Les grammaires formelles

max.silberztein@univ-fcomte.fr

Introduction

- 1. Langages, grammaires et machines
- 2. Grammaires régulières
- 3. Grammaires hors-contexte
- 4. Grammaires contextuelles
- 5. Grammaires non restreintes

Références

- Beesley K., Karttunen L. Finite-state morphology: Xerox tools and techniques. CSLI, Stanford. 2003:359-75.
- Chomsky N., 1957. Syntactic Structures. Mouton & Co Eds. N.V.
- Dalrymple, M., 2001. Lexical functional grammar (Vol. 34). Brill.
- Gazdar G., Klein E., Pullum G., Sag. I., 1985. Generalized Phrase Structure Grammar. Blackwell Publishing, Oxford, England, and Harvard University Press, Cambridge, Massachusetts.
- Gross M., Lentin A. 1970. Notions sur les grammaires formelles. Gauthier-Villars Eds., Paris.
- Pollard, Carl & Ivan A. Sag. 1994. Head-Driven Phrase Structure Grammar (Studies in Contemporary Linguistics 4). Chicago, IL: The University of Chicago Press.
- Silberztein M., 2015. La formalisation des langues : l'approche de NooJ. ISTE, Eds., Londres.

https://fr.wikipedia.org/wiki/Grammaire_formelle

https://fr.wikipedia.org/wiki/Langage_formel

https://fr.wikipedia.org/wiki/Automate_fini

1. Langages, grammaires et machines

Un triple problème

 Les langues naturelles permettent de construire une infinité de phrases. Comment décrire mathématiquement un ensemble dont la taille est infinie ?

• Comment distinguer les phrases des séquences incorrectes : Luc mange une pomme vs. Luc une mange pomme

- Convergence entre :
 - Linguistes (ex. Chomsky) : décrire les langues
 - Mathématiciens (ex. Schützenberger) : décrire les ensembles infinis
 - Informaticiens (ex. Gödel) : caractériser ce qu'on peut calculer

Langages formels, définition

Alphabet et lettres :

- Un alphabet est un ensemble fini d'éléments. Ses éléments sont les lettres de l'alphabet
- Ex. L'alphabet A_v = { a, e, i, o, u } contient 5 lettres
- Ex. L'alphabet A_F = { a, b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q, r, s, t, u, v, w, x, y, z } contient 26 lettres

Mots et langages

 A partir d'un alphabet donné, on peut construire des suites constituées de lettres : les mots

Un langage est un ensemble de mots

- Ex. L₁ = { aime, aimes, aimons, aimez, aiment }
- Ex. L₂ = { mais, ou, et, donc, or, ni, car }
- Ex. $L_3 = \{ dès, ès, grès, près, très \}$

 Comment décrire les langages qui contiennent un nombre infini de mots ?

Mot vide, langage vide, singleton du mot vide

- Le mot vide est constitué de 0 lettre
 - Les linguistes et les mathématiciens le notent ε
 - Les informaticiens le notent ""
 - Dans NooJ, on le note <E>
- Le langage vide ne contient aucun mot, ex. { } ou Ø
- Le singleton { ε } contient 1 élément : le mot vide

Ne pas confondre mot vide et langage vide

Les deux niveaux de langages

- Orthographe : Lettres → Alphabet → Mots → Langages :
 - les lettres, ex. « a », « ê », « ! », « . »
 - l'alphabet, ex. { a, b, c, ... z }
 - Le langage, ex. { mais, ou, et, donc, or, ni, car }
- Syntaxe : ALUs → Vocabulaire → Phrases → Langages
 - Unités linguistiques Atomiques (ALUs) = éléments du vocabulaire
 - Vocabulaire, ex. { lapin, laitue romaine, la, le, mange }
 - Ex. Langage = { « le lapin la mange », « le lapin mange la laitue romaine », « le lapin mange » }

Le niveau syntaxique : \rightarrow ALUs \rightarrow Vocabulaire \rightarrow Phrases \rightarrow Langages

- Unités Linguistiques Atomiques (ALUs) = éléments du vocabulaire :
 - Mots simples : ALUs qui s'écrivent sous la forme de séquences de lettres entre deux séparateurs, ex. *table*
 - **Affixes** (préfixes ou suffixes), ALUs qui s'écrivent sous la forme de séquences de lettres non délimitée par deux séparateurs, ex. *remanger*, *manifestation*
 - Mots composés : ALUs qui s'écrivent sous la forme de séquences de lettres et de séparateurs, ex. parce que, aujourd'hui, sous-marin
 - **Expressions**: ALUs qui s'écrivent sous la forme de séquences de formes potentiellement discontinues, ex. *avoir ... lieu, prendre ... une douche, accuser ... le coup, prendre ... le taureau par les cornes*

Les deux niveaux de grammaires

• Une grammaire orthographique est un ensemble fini de règles qui permet d'identifier les mots d'un langage à partir de ses lettres

• Une grammaire syntaxique est un ensemble fini de règles qui permet d'identifier les phrases d'un langage à partir de ses ALUs

Equivalence Langages ⇔ Grammaires

• Un langage est parfaitement décrit par une grammaire

• A partir d'une grammaire orthographique, on peut produire tous les mots du langage décrit par la grammaire, et on peut vérifier qu'un mot donné appartient ou non au langage décrit

 A partir d'une grammaire syntaxique, on peut produire toutes les phrases du langage décrit par la grammaire, et on peut vérifier qu'une séquence de mots donnée appartient ou non au langage décrit

Chomsky (linguiste américain) et Schützenberger (mathématicien français) ont inventé un mécanisme mathématique qui permet de décrire des langages : les grammaires génératives.

- Une grammaire générative est un ensemble fini de règles de réécriture
- Chaque règle de réécriture est constituée d'un membre gauche, d'une flèche et d'un membre droit : $\alpha \rightarrow \beta$
- Exemple d'une grammaire générative syntaxique :

```
PHRASE \rightarrow GN voit GN
GN \rightarrow le NOM
GN \rightarrow un NOM
NOM \rightarrow chat
NOM \rightarrow chien
```

- Chaque membre gauche ou droit d'une règle de réécriture est une séquence de symboles auxiliaires et/ou de symboles terminaux : lettres (orthographe) ou ALU (syntaxe)
- Il doit y avoir un symbole auxiliaire initial, ex. PHRASE

```
PHRASE \rightarrow GN voit GN
GN \rightarrow le NOM
GN \rightarrow un NOM
NOM \rightarrow chat
NOM \rightarrow chien
```

- Les noms des symboles auxiliaires sont arbitraires et n'ont aucune importance
- Les symboles auxiliaires peuvent être supprimés par des optimisations automatiques, sans rien changer au langage décrit
- Inutile donc de développer des théories linguistiques sur l'existence ou non de symboles auxiliaires (ex. VP pour groupes prépositionnels)...

