

# Automates, langages et compilation

## Grammaires

---

Isabelle Ryl

2024 – 2025

Cours de L3 - Université Paris Dauphine-PSL

# 1. Grammaires

## 1.1 Présentation

## 1.2 Grammaires régulières

## 1.3 Grammaires context-free

## 1.4 Transformation de Grammaires context-free

## 1.5 Lemme d'itération

# Grammaires

---

# Grammaires

---

## Présentation

# Exemple

Revenons sur l'exemple du premier cours

- les expressions régulières. Ex :

$$(b + c)^* a(b + c)^* a(a + b + c)^*$$

- les grammaires. Ex :

$$S \longrightarrow TaTaT$$

$$T \longrightarrow Ta \mid Tb \mid Tc \mid \varepsilon$$

Exemple de dérivation :

$$\begin{aligned} S &\longrightarrow TaTaT \\ &\longrightarrow \textcolor{red}{T}aaTaT \\ &\longrightarrow Taa\textcolor{red}{T}baT \\ &\longrightarrow TaaTba\textcolor{red}{T}c \\ &\longrightarrow aaTbaTc \\ &\longrightarrow aaTbac \\ &\longrightarrow aabac \end{aligned}$$

## Similitude avec le langage naturel (1/2)

$\Sigma = \{le, la, chien, balle, maitre, lance, attrappe\}$

Nous pouvons imaginer décomposer une phrase  $P$ , en sujet-verbe-complément. Nous pourrions alors définir les catégories syntaxiques de la manière suivante :

$P \longrightarrow \text{SUJET VERBE COMPLEMENT}$

$\text{VERBE} \longrightarrow lance \mid attrappe$

Le sujet comme le complément sont des groupes nominaux composés d'un article et d'un nom :

$\text{SUJET} \longrightarrow \text{GN}$

$\text{COMPLEMENT} \longrightarrow \text{GN}$

$\text{GN} \longrightarrow \text{ARTICLE NOM}$

$\text{ARTICLE} \longrightarrow le \mid la$

$\text{NOM} \longrightarrow chien \mid balle \mid maitre$

## Similitude avec le langage naturel (2/2)

Exemple de dérivation :

*P* → *SUJET VERBE COMPLEMENT*  
→ *SUJET lance COMPLEMENT*  
→ *GN lance COMPLEMENT*  
→ *GN lance GN*  
→ *GN lance ARTICLE NOM*  
→ *GN lance la NOM*  
→ *ARTICLE NOM lance la NOM*  
→ *ARTICLE maitre lance la NOM*  
→ *le maitre lance la NOM*  
→ *le maitre lance la balle*

## Similitude avec le langage naturel (2/3)

Mais ... :

- P* → *SUJET VERBE COMPLEMENT*
- *SUJET lance COMPLEMENT*
- *GN lance COMPLEMENT*
- *GN lance GN*
- *GN lance ARTICLE NOM*
- *GN lance la NOM*
- *ARTICLE NOM lance la NOM*
- *ARTICLE chien lance la NOM*
- *la chien lance la NOM*
- *la chien lance la maitre*

➡ seule la syntaxe est vérifiée



## Grammaire - Définition

Une grammaire est donc un mécanisme de génération des langages, nous verrons qu'il est plus général que les expressions rationnelles et les automates finis

Les grammaires ont été introduites dans plusieurs domaines en particulier pour formaliser le langage naturel et sa grammaire par Noam Chomsky, mais sont très utilisées en informatique notamment en compilation

Pour compiler un langage, il faut vérifier que les termes utilisés sont corrects par l'analyse lexicale qui utilisera les langages reconnaissables puis vérifier que la structure du programme est correcte, *i.e.* que la syntaxe est correct par l'analyse syntaxique, en vérifiant que la grammaire du langage de programmation est respectée

# Grammaire de langage de programmation - Exemple

Exemple tiré du manuel de référence Python\*

```
if_stmt      ::=  "if" assignment_expression ":" suite  
                ("elif" assignment_expression ":" suite)*  
                ["else" " ":" suite]  
  
while_stmt   ::=  "while" assignment_expression ":" suite  
                ["else" " ":" suite]
```

(\*) <https://docs.python.org/fr/3/reference/index.html>

Une grammaire est un quadruplet  $\mathcal{G} = (\Sigma, V, S, R)$  tel que :

- $\Sigma$  est l'alphabet sur lequel le langage est défini, appelé alphabet terminal
- $V$  est l'alphabet des symboles utilisés au cours de la génération, appelé alphabet non terminal, ou les « variables »  $\Sigma \cap V = \emptyset$
- $S \in V$  est le symbole de départ, appelé également axiome
- $R$  est un ensemble de règles de réécriture de la forme  $\alpha \rightarrow \beta$  avec  $\alpha \in (\Sigma \cup V)^+$  et  $\beta \in (\Sigma \cup V)^*$

Idée. Une règle de réécriture permet de réécrire un mot  $u$  en un mot  $v$  en appliquant une règle  $f_1 \rightarrow f_2$  si on trouve le facteur  $f_1$  dans le mot  $u$ , dans ce cas on « remplace » l'instance de  $f_1$  trouvée dans  $u$  par  $f_2$  pour former  $v$

Une **règle de réécriture**  $r : f_1 \rightarrow f_2$  s'applique pour réécrire un mot  $u$  en un mot  $v$ , si et seulement si  $\exists u_1, u_2$  tels que  $u = u_1 f_1 u_2$  et  $v = u_1 f_2 u_2$ , on note  $u \xrightarrow{r} v$

Soit  $\mathcal{R}$  un système de règles de réécriture, une **dérivation** d'un mot  $u$ , notée  $u \xrightarrow[\ast]{\mathcal{R}} v$ , est la clôture par application des règles de réécriture :

$$\exists u = u_1, u_2, \dots, u_{n-1}, u_n = v \text{ avec } \forall 1 \leq i < n, \exists r_i \in \mathcal{R} \text{ telle que } u_i \xrightarrow{r_i} u_{i+1}$$

# Langage engendré

Soient  $\mathcal{G} = (\Sigma, V, S, R)$  une grammaire,  $\alpha \in (\Sigma \cup V)^+$  et  $\beta \in (\Sigma \cup V)^*$

La grammaire  $\mathcal{G}$  permet de **dériver**  $u$  en  $v$ , si  $u \xrightarrow[\ast]{R} v$ , et on note  $u \xrightarrow[\ast]{\mathcal{G}} v$

