

# Langages et Automates

Cours : CM1 Introduction

3 x 45 minutes

Clément AGRET

5 min pause

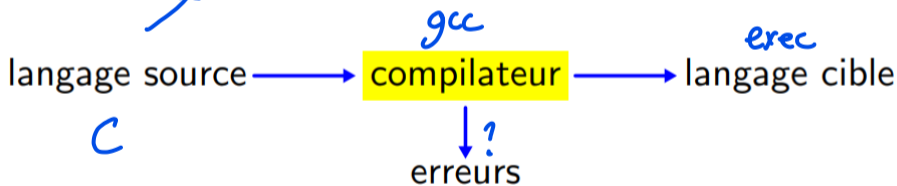
12h30 → 13h15  
13h20 → 14h05  
14h10 → 14h55

) ok

# Compilateur

Un compilateur est un utilitaire de traduction. Qui traduit un programme d'un langage source vers un langage cible, en signalant d'éventuelles erreurs.

Programmeur (Humain)	Système d'exploitation
<u>C/C++</u> , Pascal, Algol, FORTRAN, assembleur, ..	Exécutable(Binaire)
<u>for (i = 0; i &lt;= 10; i++) s += 1;</u>	1000111010100



# Compilateur

1954/57 → FORTRAN (Mathematical FORmula TRANslating System)

Plus tard, → PASCAL

Nombreux autres compilateurs ont suivi → YACC, (Yet Another Compiler of Compiler).

1955/65 → Des linguistes, philosophes et mathématiciens ont défriché la partie théorique en proposant une description et une classification des langages et des grammaires.

# Les tâches d'analyse d'un compilateur.

$$A + B = C \quad \text{trop-}$$

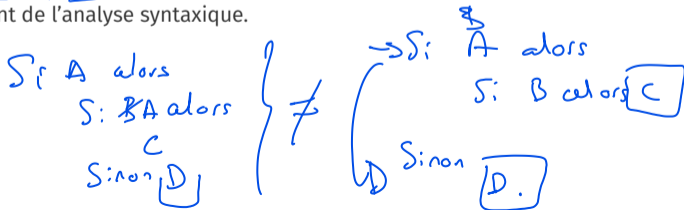
$(A + B) = C$

Handwritten notes: A blue arrow points to the opening parenthesis of  $(A + B)$ . Another blue arrow points to the closing parenthesis, with the word "trop-" written above it. The entire expression  $(A + B) = C$  is underlined in blue.

Les premières tâches d'un compilateur sont de faire :

- 1) l'analyse lexicale : reconnaître les éléments constitutifs de la chaîne entrée, c'est-à-dire du code source, et en dresser la liste.
- 2) l'analyse syntaxique : vérifier la conformité avec les règles de constitution du code. Par exemple, l'expression  $(A+B) = C$  est syntaxiquement incorrecte (parentheses...).
- 3) l'analyse sémantique : analyser le sens et fixer une interprétation.

Par exemple dans si A alors si B alors C sinon D, choisir à quel **si** se rapporte le **sinon**. Dans ce qui suit, nous nous occuperons essentiellement de l'analyse syntaxique.



# La notion de grammaire.

$$\Sigma = \{L, E, V, I, U, X, C, H, A, T$$

$R, U\}$

espace  $\hookrightarrow$  Space.

Considérons, par exemple, la phrase suivante :

LE VIEUX CHAT ATTRAPE LE PETIT RAT

Le but est de construire une grammaire qui permette de produire cette phrase.

Il faut tout d'abord préciser les « ingrédients » nécessaires :

Ce sera l'alphabet des symboles ou lexique.

Dans notre cas, nous pouvons prendre comme alphabet l'ensemble:

$$\Sigma = \{LE, VIEUX, PETIT, CHAT, RAT, ATTRAPE\}$$

$\Sigma$  ok?

## Analyse syntaxique.

$\Sigma$  du langage C :  $\Sigma = \{\text{int, char, for, do, ...}\}$

Noter que le terme «alphabet» n'a pas ici le sens habituel A, B, C,...

L'alphabet contient les «briques de base» avec lesquelles on peut former tout ce qu'on veut former.

Dans un langage de programmation, les mots réservés, balises, etc...

Comme : program, real, <head>, </head>, sont dans l'alphabet de symboles.

→ français

# Analyse syntaxique.



FRANÇAIS  $\neq$  C++  $\neq$  Chinois.

Voyons maintenant comment les symboles sont assemblés.

Dans la structure de la phrase, on peut distinguer :

- . un groupe sujet
- . un verbe
- . un groupe complément d'objet (CO)

Les **groupes sujet** et **CO** sont eux-mêmes des groupes nominaux ; un groupe nominal est formé d'un article suivi d'un nom, lui-même précédé ou suivi d'adjectifs.