• En partant du symbole auxiliaire initial PHRASE, on effectue n'importe quelle réécriture, c'est-à-dire qu'on remplace n'importe quel membre gauche d'une règle par son membre droit, jusqu'à épuisement des possibilités de réécriture. Par exemple :

PHRASE → GN voit GN
GN → le NOM
GN → un NOM
NOM → chat
NOM → chien

PHRASE

• En partant du symbole auxiliaire initial PHRASE, on effectue n'importe quelle réécriture, c'est-à-dire qu'on remplace n'importe quel membre gauche d'une règle par son membre droit, jusqu'à épuisement des possibilités de réécriture. Par exemple :

PHRASE => GN voit GN

PHRASE → GN voit GN
GN → le NOM
GN → un NOM
NOM → chat
NOM → chien

• En partant du symbole auxiliaire initial PHRASE, on effectue n'importe quelle réécriture, c'est-à-dire qu'on remplace n'importe quel membre gauche d'une règle par son membre droit, jusqu'à épuisement des possibilités de réécriture. Par exemple :

PHRASE => GN voit GN => le NOM voit GN

PHRASE → GN voit GN
GN → le NOM
GN → un NOM
NOM → chat
NOM → chien

• En partant du symbole auxiliaire initial PHRASE, on effectue n'importe quelle réécriture, c'est-à-dire qu'on remplace n'importe quel membre gauche d'une règle par son membre droit, jusqu'à épuisement des possibilités de réécriture. Par exemple :

PHRASE → GN voit GN
GN → le NOM
GN → un NOM
NOM → chat
NOM → chien

PHRASE => GN voit GN => le NOM voit GN

=> le NOM voit un NOM

• En partant du symbole auxiliaire initial PHRASE, on effectue n'importe quelle réécriture, c'est-à-dire qu'on remplace n'importe quel membre gauche d'une règle par son membre droit, jusqu'à épuisement des possibilités de réécriture. Par exemple :

PHRASE => <u>GN</u> voit GN => le <u>NOM</u> voit GN => le <u>NOM</u> voit un <u>NOM</u> => le <u>NOM</u> voit un chat PHRASE → GN voit GN
GN → le NOM
GN → un NOM
NOM → chat
NOM → chien

• En partant du symbole auxiliaire initial PHRASE, on effectue n'importe quelle réécriture, c'est-à-dire qu'on remplace n'importe quel membre gauche d'une règle par son membre droit, jusqu'à épuisement des possibilités de réécriture. Par exemple :

PHRASE \rightarrow GN voit GN GN \rightarrow le NOM GN \rightarrow un NOM NOM \rightarrow chat NOM \rightarrow chien

PHRASE => GN voit GN => le NOM voit GN

- => le NOM voit un NOM => le NOM voit un chat
- => le chien voit un chat

Il n'y a plus de symbole auxiliaire ; on ne peut plus effectuer de réécriture ; la séquence *le chien voit un chat* appartient au langage décrit par la grammaire.

Exercice: grammaire, dérivation, langage

• Trouver d'autres dérivations potentielles pour la grammaire :

```
PHRASE → GN voit GN
GN → le NOM
GN → un NOM
NOM → chat
NOM → chien
```

• Combien de phrases cette grammaire peut-elle produite ? Quelle est la taille du langage décrite par cette grammaire ?

Exercice: grammaire, dérivation, langage

Si on calcule toutes les dérivations possibles pour la grammaire précédente, on obtient l'ensemble des 16 phrases suivantes :

{ "un chat voit un chat", "un chat voit un chien", "un chat voit le chat", "un chat voit le chien", "le chat voit un chat", "le chat voit le chat", "le chat voit le chien", "un chien voit un chat", "un chien voit un chien", "un chien voit le chat", "un chien voit le chien", "le chien voit un chat", "le chien voit un chien", "le chien voit le chien" }

PHRASE \rightarrow GN voit GN
GN \rightarrow le NOM
GN \rightarrow un NOM
NOM \rightarrow chat
NOM \rightarrow chien

Traitement d'une grammaire par une machine : analyse

• En partant d'une séquence quelconque, on effectue n'importe quelle réécriture inverse, c'est-à-dire qu'on remplace n'importe quel membre droit d'une règle par son membre gauche, jusqu'à épuisement des possibilités de réécriture. Par exemple :

PHRASE \rightarrow GN voit GN
GN \rightarrow le NOM
GN \rightarrow un NOM
NOM \rightarrow chat
NOM \rightarrow chien

le <u>chien</u> voit un chat => le **NOM** voit un <u>chat</u> => le **NOM** voit **GN** => <u>GN voit GN</u> => PHRASE

On a pu produire PHRASE à partir de la séquence *le chien voit un chat*; cette séquence est donc une phrase du langage décrit par la grammaire

Traitement d'une grammaire par une machine : analyse

• En partant d'une séquence quelconque, on effectue n'importe quelle réécriture inverse, c'est-à-dire qu'on remplace n'importe quel membre droit d'une règle par son membre gauche, jusqu'à épuisement des possibilités de réécriture. Par exemple :

le chien voit un cheval => le <u>NOM</u> voit un cheval => le <u>NOM</u> voit un cheval => le <u>NOM</u> voit un cheval => <u>GN</u> voit un cheval

PHRASE \rightarrow GN voit GN GN \rightarrow le NOM GN \rightarrow un NOM NOM \rightarrow chat NOM \rightarrow chien

On ne peut plus effectuer de réécriture inverse ; on n'a pas pu produire PHRASE à partir de la séquence *le chien voit un cheval* ; cette séquence n'est donc pas une phrase du langage décrit par la grammaire

Langages / Grammaires

• Un langage correspond à une grammaire

- Une machine peut fonctionner dans deux sens :
 - Un **générateur** part d'un symbole auxiliaire initial (ex. **PHRASE**), et effectue des réécritures jusqu'à produire des phrases du langage
 - Un **analyseur** vérifie si une séquence donnée appartient à un langage en essayant d'effectuer les réécritures inverses, jusqu'à produire le symbole auxiliaire initial (ex. **PHRASE**)

Hiérarchie de Chomsky-Schützenberger

• Quatre types de grammaires génératives (et donc de langages, et donc de machines) :

- Grammaires régulières (Type 3)
- Grammaires hors contexte (Type 2)
- Grammaires contextuelles (Type 1)
- Grammaires non restreintes (Type 0)

Hiérarchie de Chomsky-Schützenberger

• Quatre types de grammaires génératives, de langages, et de machines :

Grammaires régulières Langages réguliers (ou rationnels) Automates finis

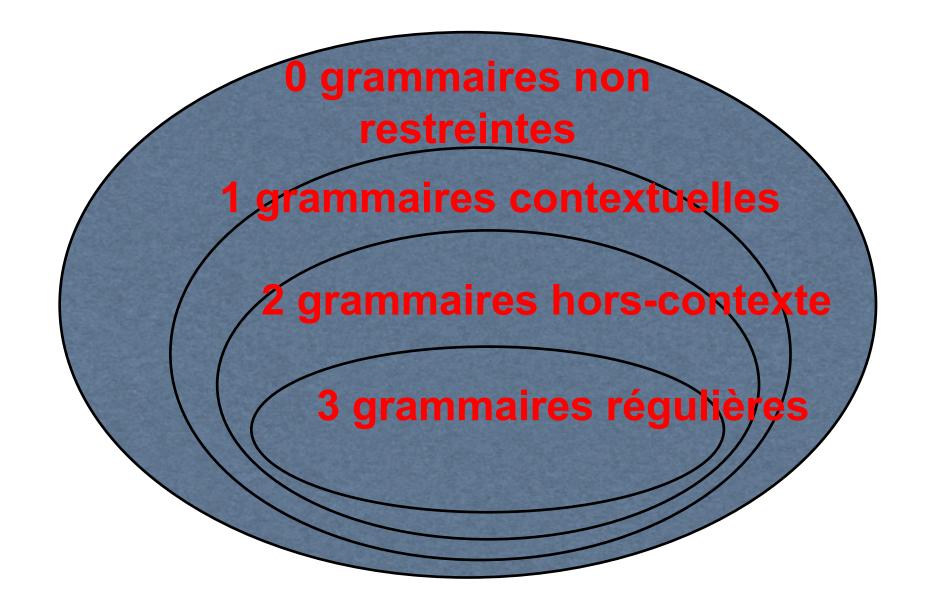
Grammaires hors contexte Langages algébriques Automates à pile