Le **langage engendré par une grammaire** (ou généré)  $\mathcal{G}$  est  $\mathcal{L}(\mathcal{G}) \in \Sigma^*$  est l'ensemble des mots de l'alphabet terminal qui peuvent être dérivés à partir de  $S$ , soit

$$\mathcal{L}(\mathcal{G}) = \{u \in \Sigma^* \mid S \xrightarrow[\ast]{\mathcal{G}} u\}$$

Un langage  $\mathcal{L} \in \Sigma^*$  est dit **recursivement énumérable** s'il existe une grammaire  $\mathcal{G}$  telle que  $\mathcal{L}(\mathcal{G}) = \mathcal{L}$

# Hiérarchie de Chomsky 1/3

Classification introduite par Noam Chomsky, basée sur la forme des règles des grammaires, qui permet d'obtenir des classes ordonnées par inclusion des langages engendrés

Type 0. Pas de contrainte sur la forme des règles

Type 1. Grammaires contextuelles (ou sensibles au contexte) : les règles sont de la forme

$$\alpha T \beta \longrightarrow \alpha \gamma \beta$$

où  $T \in V$ ,  $\alpha, \beta \in (\Sigma \cup V)^*$  et  $\gamma \in (\Sigma \cup V)^+$

- une seule variable est réécrite à chaque règle (ici  $T$ )
- les « contextes » gauche et droit ne sont pas modifiés

## Hiérarchie de Chomsky 2/3

Type 2. Grammaires hors contexte (ou context-free ou algébrique) : les règles sont de la forme

$$T \longrightarrow \alpha$$

où  $T \in V$ ,  $\alpha \in (\Sigma \cup V)^*$

- une seule variable est réécrite à chaque règle (ici  $T$ )
- chaque variable peut être réécrite indépendamment du contexte

Type 3. Grammaires régulières, la forme des règles est contrainte selon deux cas :

- régulière à droite  $T \longrightarrow aU$  ou  $T \longrightarrow a$
- régulière à gauche  $T \longrightarrow Ua$  ou  $T \longrightarrow a$

où  $T, U \in V$ ,  $a \in \Sigma$

Les différents types de grammaire sont associés à différents types de langages, qui sont reconnus par différents types de « machines »

- Type 3, grammaires régulières, langages dits réguliers que nous avons vus sous le nom de rationnels ou reconnaissables, automates finis
- Type 2, grammaires hors contexte, langages hors contexte, automates à pile
- Type 1, grammaires contextuelles, langages contextuels, machines de Turing
- Type 0, langages décidables (reconnus par une machine en un temps fini), machines de Turing



## Remarques

- Une grammaire engendre un et un seul langage
- Un langage peut être engendré par plusieurs grammaires différentes

Deux grammaires sont **équivalentes** si elles engendrent le même langage

Attention, deux grammaires distinctes pouvant être équivalentes, un langage de type  $n$  peut être engendré par une grammaire de type  $m \leq n$ . Pour montrer qu'un langage n'est pas de type  $n$  il faut montrer que toute grammaire qui l'engendre est de type  $< n$

# Exemples

Construire une grammaire qui engendre le langage  $a^n b^n$ .

$$S \longrightarrow aSb \mid \varepsilon$$

Exemples de dérivation :

$$S \longrightarrow aSb$$

$$\longrightarrow ab$$

$$S \longrightarrow aSb$$

$$\longrightarrow aaSbb$$

$$\longrightarrow aabb$$

$$S \longrightarrow aSb$$

$$\longrightarrow aaSbb$$

$$\longrightarrow aaaSbbb$$

$$\longrightarrow aaabbb$$

➡  $a^n b^n$  peut être engendré par une grammaire

# Exemples

Quel est le langage engendré par la grammaire :

$$S \longrightarrow aSa \mid bSb \mid cSc \mid a \mid b \mid c \mid \varepsilon$$

Exemple de dérivation :

$$\begin{aligned} S &\longrightarrow aSa \\ &\longrightarrow abSba \\ &\longrightarrow abbSbba \\ &\longrightarrow abbcScbba \\ &\longrightarrow abbccbba \end{aligned}$$

➡ Les palindromes sur  $\{a, b, c\}$

# Montrer qu'un langage est généré par une grammaire

Pour montrer qu'une grammaire  $\mathcal{G}$  engendre un langage  $\mathcal{L}$ , il faut montrer que :

- $\mathcal{L}(\mathcal{G}) \subseteq \mathcal{L}$  ➡ tout mot engendré par  $\mathcal{G}$  est un mot de  $\mathcal{L}$
- $\mathcal{L} \subseteq \mathcal{L}(\mathcal{G})$  ➡ tout mot de  $\mathcal{L}$  peut être engendré par  $\mathcal{G}$

## Exemple

Reprenons l'exemple ci-dessus, considérons la grammaire  $\mathcal{G} = (\Sigma, V, S, R)$  avec  $\Sigma = \{a, b, c\}$ ,  $V = \{S\}$  et  $R = \{S \rightarrow aSa \mid bSb \mid cSc \mid a \mid b \mid c \mid \varepsilon\}$  et  $\mathcal{L}$  le langage des palindromes sur  $\Sigma$  et montrons que  $\mathcal{G}$  génère  $\mathcal{L}$

## Exemple Palindromes - Preuve (1/2)

Montrons  $\mathcal{L} \subseteq \mathcal{L}(\mathcal{G})$  par récurrence sur la longueur des mots de  $\mathcal{L}$

Soit  $u \in \mathcal{L}$

- si  $|u| \leq 1$  alors  $u \in \{\varepsilon, a, b, c\}$  et  $u \in \mathcal{L}(\mathcal{G})$ , supposons donc que l'inclusion est vraie pour tout  $u$  tel que  $|u| \leq n$
- si  $|u| > 1$  alors  $u = ava$  ou  $u = bvb$  ou  $u = cvc$  avec  $v$  un palindrome et  $|v| \leq n$ . Par hypothèse de récurrence  $v \in \mathcal{L}(\mathcal{G})$  et donc  $S \xrightarrow[*]{\mathcal{G}} v$ , nous pouvons donc dériver :

$$\begin{aligned} S &\rightarrow aSa \xrightarrow[*]{\mathcal{G}} ava \\ \text{et } S &\rightarrow bSb \xrightarrow[*]{\mathcal{G}} bvb \\ \text{et } S &\rightarrow cSc \xrightarrow[*]{\mathcal{G}} cvc \end{aligned}$$

et donc  $u \in \mathcal{L}(\mathcal{G})$

## Exemple Palindromes - Preuve (2/2)

Montrons  $\mathcal{L}(\mathcal{G}) \subseteq \mathcal{L}$  par récurrence sur la longueur des mots de  $\mathcal{L}(\mathcal{G})$