# Analyse syntaxique.

phrase → LA VIEUX CHAT MANGE LE PETIT RAT  
LE VIEUX RAT ATTRAPE LE VIEUX RAT

Voici un exemple de (petite) grammaire pouvant produire notre phrase; elle a 7 règles de grammaire (on dit aussi de production, ou de réécriture) :

1. <phrase> → <groupe sujet> <verbe> <groupe CO>

2. <groupe sujet> → <groupe nominal>

3. <groupe nominal> → <article> <adjectif> <nom>

4. <article> → LE, LA

5. <nom> → RAT, CHAT

6. <adjectif> → VIEUX, PETIT

7. <verbe> → ATTRAPE, MANGE

syntaxe

OK  
Non OK

→ EASY  
→ HARD

?

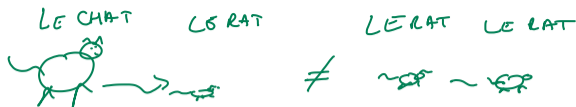
⌘ : <groupe CO> → <article adjectif, nom>



# Suite

- Le point de départ s'appelle l'axiome. Dans notre exemple, c'est <phrase>.
- Les **variables** sont les «ingrédients» qui peuvent encore être remplacés par d'autres (<phrase>, <article>, <groupe CO>, etc.).
- Les règles 1 à 3 sont des règles syntaxiques (la première règle concerne toujours l'axiome). Les règles 4 à 7 sont des règles complètement terminales (ou : lexicales).
- Le langage engendré par la grammaire est l'ensemble de toutes les phrases que l'on peut produire à partir de l'axiome en utilisant des règles de grammaire, une phrase étant une chaîne de symboles.

## Exemples de phrases



Phrases « grammaticalement correctes », c'est-à-dire des chaînes de symboles que l'on peut produire à partir de l'axiome, en utilisant les règles précédentes :

[LE RAT ATTRAPE LE VIEUX RAT] *c'est ok ?*

On peut associer à chaque phrase ainsi produite un arbre de dérivation où figurent les variables auxquelles on a appliqué des règles de grammaire.

Nous avons ajouté aux nœuds de cet arbre de dérivation, pour faciliter la compréhension, les numéros des règles qui ont été appliquées aux variables.

# Arbre de dérivation

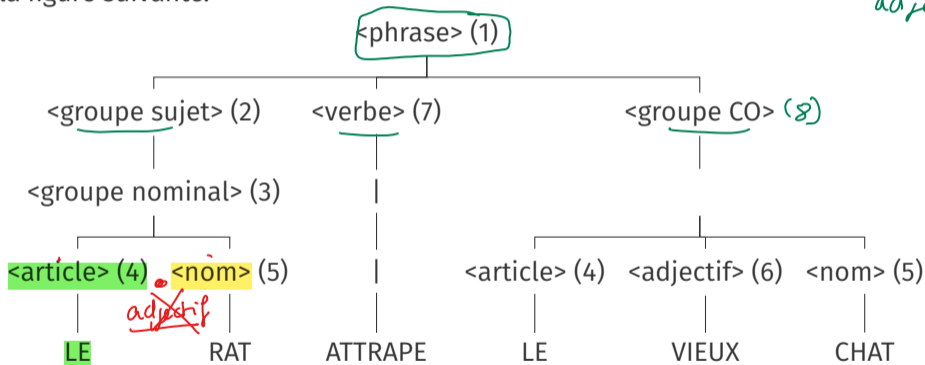


article nom adjectif  
≠

nom article adjectif

adjectif nom article  
≠

Un arbre de dérivation pour la phrase LE RAT ATTRAPE LE PETIT CHAT est représenté dans la figure suivante.



[LE ~~VIEUX~~ RAT ATTRAPE LE PETIT CHAT] est-ce correct ? Non

≠ nom

# Définitions.

a partir d'un alphabet de symbole on peut créer une phrase.

L'alphabet des symboles est un ensemble fini. On le notera en général  $\Sigma$  dans la suite du cours.

$$1 \quad \Sigma = \{\text{LE, VIEUX, PETIT, CHAT, RAT, ATTRAPE}\}$$
$$2 \quad \Sigma = \{0, 1, 2, \dots, 9, +, \times, -, /, (, )\}$$

Le second alphabet permet d'écrire les expressions arithmétiques sur les nombres entiers.

Il y a notamment les quatre opérations et les parenthèses.

Le premier problème sera de pouvoir distinguer si une expression est « correctement écrite ».

## Exemple

$2 \times (31 - 6) + 8$  est correcte, mais  $2(31 - 6) + 8$  ou  $(21+) \times 4$  ne le sont pas.

Une chaîne est une suite finie de symboles.

La longueur d'une chaîne est le nombre de ses symboles.

OK?

# Exemple

$$\Sigma \quad \Sigma^0 = \{\epsilon\} \quad \epsilon^*$$

**LE LE CHAT** est une chaîne de longueur 3 sur le premier alphabet. Les expressions arithmétiques (correctes ou non !) sont des chaînes sur le second alphabet. La chaîne vide, notée  $\lambda$  ne contient aucun symbole :  $\epsilon$ , sa longueur est nulle.