Grammaires contextuelles Langages contextuels Automates linéairement bornés

Grammaires non restreintes Langages récursivement énumérables Machines de Turing

- Les grammaires hors-contexte sont plus puissantes que les grammaires régulières : elles peuvent décrire des langages que les grammaires régulières ne peuvent pas décrire.
- Les grammaires contextuelles sont plus puissantes que les grammaires hors-contexte
- les grammaires non restreintes sont plus puissantes que les grammaires contextuelles.

Hiérarchie de Chomsky-Schützenberger



Besoin de formalismes différents

• Les gramaires génératives sont beaucoup trop lourdes pour être utilisées pour décrire des langues naturelles. Par exemple, pour décrire la conjugaison des verbes du premier groupe comme *aider*, il faudrait 6 règles pour chaque temps, et plusieurs milliers de règles <u>racine</u>:

 Conjugaison → Présent
 Présent → racine e
 racine → aid

 Présent → racine es
 racine → aim

 Présent → racine e
 racine → amus

 Présent → racine ons
 ...

 Présent → racine ez
 Présent → racine ent

Besoin de formalismes adaptés

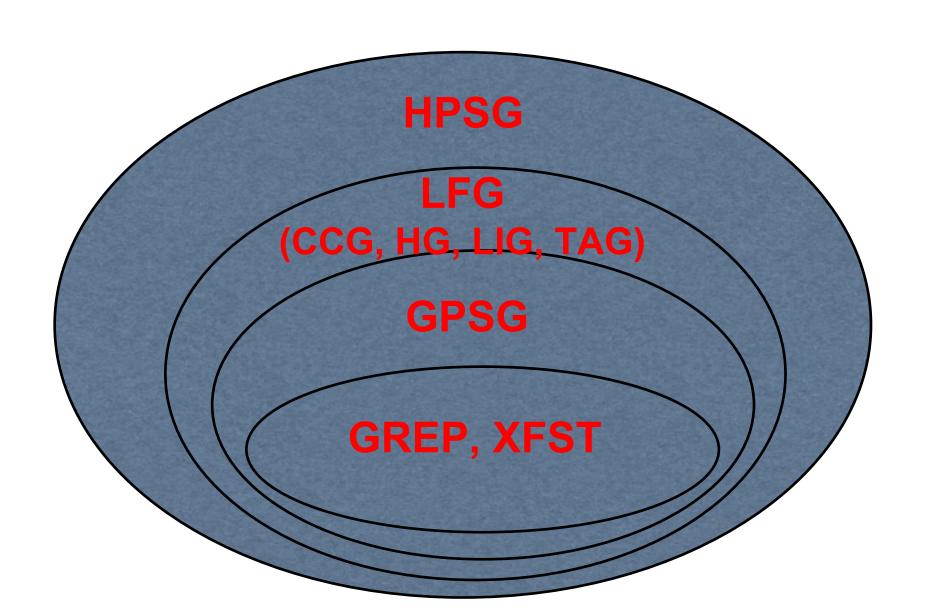
Les linguistes et linguistes-informaticiens ont donc conçus des formalismes plus adaptés à la description de gramaires.

Par exemple, les grammaires régulières utilisées en informatique :

[A-Z].*ations?

retrouver tous les fichiers dont le nom commence par une lettre majuscule non accentuée, suivie par une séquence de caractères quelconque, suivie par la séquence « ation », suivie par un « s » facultatif

Les formalismes utilisés en TAL

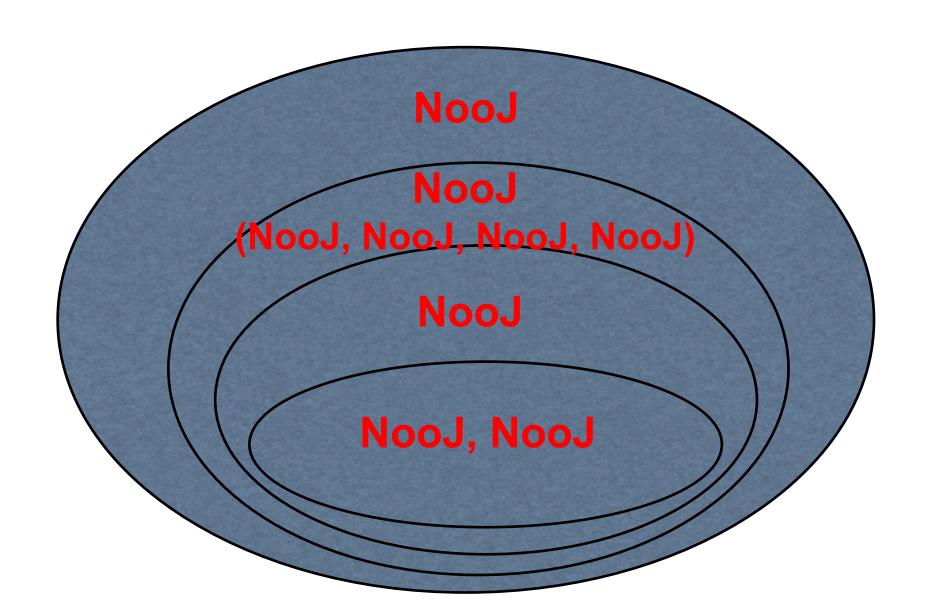


Différents formalismes : problèmes

- CCG, HG, IG, TAG: pas aussi puissants que LFG, mais les analyseurs sont plus efficaces: un compromis qui n'a pas forcément sa place en linguistique.
- Ces formalismes sont incompatibles entre eux: une règle XFST ne pourra pas être insérée dans une grammaire LFG
- Chaque grammaire doit donc être « totale », *i.e.*, doit couvrir tous les phénomènes linguistiques présents dans la phrase à représenter
- Chaque formalisme est très bien adapté à un type de phénomène linguistique... mais pas aux autres phénomènes

⇒ NooJ a été développé pour offrir les quatre types de grammaires dans un formalisme unifié.