Soit  $u \in \mathcal{L}(\mathcal{G})$

- si  $|u| \leq 1$  alors  $u \in \{\varepsilon, a, b, c\}$  et  $u \in \mathcal{L}$ , supposons donc que l'inclusion est vraie pour tout  $u$  tel que  $|u| \leq n$
- si  $|u| > 1$  alors la dérivation qui engendre  $u$  commence par l'application d'une règle  $S \rightarrow aSa$  ou  $S \rightarrow bSb$  ou  $S \rightarrow cSc$ , les cas étant strictement identiques, supposons qu'il s'agit de la première. Alors la dérivation s'écrit  $S \rightarrow aSa \xrightarrow[*]{\mathcal{G}} u = ava$  avec  $S \xrightarrow[*]{\mathcal{G}} v$  et  $|v| < n$ , donc par hypothèse de récurrence  $v \in \mathcal{L}$ , donc  $v$  est un palindrome et donc  $ava = u$  également :  $u \in \mathcal{L}$

# Notation BNF

La syntaxe de nombreux langages de programmation est décrite à l'aide de la notation de Backus-Naur (ou Backus-Naur form) :

- la dérivation  $\longrightarrow$  est notée  $::=$
- les non terminaux sont notés entre  $<$  et  $>$
- les alternatives en partie droite sont séparées par des  $|$
- les caractères spéciaux sont entourés de  $'$

Exemple. Grammaire qui engendre les sommes de nombres

$S \longrightarrow S + S \mid NB$

$NB \longrightarrow NB CH \mid CH$

$CH \longrightarrow 0 \mid 1 \mid \dots \mid 9$

$\langle \text{sum} \rangle ::= \langle \text{sum} \rangle ' + ' \langle \text{sum} \rangle \mid \langle \text{num} \rangle$

$\langle \text{num} \rangle ::= \langle \text{num} \rangle \langle \text{dig} \rangle \mid \langle \text{dig} \rangle$

$\langle \text{dig} \rangle ::= 0 \mid 1 \mid \dots \mid 9$

# Grammaires

---

## Grammaires régulières



Un langage est dit **régulier** si et seulement si il existe une grammaire régulière qui l'engendre. On note  $Reg(\Sigma^*)$  la classe des langages réguliers sur  $\Sigma$

## Théorème

$$Reg(\Sigma^*) = Rat(\Sigma^*) = Rec(\Sigma^*)$$

## Exemple

Pour le langage  $(b + c)^* a(b + c)^* a(a + b + c)^*$  nous avons donné la grammaire

$$S \longrightarrow TaTaT$$

$$T \longrightarrow Ta \mid Tb \mid Tc \mid \varepsilon$$

Cette grammaire est context-free, or le langage est rationnel, nous pouvons proposer une autre grammaire

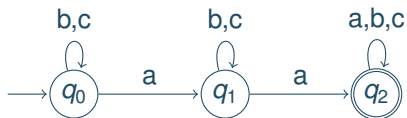
$$S \longrightarrow bS \mid cS \mid aT$$

$$T \longrightarrow bT \mid cT \mid aU$$

$$U \longrightarrow aU \mid bU \mid cU \mid \varepsilon$$

# Idée de la construction

$$\begin{aligned} S &\longrightarrow bS \mid cS \mid aT \\ T &\longrightarrow bT \mid cT \mid aU \\ U &\longrightarrow aU \mid bU \mid cU \mid \varepsilon \end{aligned}$$



Les variables et les états jouent des rôles très similaires en gardant trace de la progression dans la construction/reconnaissance du motif.

Pour le mot *cabbac* par exemple

$$\begin{aligned} q_0 &\xrightarrow{c} q_0 \xrightarrow{a} q_1 \xrightarrow{b} q_1 \xrightarrow{b} q_1 \xrightarrow{a} q_2 \xrightarrow{c} q_2 \\ S &\xrightarrow{S} cS \xrightarrow{S} caT \xrightarrow{T} cabT \xrightarrow{T} cabbT \xrightarrow{T} cabbaU \xrightarrow{U} cabbacU \xrightarrow{U} cabbac \end{aligned}$$

## Preuve $Reg \subseteq Rec$ (1/2)

Soit  $\mathcal{G} = (\Sigma, V, S, R)$  une grammaire régulière à droite. Posons  $\mathcal{A} = (\Sigma, V \cup \{F\}, S, \{F\}, \delta)$  un automate fini tel que :

$$\delta(T, a) = U \Leftrightarrow (T \rightarrow aU) \in R$$

$$\delta(T, a) = F \Leftrightarrow (T \rightarrow a) \in R$$

et montrons que  $\mathcal{L}(\mathcal{G}) = \mathcal{L}(\mathcal{A})$  par récurrence sur la longueur des mots produits par la grammaire :

$$(S \xrightarrow[*]{\mathcal{G}} uT) \Leftrightarrow (S \xrightarrow[*]{u} T \text{ calcul de } \mathcal{A})$$

Par construction de  $\delta$  l'équivalence est vraie pour  $|u| = 1$ . Supposons l'équivalence vraie pour tout  $u$  tel que  $|u| < n$  et considérons  $u$  tel que  $|u| = n$  alors  $u = va$  avec  $a \in \Sigma$  et  $|v| < n$ .

Dans le sens  $\Rightarrow$  : soit la dérivation  $S \xrightarrow[*]{\mathcal{G}} vaT$  comme  $\mathcal{G}$  est régulière à droite, il existe une dérivation  $S \xrightarrow[*]{\mathcal{G}} vU \xrightarrow{\mathcal{G}} vaT$ , par HR,  $S \xrightarrow[*]{v} U$  est donc un calcul de  $\mathcal{A}$  par construction de  $\delta$  comme  $U \rightarrow aT \in R$ ,  $\delta(U, a) = T$  et donc  $S \xrightarrow[*]{v} U \xrightarrow{a} T$  est un calcul de  $\mathcal{A}$

## Preuve $Reg \subseteq Rec$ (2/2)

Dans le sens  $\Leftarrow$  : il existe un calcul de  $\mathcal{A}$  tel que  $S \xrightarrow[*]{u} T$ , nous pouvons exhiber la dernière transition  $S \xrightarrow[*]{v} U \xrightarrow{a} T$  par HR  $S \xrightarrow[*]{\mathcal{G}} vU$  est donc une dérivation. Par construction de  $\delta$  si  $(\delta(U, a) = T) \Leftrightarrow (U \rightarrow aT) \in R$  donc la dérivation se poursuit  $S \xrightarrow[*]{\mathcal{G}} vU \xrightarrow{\mathcal{G}} vaT$  soit  $S \xrightarrow[*]{\mathcal{G}} uT$

