

*La classe HC n'est pas fermée pour l'intersection.*

#### **Démonstration :**

Soient

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

$$L_2 = \{a^n b^n c^m \mid m \geq 0, n \geq 0\}.$$

Alors

$$L_1 \cap L_2 = \{a^n b^n c^n \mid n \geq 0\}.$$

Or, on a vu au théorème 4.99 que ce langage n'est pas hors contexte.  $\square$

### **Théorème 4.102**

*La classe des langages hors contexte n'est pas fermée pour la complémentation.*

#### **Démonstration :**

Supposons au contraire que la classe est fermée afin d'arriver à une contradiction.

Soient  $L_1 \in \text{HC}$  et  $L_2 \in \text{HC}$ .

En considérant

$$L_1 \cap L_2 = \overline{\overline{L_1} \cup \overline{L_2}}$$

et que la classe HC est fermée pour l'union (théorème 4.100), on obtient une contradiction avec le fait que HC n'est pas fermée pour l'intersection (théorème 4.101). □

### **Théorème 4.103**

*La classe HC est fermée sous la concaténation*

**Démonstration :**

Soient  $L_1$  et  $L_2$  deux langages hors-contexte ; nous devons prouver que  $L_1 \circ L_2 \in \text{HC}$ . Soient  $G_1 = (V_1, \Sigma, R_1, S_1)$  et  $G_2 = (V_2, \Sigma, R_2, S_2)$  des grammaires engendrant  $L_1$  et  $L_2$  respectivement. Définissons  $G = (V, \Sigma, R, S)$  pour  $L_1 \circ L_2$  comme suit :

$$V = V_1 \cup V_2 \cup \{S\}$$

$$R = R_1 \cup R_2 \cup \{S \rightarrow S_1 S_2\}.$$



### **Théorème 4.104**

*La classe HC est fermée sous l'étoile.*

**Démonstration :**

Soit  $L \in \text{HC}$ , et soit  $G = (V, \Sigma, R, S)$  une grammaire en FNC engendrant  $L$ . Alors, définissons la grammaire  $G' = (V', \Sigma, R', S')$  engendrant  $L^*$  comme suit :

$$V' = V$$

$$S' = S$$

$$R' = R \cup \{S \rightarrow \varepsilon, S \rightarrow SS\}$$

Notons que vu que  $G$  est en FNC,  $S$  ne peut être utilisé qu'au début de la dérivation, et les deux nouvelles règles ajoutées ne peuvent pas être utilisées « à l'intérieur » d'une dérivation dans  $G$ . □



### **Théorème 4.105**

*Si  $L_1 \in \text{HC}$  et  $L_2 \in \text{REG}$ , alors  $L_1 \cap L_2 \in \text{HC}$ .*

### Démonstration :

Soit  $G = (V, \Sigma, R, S)$  une grammaire hors-contexte en FNC pour  $L_1$  et  $M = (Q, \Sigma, \delta, q_0, F)$  un AFD reconnaissant  $L_2$ .

Nous allons construire une grammaire  $G' = (V', \Sigma, R', S')$  engendrant  $L_1 \cap L_2$ .

Soient :

$$V' = (V \times Q \times Q) \cup \{\langle \text{INIT} \rangle\}$$

$$S' = \langle \text{INIT} \rangle$$

Le but sera de créer  $R'$  de façon à ce que la variable  $(A, q_i, q_j)$  puisse engendrer les mots pouvant être engendrés par  $G$  à partir de  $A$ , et reconnus par  $M$  en passant de l'état  $q_i$  à  $q_j$ .

Construisons  $R'$  de la façon suivante :

- Pour tout  $q \in F$  : on ajoute la règle

$$\langle \text{INIT} \rangle \rightarrow (S, q_0, q)$$

- Si la règle  $S \rightarrow \varepsilon$  existe dans  $G$  et que  $q_0 \in F$ , alors on ajoute la règle  $\langle \text{INIT} \rangle \rightarrow \varepsilon$ .
- Pour toute règle de forme  $A \rightarrow a$  et tout  $q \in Q$ , on ajoute la règle

$$(A, q, \delta(q, a)) \rightarrow a$$

- Pour toute règle de forme  $A \rightarrow BC$  et pour tout  $q_1, q_2, q_3 \in Q$ , on ajoute :

$$(A, q_1, q_3) \rightarrow (B, q_1, q_2)(C, q_2, q_3)$$



# Les langues naturelles sont-elles hors-contexte ?

L'idée originale de Chomsky était d'utiliser les grammaires hors-contexte pour décrire les langues naturelles.

Dans plusieurs langues, il existe des phrases sur lesquelles on peut utiliser le lemme du pompiste pour infirmer cette théorie.

Considérons la phrase :

Pierre, Jean, Pierrette, Jeannette, [...] sont respectivement chauve, roux,  
myope, presbyte, [...]

On retrouve la structure :

$\langle \text{NM} \rangle \langle \text{NM} \rangle \langle \text{NF} \rangle \langle \text{NF} \rangle \dots$  sont respectivement  $\langle \text{AM} \rangle \langle \text{AM} \rangle \langle \text{AF} \rangle \langle \text{AF} \rangle \dots$

La structure  $\{\langle \text{NM} \rangle, \langle \text{NF} \rangle\}^*$  sont respectivement  $\{\langle \text{AM} \rangle, \langle \text{AF} \rangle\}^*$   
ne peut être valide que si le nombre de noms masculins (resp. féminins)  
est égal au nombre d'adjectifs masculins (resp. féminins).