Le formalisme de NooJ



2. Grammaires régulières

Grammaires génératives régulières (type 3) Une grammaire générative est régulière si :

• Les membres à gauche des règles de réécriture ne contiennent qu'un et un seul symbole auxiliaire

- Les membres à droite contiennent :
 - Soit \mathcal{E} , ex. $GN \rightarrow \mathcal{E}$
 - Soit une seule ALU, ex. GN → Luc
 - Soit une seule ALU suivie d'un seul symbole auxiliaire, ex. GN → le NOM

• La grammaire suivante est une grammaire régulière :

```
PHRASE → Jean GV
GV → voit GN
GN → un NOM
NOM → chat
NOM → chien
```

• La grammaire suivante est une grammaire régulière :

```
PHRASE → le SUITE
  PHRASE → un SUITE
SUITE → chat GVERBAL
SUITE → chien GVERBAL
GVERBAL → voit OBJET
   OBJET \rightarrow le NOM
   OBJET → un NOM
     NOM \rightarrow chat
     NOM \rightarrow chien
```

• La grammaire suivante n'est pas une grammaire régulière :

```
PHRASE \rightarrow GN voit GN
GN \rightarrow le NOM
GN \rightarrow un NOM
NOM \rightarrow chat
NOM \rightarrow chien
```

• La grammaire suivante n'est pas une grammaire régulière :

```
PHRASE \rightarrow GN voit GN
GN \rightarrow le NOM
GN \rightarrow un NOM
NOM \rightarrow chat
NOM \rightarrow chien
```

Grammaires régulières Formalisme alternatif

- Une grammaire régulière est définie à partir d'un alphabet, ex. : { a, b, c }
- n'importe quelle lettre constitue une grammaire régulière, ex. : b
- le mot vide est une grammaire régulière, ex. : ε
- Si X est une grammaire régulière, alors (X) est une grammaire régulière
- Si X est une grammaire régulière, alors X* est une grammaire régulière
- Si X et Y sont deux grammaires régulières, alors X Y est une grammaire régulière
- Si X et Y sont deux grammaires régulières, alors X | Y est une grammaire régulière

(lettres en noir, mots en bleu, grammaires en vert)

- Sur l'alphabet {a, b, c} :
- la grammaire régulière b représente le langage { b }, i.e., l'ensemble qui contient le seul mot b
- Attention: ne pas confondre la lettre a avec la grammaire a ou le mot a!
- la grammaire régulière ε représente le langage { ε }, i.e., l'ensemble qui contient le mot ε
- Attention : ne pas confondre le langage { € } avec le langage vide Ø
- la grammaire régulière (X) représente le même langage que celui représenté par la grammaire X, par ex., b = (b)

- Sur l'alphabet {a, b, c}:
- Concaténation : Si X et Y sont deux grammaires régulières, alors la grammaire X Y représente l'ensemble de tous les mots qui commencent par un mot du langage reconnu par X suivi par un mot du langage reconnu par Y

Exemples:

- la grammaire X = a reconnait le langage { a }
- la grammaire Y = c reconnait le langage { c }
- la grammaire X Y reconnait le langage { ac }
- la grammaire X Y X reconnait le langage { aca }
- la grammaire X b Y c reconnait le langage { abcc }

- Sur l'alphabet {a, b, c}:
- **Disjonction** : Si X et Y sont deux grammaires régulières, alors la grammaire X | Y représente l'ensemble de tous les mots du langage reconnu par X et aussi tous les mots du langage reconnu par Y

Exemples:

- la grammaire X = a reconnait le langage { a }
- la grammaire Y = c reconnait le langage { c }
- la grammaire X | Y reconnait le langage { a, c }
- la grammaire Y | X reconnait le langage { a, c }
- la grammaire X Y | c reconnait le langage { ac, c }

- Sur l'alphabet {a, b, c}:
- Opération de Kleene : Si X est une grammaire régulière, alors la grammaire X* représente l'ensemble de tous les mots que l'on peut contruire avec des mots du langage reconnu par X

Exemples:

- la grammaire a reconnait le langage { a }
- la grammaire a* reconnait le langage { ε, a , aa , aaa, aaaaa, aaaaa, ... }
- la grammaire (aba)* reconnait le langage { ε, aba, abaaba, abaabaaba, ... }
- la grammaire (a|b)* reconnait le langage { ε, a, b, ab, aa, bb, ba, aaaba, ... }

Utilité des parenthèses

(priorité de l'opérateur de concaténation par rapport à l'opérateur de disjonction)

• En arithmétique, on a :

$$2 \times 3 + 4 = 6 + 4 = 10$$
 \neq
 $2 \times (3 + 4) = 2 \times 7 = 14$

• En algèbre, on a :

$$2(a + b) = 2a + 2b \neq 2a + b$$

• En linguistique, on a de même :

la chaise | table

Sur l'alphabet { a, b, c }, construire la grammaire qui représente le langage qui contient :

• tous les mots de quatre lettres

Sur l'alphabet { a, b, c }, construire la grammaire qui représente le langage qui contient :

• tous les mots de quatre lettres

(a|b|c) (a|b|c) (a|b|c) (a|b|c)

Sur l'alphabet { a, b, c }, construire la grammaire qui représente le langage qui contient :

• tous les mots qui se terminent par la lettre « c »

Sur l'alphabet { a, b, c }, construire la grammaire qui représente le langage qui contient :

• tous les mots qui se terminent par la lettre « c » :

(a|b|c)*c

Sur l'alphabet { a, b, c }, construire la grammaire qui représente le langage qui contient :

• tous les mots qui ne contiennent pas la lettre « b »

Sur l'alphabet { a, b, c }, construire la grammaire qui représente le langage qui contient :

• tous les mots qui ne contiennent pas la lettre « b »

(a|c)*

Sur l'alphabet { a, b, c }, construire la grammaire qui représente le langage qui contient :

• tous les mots qui contiennent une lettre « b » et une seule

Sur l'alphabet { a, b, c }, construire la grammaire qui représente le langage qui contient :

• tous les mots qui contiennent une lettre « b » et une seule

(a|c)* b (a|c)*

Sur l'alphabet { a, b, c }, construire la grammaire qui représente le langage qui contient :

• tous les mots qui contiennent le mot « bac »

Sur l'alphabet { a, b, c }, construire la grammaire qui représente le langage qui contient :

• tous les mots qui contiennent le mot « bac »

 $(a|b|c)^*$ bac $(a|b|c)^*$

Sur l'alphabet { a, b, c }, construire la grammaire qui représente le langage qui contient :

 tous les mots qui commencent par un « a » et se terminent par un « b »

Sur l'alphabet { a, b, c }, construire la grammaire qui représente le langage qui contient :

 tous les mots qui commencent par un « a » et se terminent par un « b »

Grammaires régulières augmentées

- Le but n'est pas seulement de reconnaître des séquences : on veut aussi produire des résultats. Le résultat peut être une analyse linguistique, une traduction, une représentation sémantique, etc.
- On utilise deux alphabets : l'alphabet de lecture (input) et l'alphabet d'écriture (output)
- La grammaire met en relation les séquences lues (input) avec les séquences écrites (output)
- Exemples :
 - Remplacer dans un texte toutes les occurrences de « a » par des « b » : a/b
 - Mettre toutes les occurrences de « monday » en majuscule : monday/Monday
 - Traduire l'expression « quelle heure est-il ? » en anglais : quelle/what heure/time est/is-il/it ?/?