Montrons à présent que  $u \in \mathcal{L}(\mathcal{G}) \Leftrightarrow u \in \mathcal{L}(\mathcal{A})$  avec  $u = va$ . Nous avons montré qu'il existe  $T$  tel que

$$S \xrightarrow[*]{v} T \text{ et } S \xrightarrow[*]{\mathcal{G}} vT$$

alors si  $u \in \mathcal{L}(\mathcal{G})$ ,  $vT \xrightarrow[*]{\mathcal{G}} va$  et par construction de  $\delta$ ,  $\delta(T, a) = F$  et réciproquement

Soit  $\mathcal{A} = (\Sigma, Q, q_0, F, \delta)$  un automate fini déterministe. Posons  $\mathcal{G} = (\Sigma, Q, q_0, R)$  telle que :

$$\begin{aligned} R &= \{T \rightarrow aU \mid T, U \in Q, a \in \Sigma, \delta(T, a) = U\} \\ &\cup \{T \rightarrow a \mid T \in Q, a \in \Sigma, \delta(T, a) \in F\} \end{aligned}$$

Exercice : montrer que  $\mathcal{L}(\mathcal{G}) = \mathcal{L}(\mathcal{A})$ .

# Grammaires

---

## Grammaires context-free

## Définition - Rappel

Grammaires hors contexte (ou context-free ou algébrique) : les règles sont de la forme

$$T \longrightarrow \alpha$$

où  $T \in V$ ,  $\alpha \in (\Sigma \cup V)^*$

- une seule variable est réécrite à chaque règle (ici  $T$ )
- chaque variable peut être réécrite indépendamment du contexte

*i.e.* À tout moment d'une dérivation un non-terminal peut être réécrit indépendamment du contexte dans lequel il est placé

Un langage est dit **hors contexte** ou context-free ou algébrique si et seulement si il existe une grammaire hors-contexte qui l'engendre



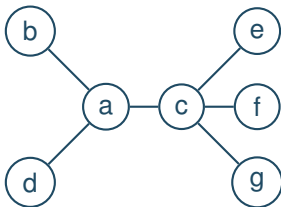
# Arbre de dérivation

Les grammaires context-free permettent de dériver un mot de nombreuses manières, nous allons donc utiliser une représentation graphique qui prend en compte le fait que l'ordre des dérivations ne change pas le résultat

Un arbre est un **arbre de dérivation** d'une grammaire context-free  $\mathcal{G} = (\Sigma, V, S, R)$  si et seulement si :

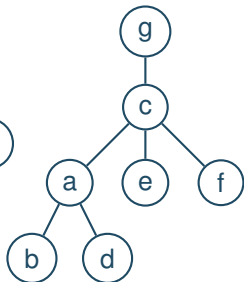
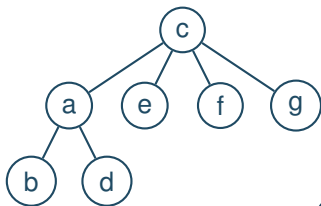
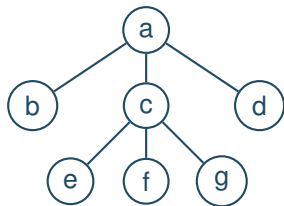
## Rappel - Arbre

- Un arbre est un graphe non orienté acyclique et connexe
- Deux sommets sont connectés par un chemin unique
- Un arbre a une arête de moins que son nombre de sommets
- Un arbre est connexe, mais ne l'est plus si on enlève une seule arête, il est acyclique mais ne l'est plus si on ajoute une seule arête



## Rappel - Arbre

Un sommet est distingué et appelé **racine**



- Il existe un chemin unique entre la racine  $R$  et un nœud  $N$  quelconque
- Si on oriente le chemin de  $R$  vers  $N$ , le prédécesseur de  $N$  est appelé père de  $N$ , ses successeurs, ses fils
- L'arité d'un nœud est son nombre de fils
- Les nœuds d'arité 0 sont appelés feuilles et les nœuds d'arité  $> 0$  et différents de la racine, nœuds internes

# Arbre de dérivation

Les grammaires context-free permettent de dériver un mot de nombreuses manières, nous allons donc utiliser une représentation graphique qui prend en compte le fait que l'ordre des dérivations ne change pas le résultat

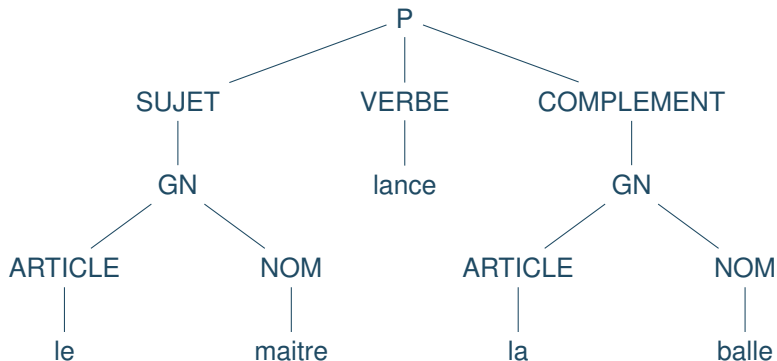
Un arbre est un **arbre de dérivation** d'une grammaire context-free  $\mathcal{G} = (\Sigma, V, S, R)$  si et seulement si :

- l'étiquette de la racine est  $S$
- chaque nœud interne est étiqueté par un symbole de  $V$
- chaque feuille est étiquetée par un symbole de  $\Sigma$
- chaque nœud interne étiqueté par un symbole  $T$  a pour fils des symboles  $\alpha_1, \dots, \alpha_n$  (l'ordre est important) si et seulement si  $T \rightarrow \alpha_1 \dots \alpha_n \in R$

➡ La concaténation des feuilles de l'arbre de gauche à droite constitue le mot produit