→  $\Sigma^n$  : L'ensemble de toutes les chaînes de longueur n.

$\Sigma^0 = \sigma$  : Attention, cet ensemble n'est pas vide : il contient la chaîne vide  $\epsilon$  !

→  $\Sigma^*$  : C'est  $\Sigma^n$  pour  $n \geq 0$  C'est donc l'ensemble de toutes les chaînes, chaîne vide comprise.

→  $\Sigma^+$  : Est la réunion des  $\Sigma^n$  pour  $n \geq 1$ .

$$\Sigma^n = \{ \Sigma^0 + \Sigma^1 + \Sigma^2 + \dots \}$$

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

$$n=1 \quad \Sigma^1 = \{a, b\}$$

$$n=2 \quad \Sigma^2 = \{aa, ab, ba, bb\}$$

$$\Sigma^3 = \{aaa, aab, aba, \dots, bbb\}$$

taille  $\Sigma^2$  |  $|\Sigma^2| = 4$  |  $|\Sigma^3| = 2^3 = 8 + 1$   
5 OK  $\epsilon$

# Concaténation

$$\Sigma = \{a, b\} \quad \Sigma = \{\varepsilon, a, b\}.$$

$$u = aba$$
$$v = bbb.$$

$$Q? \quad u.v = ?$$

$$uv = \underline{ababbb}$$

On peut concaténer des chaînes de symboles, c'est-à-dire les «coller» les unes derrière les autres, de la même façon que plusieurs textes peuvent être assemblés les uns derrière les autres pour former un nouveau texte.

Si  $u$  et  $v$  sont deux chaînes de symboles, leur concaténation sera notée  $u.v$  ou plus simplement  $uv$ .

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

# Exemples.



Alphabet  $\neq$  Langage

Pouvez vous me donner un mot que je n'ai pas listé qui est dans l'ensemble  $L_2$

Sur l'alphabet  $\Sigma = \{a, b, c, d\}$ , les ensembles suivants sont des langages :

$L_1 = \{ac, abbc\}$

$e$

$L_2$  de toutes les chaînes qui commencent par  $a$ .

$$L_2 = \{a, aa, ab, ac, ad, \dots \\ aaaaaa \dots - \\ abd \text{ (ok)}\}$$

$L_3$  de toutes les chaînes de longueur paire.

## Notation.

Afin de faciliter la lecture, nous adopterons en général une notation  $a, b, c, \dots$  pour des symboles de  $\Sigma$  et  $u, v, w, x, \dots$  pour des chaînes de symboles. (mots)

Symboles  
 $a, b, c$

mots  
 $u = aab-$   
 $v = bab$   
 $w = aacb-$

$$L_3 = \{aa, ab, ac, ad, \\ bb, \dots aaaa-\}$$

# Opération sur les langages.

Un alphabet de symboles  $\Sigma$  étant fixé, voyons les opérations que nous appliquerons aux langages sur  $\Sigma$ .

1. Opérations ensemblistes usuelles (réunion, intersection, complémentaire).
2. Concaténation de langages.

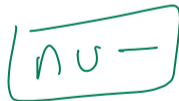


$u$  et  $v$

$$u \cdot v \neq v \cdot u$$

$$u \cap v \quad \bar{u}$$

$$u \cup v$$



ok?



# Opération ensemblistes

$\Sigma^*$      $\Sigma^n$  et     $\Sigma^+$

Les langages sur  $\Sigma$  ne sont rien d'autre que des sous-ensembles de  $\Sigma^*$ .  
Les opérations habituelles que l'on connaît sur les sous-ensembles s'appliquent à ces langages.

# Concaténation de langages

$$uv = uv$$

aba  
bbb

$$uv = ababbb$$

Comme nous savons concaténer des chaînes de symboles, nous pourrions également concaténer deux langages, c'est-à-dire collecter dans un nouvel ensemble, noté  $L_1.L_2$ , toutes les chaînes que l'on peut obtenir en concaténant une chaîne de  $L_1$  avec une chaîne de  $L_2$ .

$$L_1 = \{aa, ab\}$$

$$L_2 = \{ba, bb\}$$

$$L_1 \cdot L_2 = \{aaba, aabb, abba, abbb\}$$

OK?

## Concaténation $L1.L2$

$$L = \{a, b, ab\}$$

$$L^2 = \{aa, ab, aab, ba, bb, bab, aba, abb, abab\}$$

Par récursivité, nous pourrions concaténer aussi un nombre fini de langages. En particulier, si  $L$  est un langage, on notera :

$$L^0 = \{\lambda\}$$

$$L^1 = L$$

$$L^2 = L.L$$

Puis on définit de manière récursive

$$L^n = L^{n-1}.L \text{ (Concaténation de } n \text{ exemplaires de } L)$$

$$L^3 = L^2.L$$

$$L^3 = \{aaa, aba, aaba, \dots\}$$