• Associer toutes les formes du nom *cousin* aux propriétés +m, +f, +s et +p

• Associer toutes les formes du nom cousin aux propriétés +m, +f, +s et +p

• Associer toutes les formes conjuguées du verbe *voler* au présent aux propriétés +1, +2, +3, +s et +p

 Associer toutes les formes conjuguées du verbe voler au présent aux propriétés +1, +2, +3, +s et +p

$$vol(e/+1+s | es/+2+s | e/+3+s | ons/+1+p | ez/+2+p | ent/+3+p)$$

• Associer toutes les formes conjuguées du verbe *voler* à l'imparfait aux propriétés +1, +2, +3, +s et +p

 Associer toutes les formes conjuguées du verbe voler à l'imparfait aux propriétés +1, +2, +3, +s et +p

vol (ais/
$$+1+s$$
 | ais/ $+2+s$ | ait/ $+3+s$ | ions/ $+1+p$ | iez/ $+2+p$ | aient/ $+3+p$)

• Associer toutes les formes conjuguées du verbe *voler* à l'imparfait et au conditionnel présent aux propriétés +1, +2, +3, +s et +p

```
vol (\varepsilon|er) (ais/+1+s | ais/+2+s | ait/+3+s | ions/+1+p | iez/+2+p | aient/+3+p )
```

• Associer toutes les variantes orthographiques et morphologiques du nom *tsar* aux propriétés +m, +f, +s et +p

 Associer toutes les variantes orthographiques et morphologiques du nom tsar aux propriétés +m, +f, +s et +p

```
(c|t) (s|z) ar (\varepsilon/+m | ine/+f) (\varepsilon/+s | s/+p)
```

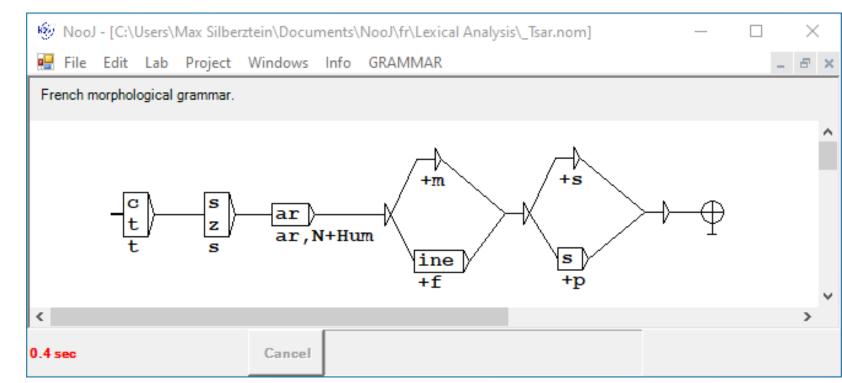
Grammaire régulières dans NooJ

• Notation dans NooJ (le mot vide ϵ s'écrit ici <E>) (c|t) (s|z) ar (<E>/+m | ine/+f) (<E>/+s | s/+p)

• Notation graphique équivalente dans NooJ:

Ce graphe reconnait et annote les 16 formes du mot *tsar*. Par ex., la forme *tsarines* sera annotée : tsar,N+Hum+f+p

(tsarines est un forme du mot tsar, c'est un nom humain féminin pluriel).



• Reconnaître toutes les variantes orthographiques du terme *Moyen-Âge*, et les normaliser

• Reconnaître toutes les variantes orthographiques du terme *Moyen-Âge* et les normaliser

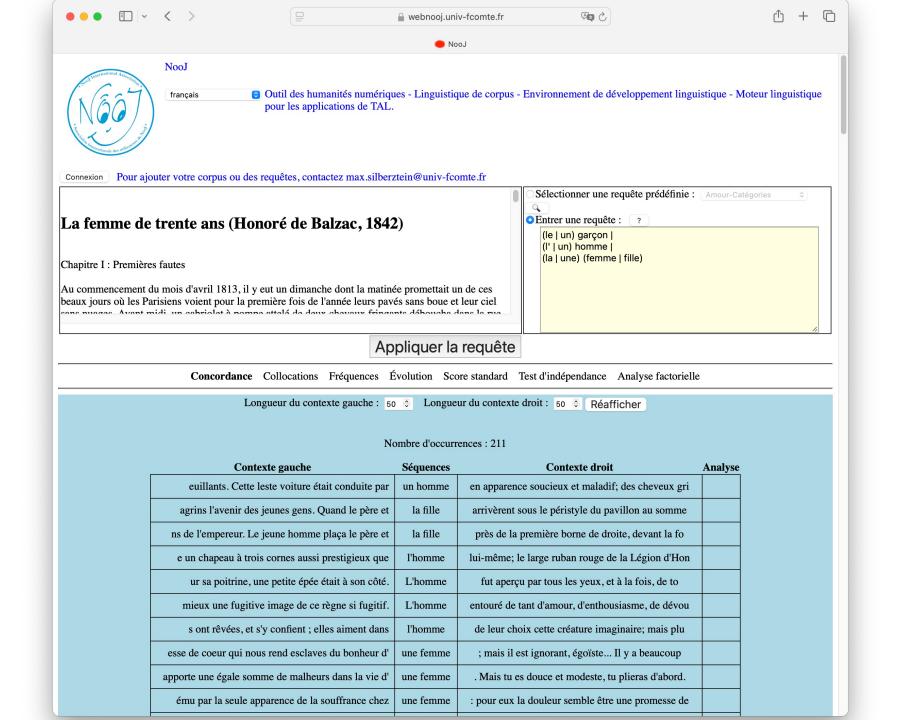
(M | m)/M oyen/oyen (- |
$$\epsilon$$
)/- (Â|A|â)/Â ge/ge | MOYEN/Moyen (- | ϵ)/- (Â|A)/Â GE/ge

• Reconnaître toutes les variantes orthographiques du terme *Moyen-Âge* et les normaliser

(M | m)/M oyen (- |
$$\epsilon$$
)/- (Â|A|â)/Â ge | MOYEN/Moyen (- | ϵ)/- (Â|A)/Â GE/ge = ((M | m) oyen (- | ϵ) (Â|A|â) ge | MOYEN (- | ϵ) (Â|A) GE/ge)/Moyen-Âge