# Arbre de dérivation - Exemple Phrase

Mot : le maitre lance la balle



## Arbre de dérivation - Exemple $a^n b^n$

$$S \longrightarrow aSb \mid \varepsilon$$

$$S \longrightarrow aSb$$

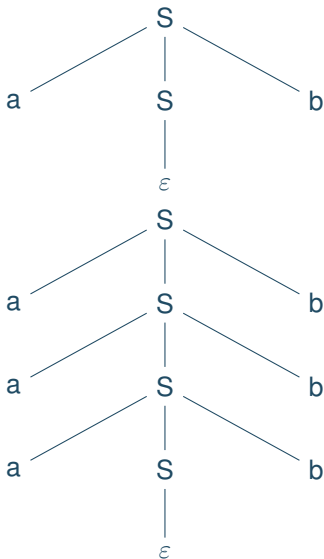
$$\longrightarrow ab$$

$$S \longrightarrow aSb$$

$$\longrightarrow aaSbb$$

$$\longrightarrow aaaSbbb$$

$$\longrightarrow aaabbbb$$

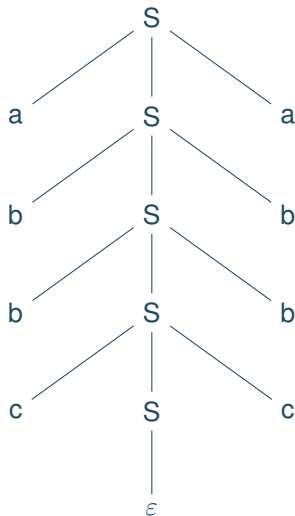


# Arbre de dérivation - Exemple palindromes

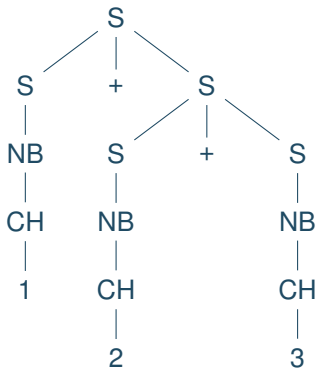
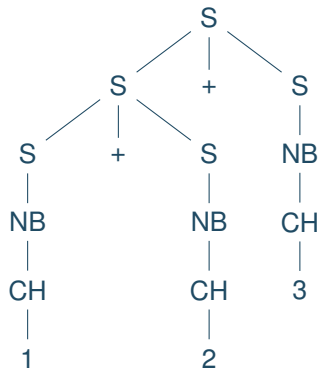
$$S \longrightarrow aSa \mid bSb \mid cSc \mid a \mid b \mid c \mid \varepsilon$$

Dérivation

$S \longrightarrow aSa$   
 $\longrightarrow abSba$   
 $\longrightarrow abbSbba$   
 $\longrightarrow abbScbba$   
 $\longrightarrow abbccbba$



# Ambiguïté - Exemple

$$S \rightarrow S + S \mid NB$$
$$NB \rightarrow NB\ CH \mid CH$$
$$CH \rightarrow 0 \mid 1 \mid \dots \mid 9$$


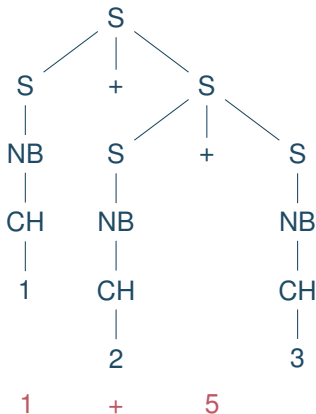
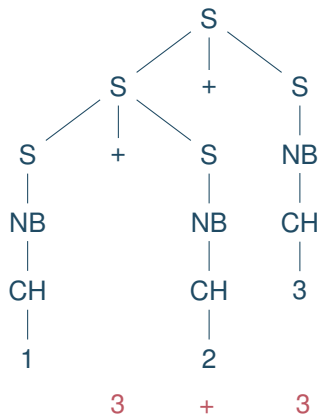


# Ambiguïté - Exemple

$S \rightarrow S + S \mid NB$

$NB \rightarrow NB\ CH \mid CH$

$CH \rightarrow 0 \mid 1 \mid \dots \mid 9$

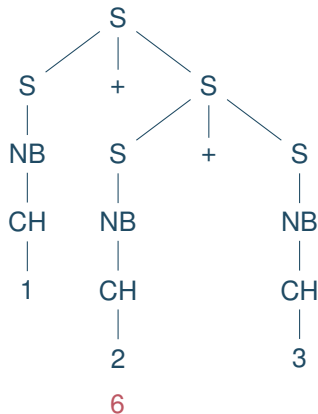
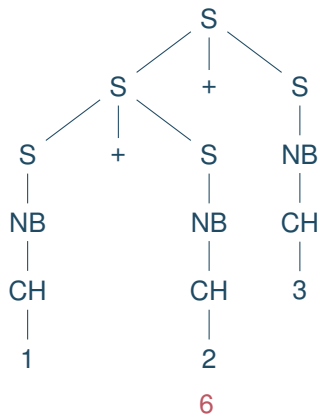


# Ambiguïté - Exemple

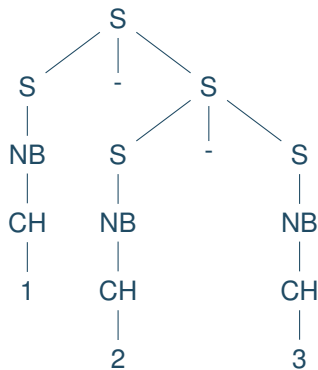
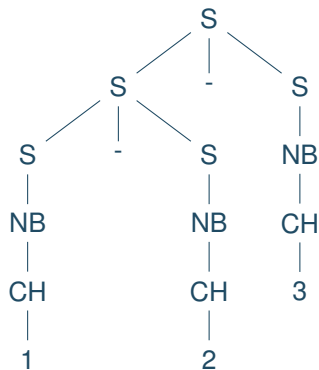
$S \rightarrow S + S \mid NB$

$NB \rightarrow NB\ CH \mid CH$

$CH \rightarrow 0 \mid 1 \mid \dots \mid 9$



# Ambiguïté - Gênant ?

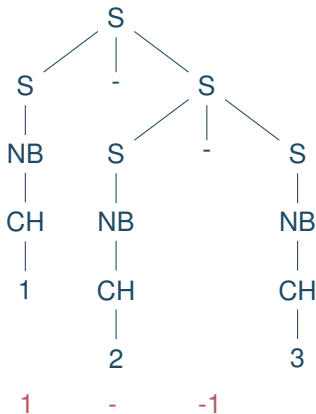
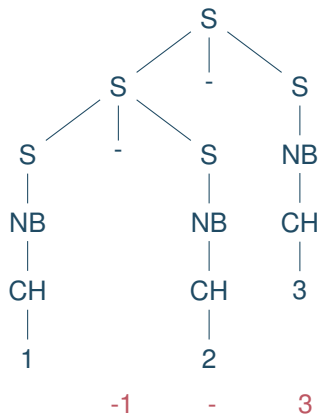
$$S \rightarrow S + S \mid NB$$
$$NB \rightarrow NB\ CH \mid CH$$
$$CH \rightarrow 0 \mid 1 \mid \dots \mid 9$$