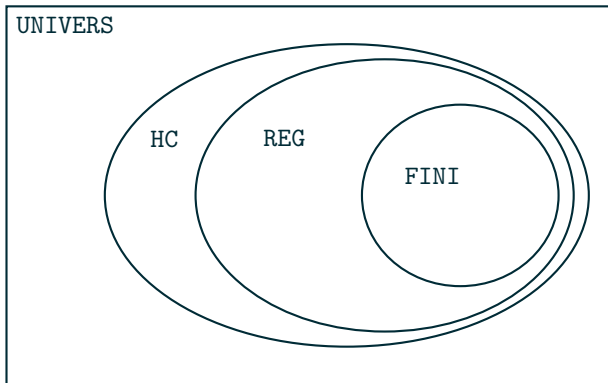
Définissons donc le langage régulier suivant :

$$L = \{w \mid w = \langle \text{NM} \rangle^n \langle \text{NF} \rangle^m \text{ sont respectivement } \langle \text{AM} \rangle^\ell \langle \text{AF} \rangle^k, \\ n, m, \ell, k \geq 0\}$$

Si le français est hors-contexte, l'intersection du français avec  $L$  est également hors-contexte selon le théorème 4.105. Or, cette intersection est égale à :

$$L' = \{w \mid w = \langle \text{NM} \rangle^n \langle \text{NF} \rangle^m \text{ sont respectivement } \langle \text{AM} \rangle^n \langle \text{AF} \rangle^m, \\ n, m \geq 0\}$$

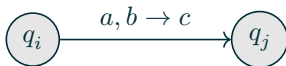
et on peut facilement démontrer avec le lemme du pompiste que ce langage n'est pas hors-contexte.



# Automates à pile

Il existe aussi un type d'automate qui reconnaît les langages hors-contexte : les automates à pile (AP).

Un AP est essentiellement un AFN avec une pile pouvant contenir des lettres d'un alphabet  $\Gamma$  potentiellement différent de celui du langage. Chaque transition permet de dépiler et d'empiler une lettre. On a donc des transitions comme :



Dans cette transition, on lit le caractère  $a \in \Sigma_\epsilon$ , on dépile  $b \in \Gamma_\epsilon$ , et on empile  $c \in \Gamma_\epsilon$ . Ici,  $\Gamma_\epsilon = \Gamma \cup \{\epsilon\}$ , voulant dire que ces symboles peuvent être le caractère vide  $\epsilon$ .

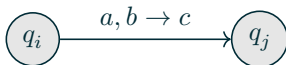


### Définition 4.106 (Automate à pile)

Un automate à pile est un tuple  $(Q, \Sigma, \Gamma, \delta, q_0, F)$ , où  $Q, \Sigma, \Gamma$  et  $F$  sont des ensembles finis, et où

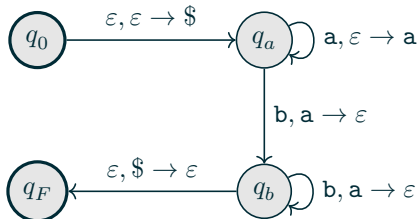
1.  $Q$  est l'ensemble des états
2.  $\Sigma$  est l'alphabet du langage
3.  $\Gamma$  est l'alphabet de la pile
4.  $\delta : Q \times Q \rightarrow \mathcal{P}(\Sigma_\varepsilon \times \Gamma_\varepsilon \times \Gamma_\varepsilon)$  est la fonction de transition
5.  $q_0 \in Q$  est l'état de départ
6.  $F \subseteq Q$  est l'ensemble des états acceptants

$\delta(q_i, q_j)$  contient l'ensemble des transitions permises pour passer de  $q_i$  à  $q_j$ . Si  $(a, b, c) \in \delta(q_i, q_j)$  alors la transition



est permise. Un mot est accepté par l'automate s'il existe une suite de transitions sur ce mot permettant d'arriver à un état acceptant.

L'automate suivant avec  $\Sigma = \{a, b\}$  et  $\Gamma = \{\$, a\}$  reconnaît  
 $L = \{a^n b^n \mid n \geq 0\}$  :



### **Théorème 4.107**

*Un langage  $L$  est hors-contexte si et seulement si il est reconnu par un automate à pile.*

Les langages contextuels sont engendrés par une grammaire dont les règles permettent un *contexte* conditionnant leur application, par exemple

$$\alpha A \beta \rightarrow \alpha \gamma \beta$$

permet la transition  $A \Rightarrow \gamma$  seulement si  $A$  est précédé de  $\alpha$  et suivi de  $\beta$ .

# Hiérarchie de Chomsky

Niveau	Classe	Règles	Exemples
Type-3	REG	$A \rightarrow a, A \rightarrow aB$	$\{a^n \mid n \geq 0\}$
Type-2	HC	$A \rightarrow \gamma$	$\{a^n b^n \mid n \geq 0\}$
Type-1	Contextuels	$\alpha A \beta \rightarrow \alpha \gamma \beta$	$\{a^n b^n c^n \mid n \geq 0\}$
Type-0	Récursivement énumérables	$\alpha \rightarrow \gamma$	Peut simuler les MT