Sur le vocabulaire français, construire

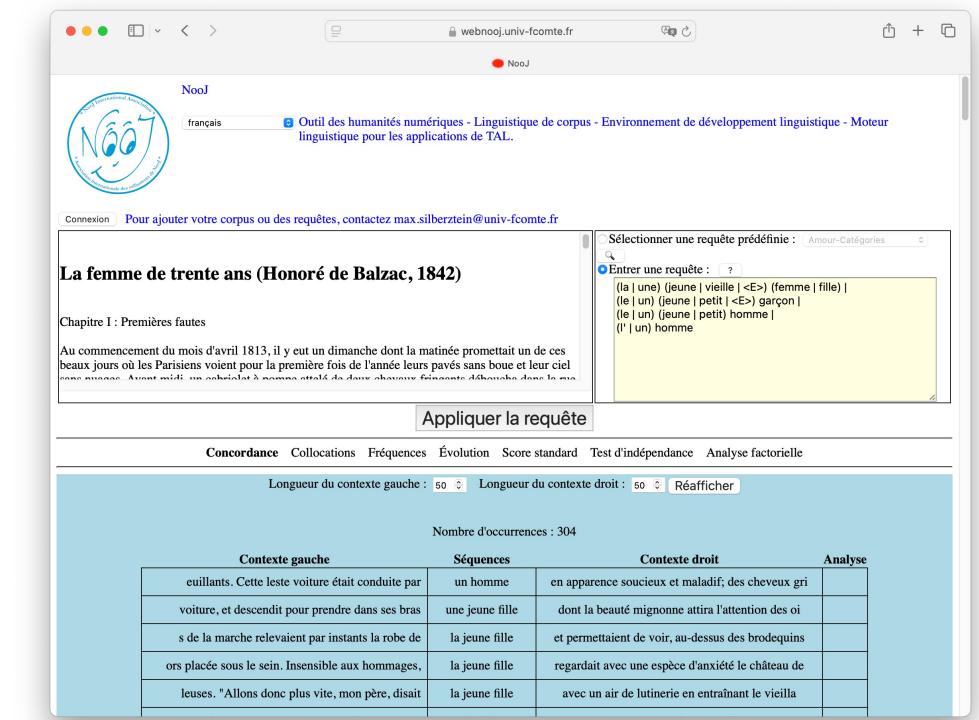
la grammaire qui reconnaît toutes les séquences constituées de l'article défini (*le* ou la) ou de l'article indéfini (*un* ou *une*) suivi par les mots *femme*, *fille*, *garçon* ou *homme*



Sur le vocabulaire français, construire la grammaire qui reconnaît toutes les séquences constituées :

- de l'article défini (*le* ou la) ou de l'article indéfini (*un* ou *une*) suivi par
- un adjectif facultatif jeune, vieille ou petit, puis par
- les mots femme, fille, garçon ou homme

la ou une,
jeune, vieille, &
fille ou femme
+
le, l' ou un,
petit, jeune, &
homme, garçon

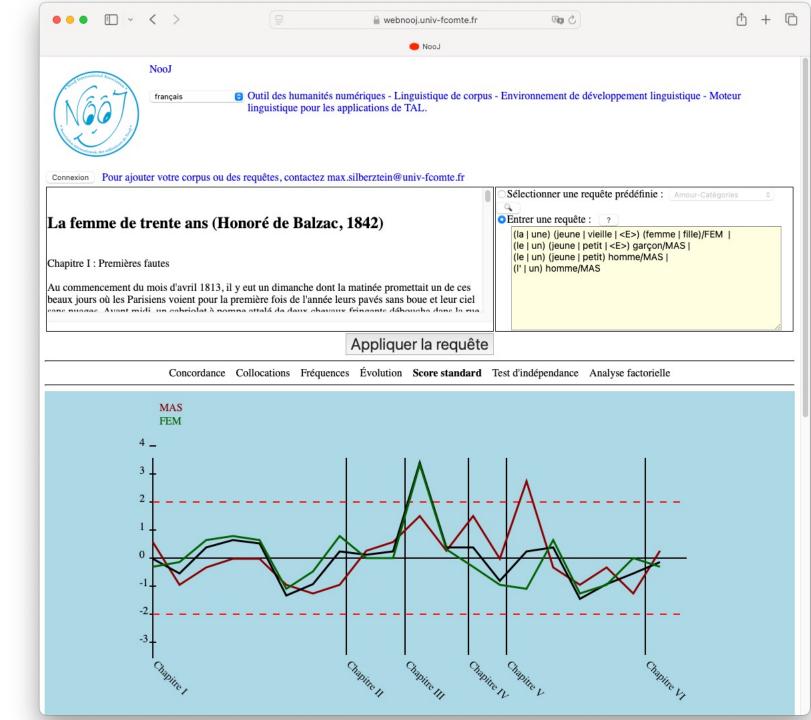


• Associer les occurrences de *femme* et *fille* avec l'étiquette FEM

• Associer les occurrences de *homme* et *garçon* avec l'étiquette MAS

• Commenter l'analyse statistique

Commenter l'analyse statistique



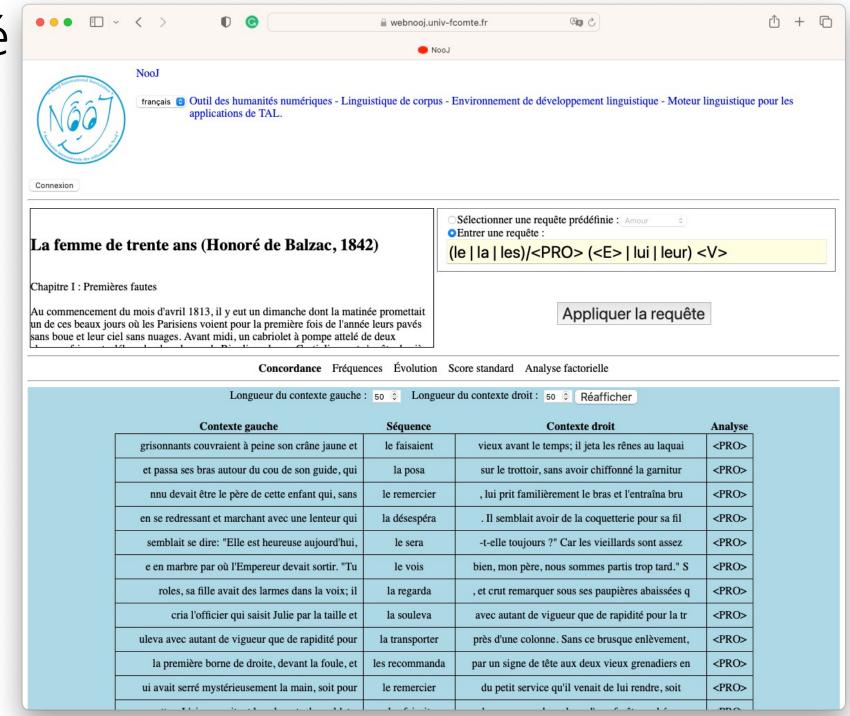
• Lever l'ambiguïté des mots *le, la* et *les* lorsqu'ils apparaissent avant un verbe, ex. *il la donne, elle ne les donne pas, ils le leur donne,* etc.

• Lever l'ambiguïté des mots *le, la* et *les* lorsqu'ils apparaissent avant un verbe, ex. *il la donne, elle ne les donne pas, ils le leur donne,* etc.