# Ambiguïté - Gênant ?

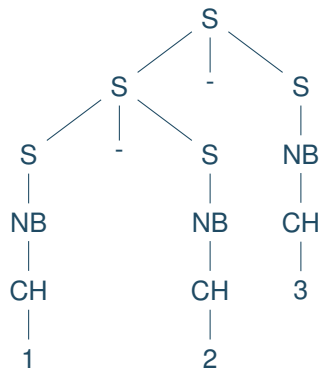
$S \rightarrow S + S \mid NB$

$NB \rightarrow NB\ CH \mid CH$

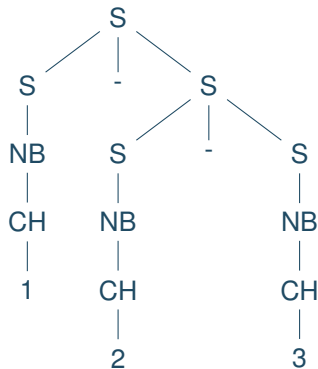
$CH \rightarrow 0 \mid 1 \mid \dots \mid 9$



# Ambiguïté - Gênant ?

$$S \longrightarrow S + S \mid NB$$
$$NB \longrightarrow NB\ CH \mid CH$$
$$CH \longrightarrow 0 \mid 1 \mid \dots \mid 9$$


-4



2

# Dérivations gauche et droite

Une dérivation **gauche** d'une grammaire hors contexte est une dérivation dans laquelle chaque pas de dérivation réécrit le non-terminal le plus à gauche

Une dérivation **droite** d'une grammaire hors contexte est une dérivation dans laquelle chaque pas de dérivation réécrit le non-terminal le plus à droite

Dans l'exemple précédent :

- dérivation gauche :

$$S \rightarrow S + S \rightarrow S + S + S \rightarrow NB + S + S \rightarrow CH + S + S \rightarrow 1 + S + S \rightarrow 1 + NB + S \rightarrow 1 + CH + S \rightarrow 1 + 2 + S \rightarrow 1 + 2 + NB \rightarrow 1 + 2 + CH \rightarrow 1 + 2 + 3$$

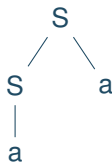
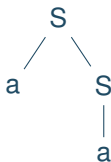
- dérivation droite :

$$S \rightarrow S + S \rightarrow S + S + S \rightarrow S + S + NB \rightarrow S + S + CH \rightarrow S + S + 3 \rightarrow S + NB + 3 \rightarrow S + CH + 3 \rightarrow S + 2 + 3 \rightarrow NB + 2 + 3 \rightarrow CH + 2 + 3 \rightarrow 1 + 2 + 3$$

# Ambiguïté

Une grammaire est **ambigüe** s'il existe un mot qui admet plusieurs arbres de dérivation

$$S \rightarrow aS + Sa + a$$



L'ambiguïté est relative à la grammaire : une grammaire ambigüe peut engendrer un langage qui peut également être engendré par une autre grammaire non ambigüe :  $S \rightarrow aS + a$


# Langage ambigu

➡ L'ambiguïté est fort embarrassante lors de l'analyse d'un langage de programmation par exemple

Un **langage ambigu** est un langage qui n'est engendré par aucune grammaire non ambiguë

## Exemple

$\{a^n b^n c^m \mid m, n \geq 0\} \cup \{a^m b^n c^n \mid m, n \geq 0\}$  est un langage ambigu

 Il n'existe pas d'algorithme permettant de décider si une grammaire est ambiguë



## Exemple de langage - Langage de Dyck

Le langage de Dyck est le langage des mots « bien parenthésés » avec  $n$  type de parenthèses,  $\Sigma = \{a_1, a'_1, \dots, a_n, a'_n\}$ , le langage de  $Dyck_n$  sur  $\Sigma$  est engendré par :

$$\begin{aligned} S &\longrightarrow ST \mid \varepsilon \\ T &\longrightarrow a_1 Sa'_1 \mid \dots \mid a_n Sa'_n \end{aligned}$$

$$\begin{aligned} S &\longrightarrow ST \longrightarrow Sa_1 Sa'_1 \longrightarrow STa_1 Sa'_1 \longrightarrow STa_1 STa'_1 \longrightarrow Sa_1 Sa'_1 a_1 STa'_1 \\ &\longrightarrow Sa_1 Sa'_1 a_1 Sa_2 Sa'_2 a'_1 \longrightarrow Sa_1 STa'_1 a_1 Sa_2 Sa'_2 a'_1 \\ &\longrightarrow Sa_1 Sa_3 Sa'_3 a'_1 a_1 Sa_2 Sa'_2 a'_1 \longrightarrow Sa_1 Sa_3 Sa'_3 a'_1 a_1 Sa_2 Sa'_2 a'_1 \\ &\longrightarrow^* a_1 a_3 a'_3 a'_1 a_1 a_2 a'_2 a'_1 \end{aligned}$$

Ce langage est algébrique, est-il rationnel ?

$$Dyck_1 \cap a_1^* a'_1{}^* = \{a_1^n a'_1{}^n \mid n \in \mathbb{N}\}$$

# Propriétés de clôture

## Proposition

Les langages algébriques sont clos par les opérations rationnelles (union, concaténation, étoile)

## Construction

Soient  $\mathcal{G} = (\Sigma, V, S, R)$  et  $\mathcal{G}' = (\Sigma, V', S', R')$  deux grammaires algébriques (avec  $V \cap V' = \emptyset$ )

- $\mathcal{L}(\mathcal{G}) \cup \mathcal{L}(\mathcal{G}')$  est engendré par la grammaire  $(\Sigma, V \cup V' \cup \{S''\}, S'', R \cup R' \cup \{S'' \rightarrow S + S'\})$
- $\mathcal{L}(\mathcal{G})\mathcal{L}(\mathcal{G}')$  est engendré par la grammaire  $(\Sigma, V \cup V' \cup \{S''\}, S'', R \cup R' \cup \{S'' \rightarrow SS'\})$
- $\mathcal{L}(\mathcal{G})^*$  est engendré par la grammaire  $(\Sigma, V \cup \{S''\}, S'', R \cup \{S'' \rightarrow S''S + \varepsilon\})$

# Grammaires

---

## Transformation de Grammaires context-free

Soit  $\mathcal{G} = (\Sigma, V, S, R)$  une grammaire, pour tout  $T \in V$  on note  $\mathcal{L}_{\mathcal{G}}(T) = \mathcal{L}(\mathcal{G}')$  avec  $\mathcal{G}' = (\Sigma, V, T, R)$

Une grammaire est dite **réduite** si

- $\forall T \in V, \mathcal{L}_{\mathcal{G}}(T) \neq \emptyset$   $\Rightarrow$  toute variable peut engendrer un mot
- $\forall T \in V, \exists \alpha, \beta \in (\Sigma \cup V)^*$  tels que  $S \xrightarrow{*} \alpha T \beta$   
 $\Rightarrow$  toute variable est atteignable à partir de  $S$

$\Rightarrow$  Il n'y a pas de variable trivialement inutile

## « Réduire » une grammaire

Proposition. Pour toute grammaire  $\mathcal{G} = (\Sigma, V, S, R)$ , il existe une grammaire  $\mathcal{G}' = (\Sigma, V', S', R')$  réduite telle que  $\mathcal{L}(\mathcal{G}) = \mathcal{L}(\mathcal{G}')$

Algorithme.

1. Calculer l'ensemble  $U$  des symboles de  $V$  qui peuvent engendrer un mot  
    ➡ Si  $S \notin U$ , la grammaire ne génère aucun mot
2. Supprimer les symboles de  $V \setminus U$  ainsi que toutes les règles dans lesquelles ils apparaissent
3. Calculer l'ensemble  $W$  des symboles de  $U$  atteignables
4. Supprimer les symboles de  $U \setminus W$  ainsi que toutes les règles dans lesquelles ils apparaissent

## « Réduire » une grammaire - Symboles productifs

*i.e.* les symboles qui peuvent engendrer un mot

Construisons l'ensemble par itérations :

- $U_0 = \Sigma$
- $\forall i, U_{i+1} = U_i \cup \{T \in V \mid T \rightarrow \alpha \text{ avec } \alpha \in U_i^*\}$
- Arrêter quand  $U_{i+1} = U_i$ . Noter que comme les  $U_i$  forment une suite croissante et que  $\Sigma \cup V$  est fini, la suite est forcément constante à partir d'un certain rang
- Alors  $U = U_i$

Construisons l'ensemble par itérations :

- $W_0 = \{S\}$
- $\forall i, W_{i+1} = W_i \cup \{T \in V \mid \exists T' \in W_i, T' \rightarrow \alpha T \beta \text{ avec } \alpha, \beta \in (\Sigma \cup V)^*\}$
- Arrêter quand  $W_{i+1} = W_i$ . Noter que comme les  $W_i$  forment une suite croissante et que  $V$  est fini, la suite est forcément constante à partir d'un certain rang
- Alors  $W = W_i$

## « Réduire » une grammaire - Exemple

$$\begin{aligned}S &\longrightarrow ST \mid a \mid \varepsilon \\T &\longrightarrow aTa \mid TX \\X &\longrightarrow b\end{aligned}$$

- Symboles productifs :  $U_0 = \{a, b\}$ ,  $U_1 = \{a, b, S, X\}$ ,  
 $U_2 = \{a, b, S, X\} = U_1$

- On obtient

$$\begin{aligned}S &\longrightarrow a \mid \varepsilon \\X &\longrightarrow b\end{aligned}$$

- Symboles accessibles :  $W_0 = \{S\}$ ,  $W_1 = \{S\} = W_0$
- On obtient la grammaire réduite

$$S \longrightarrow a \mid \varepsilon$$



Une grammaire algébrique est **propre** si elle ne contient pas :

- de règle  $T \rightarrow U$  avec  $T, U \in V$
- de règle  $T \rightarrow \varepsilon$

## Proposition

Tout langage algébrique ne contenant pas le mot vide peut être engendré par une grammaire propre

Remarque. Si on veut considérer un langage algébrique qui contient le mot vide on peut autoriser  $S \rightarrow \varepsilon$ , si  $S$  n'apparaît jamais en partie droite d'une règle

# Rendre une grammaire propre - mot vide

Soit  $\mathcal{G} = (\Sigma, V, S, R)$  une grammaire algébrique, il existe une grammaire  $\mathcal{G}' = (\Sigma, V', S', R')$  propre telle que  $\mathcal{L}(\mathcal{G}) = \mathcal{L}(\mathcal{G}')$

## Algorithme

1. Calculer l'ensemble  $U$  des symboles de  $V$  qui peuvent se dériver en  $\varepsilon$ 
  - 1.1  $U_0 = \{\varepsilon\}$
  - 1.2  $\forall i, U_{i+1} = U_i \cup \{T \in V \mid T \rightarrow \alpha \text{ avec } \alpha \in U_i^*\}$
  - 1.3 Arrêter quand  $U_{i+1} = U_i = U$
2. Pour toute règle  $T \rightarrow \alpha X \beta$  avec  $X \in U$  ajouter  $T \rightarrow \alpha \beta$
3. Supprimer toutes les règles  $T \rightarrow \varepsilon$

# Rendre une grammaire propre - variables équivalentes

Il faut à présent supprimer les règles de type  $T \rightarrow U$  avec  $T, U \in V$ .

Comme aucune variable n'engendre le mot vide,  $T \xrightarrow{*} U$  si et seulement si il existe  $T = T_0 \rightarrow T_1 \rightarrow \dots \rightarrow T_n = U$

Idée. Remplacer toute dérivation de type  $T \xrightarrow{*} U \rightarrow \alpha$  avec  $\alpha \notin V$  par une dérivation  $T \rightarrow \alpha$ .

1. Supprimer les règles  $T \rightarrow U$  avec  $T, U \in V$
2. Ajouter  $T \rightarrow \alpha$ , si  $T \xrightarrow{*} U$  et  $U \rightarrow \alpha$  avec  $T, U \in V$  et  $\alpha \notin V$

## Algorithme

- Calculer les paires de symboles de  $V$  telles que  $T \xrightarrow{*} U$  et  $U \xrightarrow{*} T$ , ce qui produit des classes d'équivalence de symboles de  $V$ . Remplacer tous les symboles d'une classe par un représentant choisi parmi eux
- Supprimer les règles de type  $T \rightarrow T$
- Pour chaque règle de type  $T \rightarrow U$ , et chaque règle  $U \rightarrow \alpha$ , on ajoute  $T \rightarrow \alpha$ , puis on supprime  $T \rightarrow U$

## Rendre une grammaire propre - Exemple (1/3)

$$\begin{aligned}S &\longrightarrow aTU \mid UT \mid b \\T &\longrightarrow UU \mid a \\U &\longrightarrow TU \mid b \mid bX \mid \varepsilon \\X &\longrightarrow Y \\Y &\longrightarrow aa\end{aligned}$$