(le|la|les)/PRONOM (ε | lui | leur) VERBE

Levée d'ambiguïté automatique

- le mot vide ε est ici noté<E>
- les catégories sont écrites entre crochets, ex. <PRO> = pronom, <V> = verbe
- grammaire incomplète car elle reconnaît "la bise"



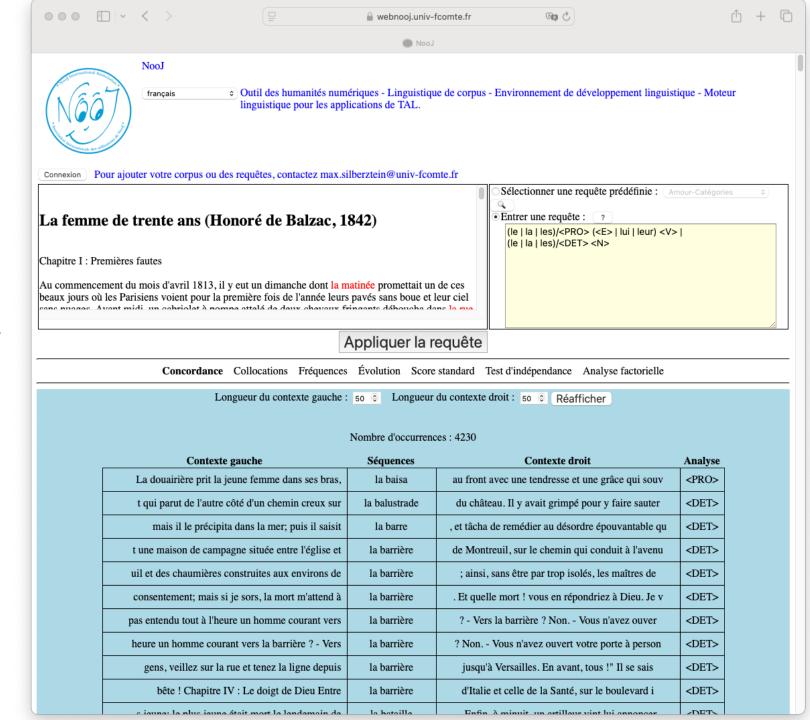
Levée d'ambiguïté automatique

"la bise" est reconnue deux fois, car "bise" peut être un verbe ou un nom.

On produit donc les deux analyses <PRO> et <DET> pour la même séquence :

cum en repos.		na m outeaux m permentor ve ne m perenes	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,
par intervalles,	la bise	sifflant à travers les branches noires des arbres	<pro></pro>
ıms au goût de	la bise	humide. Les oiseaux faisaient entendre leurs prol	<det></det>
anchaient avec	la blancheur	de son front pur. Elle avait même sur la lèvre su	<det></det>

Dans ce cas, on ne lève pas l'ambiguïté.



• Construire une grammaire pour lever l'ambiguïté des mots *par* et *pour* (préposition ou nom ?)

• Construire une grammaire pour lever l'ambiguïté des mots *par* et *pour* (préposition ou nom ?)

```
pour/PREPOSITION (\varepsilon|ne) (\varepsilon|jamais|pas|plus) (\varepsilon|me) (\varepsilon|le) (\varepsilon|leur|lui) VERBE+Infinitif | (par|pour)/PREPOSITION DETERMINANT | le/DETERMINANT pour/NOM | DETERMINANT+m+s par/NOM
```

• Construire une grammaire pour lever l'ambiguïté des mots *par* et *pour* (préposition ou nom ?)

```
pour/PREPOSITION (\varepsilon|ne) (\varepsilon|jamais|pas|plus) (\varepsilon|me) (\varepsilon|le) (\varepsilon|leur|lui) VERBE+Infinitif | (par|pour)/PREPOSITION DETERMINANT | le/DETERMINANT pour/NOM | DETERMINANT+m+s par/NOM
```

Notation NooJ:

```
pour/<PREP> (<E>|<ne>) (<E>|jamais|pas|plus) (<E>|<me>) (<E>|<le>) (<E>|leur|lui)
<V+INF> |
  (par|pour) /<PREP> <DET> |
le/<DET> pour/<N> |
  <DET+m+s> par/<N>
```

Levée d'ambiguïté morpho-syntaxique automatique

• Des dizaines de grammaires de levée d'ambiguïté peuvent être développées autour de mots fréquents comme :

```
• est : verbe, nom ou adjectif?
```

- *si* : conjonction ou nom ?
- place : nom ou verbe ?
- fait : nom ou verbe?
- *a* : nom ou verbe ?
- son, ton : déterminants ou noms ?
- plus : verbe ou adverbe ?
- rue : verbe ou nom ?
- complet : adjectif ou nom ?
- première : adjectif ou nom ?
- ...

• Construire la grammaire qui reconnait toutes les séquences de mots qui peuvent apparaître entre "il" et la forme verbale "donnera"

 Construire la grammaire qui reconnait toutes les séquences de mots qui peuvent apparaître entre "il" et la forme verbale "donnera"

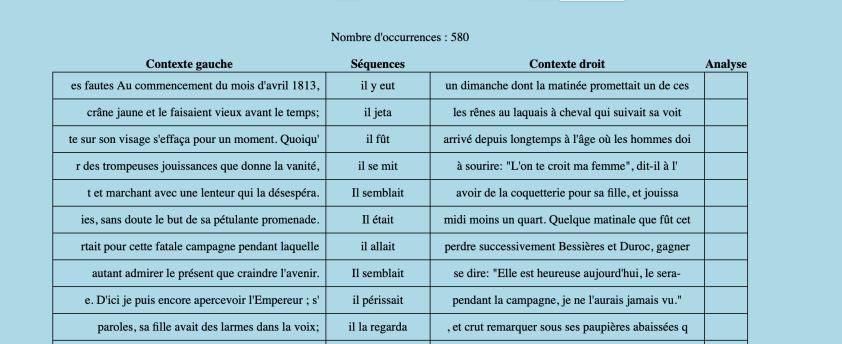
```
(ne \mid \epsilon) (le \mid la \mid les \mid \epsilon) (lui \mid leur \mid \epsilon) \mid
(ne \mid \epsilon) (me \mid te \mid se \mid nous \mid vous \mid leur) (le \mid la \mid les \mid \epsilon) \mid
        (ne | ε) (m' | t' | s' | nous | vous | leur) (en | y) |
                                     (n' | \epsilon) (en | y)
                                          donnera
```

• Généraliser à tous les verbes ; appliquer la grammaire au texte



Concordance Collocations Fréquences Évolution Score standard Test d'indépendance Analyse factorielle

Longueur du contexte gauche : 50 🗦 Longueur du contexte droit : 50 🖨 Réafficher



3. Grammaires hors-contexte

Une grammaire générative est hors-contexte si :

 Les membres à gauche des règles de réécriture ne contiennent qu'un et un seul symbole auxiliaire

 Les membres à droite contiennent n'importe quelle séquence d'ALU, de symboles et/ou de mots vides, ex. : GN → DET NOM de NOM

• La grammaire suivante est une grammaire hors contexte :

```
PHRASE → GN voit GN
GN → le NOM
GN → un NOM
NOM → chat
NOM → chien
```

• La grammaire suivante n'est pas une grammaire hors contexte :

```
PHRASE → GN voit GN
GN Singulier → le NOM
GN Pluriel → les NOM
NOM → chat
NOM → chats
```

• La grammaire suivante n'est pas une grammaire hors contexte :

```
PHRASE → GN voit GN
GN Singulier → le NOM
GN Pluriel → les NOM
NOW → chat
NOM → chats
```

Les grammaires hors contexte, formalismes alternatifs (ici : notation Backus-Naur)