Calculons l'ensemble  $U$  des symboles de  $V$  qui peuvent se dériver en  $\varepsilon$  :  
 $U_0 = \{\varepsilon\}$ ,  $U_1 = \{U\}$ ,  $U_2 = \{U, T\}$ ,  $U_3 = \{U, T, S\}$  Donc

$$\begin{aligned}S &\longrightarrow aTU \mid UT \mid b \mid aU \mid aT \mid U \mid T \\T &\longrightarrow UU \mid a \mid U \\U &\longrightarrow TU \mid b \mid bX \mid U \mid T \\X &\longrightarrow Y \\Y &\longrightarrow aa\end{aligned}$$

Calculons les classes d'équivalences de symboles :  $\{S\}$ ,  $\{T, U\}$ ,  $\{X\}$ ,  $\{Y\}$

## Rendre une grammaire propre - Exemple (2/3)

Choisissons  $T$  :

$$S \longrightarrow aTT \mid TT \mid b \mid aT \mid aT \mid T \mid T$$

$$T \longrightarrow TT \mid a \mid T$$

$$T \longrightarrow TT \mid b \mid bX \mid T \mid T$$

$$X \longrightarrow Y$$

$$Y \longrightarrow aa$$

Supprimons les redondances

$$S \longrightarrow aTT \mid TT \mid b \mid aT \mid T$$

$$T \longrightarrow TT \mid a \mid T \mid b \mid bX$$

$$X \longrightarrow Y$$

$$Y \longrightarrow aa$$

## Rendre une grammaire propre - Exemple (3/3)

Supprimons les règles de type  $T \rightarrow T$

$$\begin{aligned}S &\longrightarrow aTT \mid TT \mid b \mid aT \mid T \\T &\longrightarrow TT \mid a \mid b \mid bX \\X &\longrightarrow Y \\Y &\longrightarrow aa\end{aligned}$$

Supprimons les règles de type  $T \rightarrow U$

1. Supprimons  $S \rightarrow T$

$$\begin{aligned}S &\longrightarrow aTT \mid TT \mid b \mid aT \mid a \mid bX \\T &\longrightarrow TT \mid a \mid b \mid bX \\X &\longrightarrow Y \\Y &\longrightarrow aa\end{aligned}$$

2. Supprimons  $X \rightarrow Y$

$$\begin{aligned}S &\longrightarrow aTT \mid TT \mid b \mid aT \mid a \mid bX \\T &\longrightarrow TT \mid a \mid b \mid bX \\X &\longrightarrow aa\end{aligned}$$

# Forme normale de Chomsky

Une grammaire  $\mathcal{G} = (\Sigma, V, S, R)$  est en **forme normale de Chomsky** si toute règle est de la forme :

$$\begin{array}{ll} T \rightarrow UW & \text{avec } T, U, W \in V \\ T \rightarrow a & \text{avec } T \in V, a \in \Sigma \end{array}$$

## Remarques

- la grammaire est donc propre
- la grammaire ne peut engendrer le mot vide
- une variante autorise la règle  $S \rightarrow \varepsilon$  si  $S$  n'apparaît jamais en partie droite
- l'arbre de dérivation est donc un arbre binaire

# Mettre en forme normale de Chomsky

Soit  $\mathcal{G} = (\Sigma, V, S, R)$  une grammaire algébrique réduite et propre

La transformation s'effectue en 2 étapes

- $\forall x \in \Sigma$ , ajouter une variable  $V_x$  et une règle  $V_x \rightarrow x$ , puis remplacer  $x$  par  $V_x$  dans toutes les règles de  $R$

Toutes les règles sont désormais de la forme  $T \rightarrow T_1 \dots T_n$  ou  $T \rightarrow x$

- Pour chaque règle  $T \rightarrow T_1 \dots T_n$  avec  $n \geq 2$  ajouter de nouvelles variables  $U_1 \dots U_{n-2}$  et remplacer la règle par

$$\begin{array}{lll} T & \rightarrow & T_1 U_1 \\ U_1 & \rightarrow & T_2 U_2 \\ U_2 & \rightarrow & T_3 U_3 \\ \dots & & \\ U_{n-2} & \rightarrow & T_{n-1} T_n \end{array}$$



## Décider si un mot appartient à un langage algébrique

Soient  $\mathcal{G} = (\Sigma, V, S, R)$  une grammaire algébrique en forme normale de Chomsky et  $u \in \Sigma^*$ . Il est possible de **décider en temps fini** si  $u$  appartient à  $\mathcal{L}(\mathcal{G})$

Comment ? La partie droite  $\alpha$  de toute règle de  $R$  est telle que soit  $\alpha \in \Sigma$  soit  $\alpha \in V^2$  donc une dérivation qui engendre le mot  $u$  a une longueur d'au plus  $2|u|$ , il suffit donc d'énumérer toutes les dérivations possibles pour tester si  $u$  est engendré par la grammaire

# Grammaires

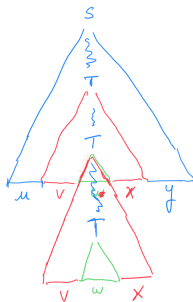
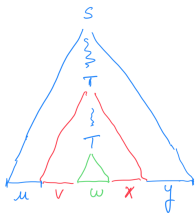
---

## Lemme d'itération

## Lemme d'itération

Soit  $\mathcal{L}$  un langage context-free, il existe un entier  $k$  tel que tout mot  $u \in \mathcal{L}$  avec  $|u| > k$  se factorise en  $u = vwxyz$  avec  $wy \neq \varepsilon$ ,  $|wxy| < k$  et pour tout  $n$ ,  $vw^nxy^n z \in \mathcal{L}$

# Lemme d'itération - Idée



## Montrer qu'un langage n'est pas context-free

Soit  $\mathcal{L} = \{a^n b^n c^n \mid n \in \mathbb{N}\}$ , supposons qu'il est context-free

Soient  $k$  l'entier donné par le lemme d'itération et  $u$  tel que  $|u| > k$ , alors  $u$  se factorise en  $u = vwxyz$  avec  $u_n = vw^n xy^n z \in \mathcal{L}$  pour tout  $n$   $wy \neq \varepsilon$ ,  $|wxy| < k$

Si  $w$  ou  $y$  sont composés d'au moins 2 lettres différentes de  $\{a, b, c\}$ , alors clairement les  $u_n$  ne peuvent appartenir à  $\mathcal{L}$  puisque la répétition des  $w$  et des  $y$  introduit une alternance de lettres impossible dans le langage

Donc  $w$  et  $y$  appartiennent à  $a^* + b^* + c^*$ , or par hypothèse,  $vwxyz$  contient autant d'occurrences de  $a$ , de  $b$  et de  $c$ , il est donc impossible que  $vw^n xy^n z$  dans lequel on a modifié le nombre d'occurrences de une ou deux des trois lettres contienne toujours le même nombre d'occurrences des trois lettres

Contradiction