• Le membre droit des règles de réécriture peut contenir des disjonctions « | »:

• Une grammaire hors-contexte contient donc un ensemble de règles nommées

Les grammaires hors contexte, formalismes alternatifs (ici : notation simplifiée)

• Le membre droit peut contenir des expressions avec disjonctions « | », des parenthèses, des opérateurs de Kleene « * » et des mots vides ϵ :

```
Unités = I | II | III | IV | V | VI | VII | VIII | IX

Dizaines = X | XX | XXX | XL | L | LX | LXX | LXXX | XC

ChiffreRomain = Dizaines (Unités | ε) | Unités
```

Compléter la grammaire des chiffres romains :

```
Unités = I | II | III | IV | V | VI | VII | VIII | IX
Dizaines = X | XX | XXX | XL | L | LX | LXX | LXXX | XC
Centaines = ...
Milliers = ...
ChiffreRomain = ...
```

• Compléter la grammaire des chiffres romains :

```
Unités = I | II | III | IV | V | VI | VII | VIII | IX

Dizaines = X | XX | XXX | XL | L | LX | LXX | LXXX | XC

Centaines = C | CC | CCC | CD | D | DC | DCC | DCCC | CM

Milliers = M | MM | MMM

ChiffreRomain = (\epsilon|Milliers) (\epsilon|Centaines) (\epsilon|Dizaines) (\epsilon|Unités)
```

Première tentative : les chiffres des milliers, des centaines, des dizaines et des unités peuvent être absents, ex. « MMMII »

• Compléter la grammaire des chiffres romains :

```
Unités = I | II | III | IV | V | VI | VII | VIII | IX

Dizaines = X | XX | XXX | XL | L | LX | LXX | LXXX | XC

Centaines = C | CC | CCC | CD | D | DC | DCC | DCCC | CM

Milliers = M | MM | MMM

ChiffreRomain = (\varepsilon|Milliers) (\varepsilon|Centaines) (\varepsilon|Dizaines) (\varepsilon|Unités)
```

PROBLÈME : cette grammaire reconnaît le mot vide :(

• Compléter la grammaire des chiffres romains :

```
Unités = I | II | III | IV | V | VI | VII | VIII | IX
Dizaines = X | XX | XXX | XL | L | LX | LXX | LXXX | XC
Centaines = C | CC | CCC | CD | D | DC | DCC | DCCC | CM
Milliers = M | MM | MMM
ChiffreRomain = (\varepsilon | Milliers) (\varepsilon | Centaines) (\varepsilon | Dizaines) Unités |
       (ε|Milliers) (ε|Centaines) Dizaines |
       (ε|Milliers) Centaines
       Milliers
```

Exercice: Grammaire hors-contexte augmentée

• En plus de reconnaître les chiffres romains, produire leur valeur en chiffres arabes.

```
Unités = I | II | III | IV | V | VI | VII | VIII | IX
Dizaines = X | XX | XXX | XL | L | LX | LXX | LXXX | XC
Centaines = C | CC | CCC | CD | D | DC | DCC | DCCC | CM
Milliers = M | MM | MMM
<u>ChiffreRomain</u> = (\varepsilon | Milliers) (\varepsilon | Centaines) (\varepsilon | Dizaines) Unités |
       (ε|Milliers) (ε|Centaines) Dizaines |
       (ε Milliers) Centaines
       Milliers
```

Exercice: Grammaires hors-contexte augmentées

• En plus de reconnaître les chiffres romains, produire leur valeur en chiffres arabes.

```
Unités = I/1 | II/2 | III/3 | IV/4 | V/5 | VI/6 | VII/7 | VIII/8 | IX/9
Dizaines = X/1 | XX/2 | XXX/3 | XL/4 | L/5 | LX/6 | LXX/7 | LXXX/8 | XC/9
Centaines = C/1 \mid CC/2 \mid CCC/3 \mid CD/4 \mid D/5 \mid DC/6 \mid DCC/7 \mid DCCC/8 \mid CM/9
Milliers = M/1 \mid MM/2 \mid MMM/3
ChiffreRomain = \varepsilon/VAL=
                    Milliers (\varepsilon/0|Centaines) (\varepsilon/0|Dizaines) (\varepsilon/0|Unités)
                    Centaines (\varepsilon/0|Dizaines) (\varepsilon/0|Unités)
                    Dizaines (ε/0 Unités)
                    Unités
```

Grammaires hors-contexte, notation de NooJ

Notation NooJ:

- utilise la notation simplifiée
- le mot vide ε est ici noté <E>
- les symboles auxiliaires sont préfixés par le caractère « : »
- chaque règle doit être terminée par un point virgule « ; »
- la règle principale doit s'appeler Main

```
NooJ - [C:\Users\Max Silberztein\Documents\NooJ\fr\Lexical Analysis\Chiffres romains.nom]
🖳 File Edit Lab Project Windows Info TEXT
                                                                                                        _ & ×
# NooJ V7
# Morphological grammar
  Language is: fr
  Special Start Rule: Main
# Special Characters: '=' '<' '>' '\' '"' ':' '|' '+' '-' '/' '$' ' ' ':' '#'
# Max Silberztein
Unités = I/1 | II/2 | III/3 | IV/4 | V/5 | VI/6 | VII/7 | VIII/8 | IX/9 ;
Dizaines = X/1 \mid XX/2 \mid XXX/3
                                           I \times L/4
|Centaines = C/1|
                         CC/2 \mid CCC/3 \mid CD/4 \mid D/5 \mid DC/6 \mid
Milliers = M/1
                    | MM/2 | MMM/3 ;
Main = \langle E \rangle / CR + VAL =
          :Milliers (\langle E \rangle / 0 \mid :Centaines) (\langle E \rangle / 0 \mid :Dizaines) (\langle E \rangle / 0 \mid :Unités)
          :Centaines (\langle E \rangle / 0 | :Dizaines) (\langle E \rangle / 0 |
                                                               :Unités)
          :Dizaines (<E>/0 | :Unités)
          :Unités ;
                       Cancel
```

• Réécrire la grammaire qui reconnait toutes les séquences de particules préverbales sous forme d'une grammaire hors-contexte

 Réécrire la grammaire qui reconnait toutes les séquences de particules préverbales sous forme d'une grammaire hors-contexte

```
Main = il PPV VERBE
PPV = (ne | \epsilon) (LE | \epsilon) (LUI | \epsilon) |
         (ne |\epsilon| (ME |\epsilon| (LE |\epsilon| )
         (ne |\epsilon| (M |\epsilon|) (y |\epsilon|) (en |\epsilon|)
         (n' \mid \epsilon) (y \mid \epsilon) (en \mid \epsilon)
LE = le | la | les
LUI = lui | leur
ME = me | te | se | nous | vous | leur
M = m' | t' | s' | nous | vous | leur
```

Construire une grammaire des groupes nominaux féminins singuliers que l'on peut construire avec des déterminants, des noms, des adjectifs.

<u>GNFémininSingulier</u> =

Construire une grammaire des groupes nominaux féminins singuliers que l'on peut construire avec des déterminants, des noms, des adjectifs.

```
ModifD = très* ADJ | de NOM

ADJ = belle | grande | jolie | petite | ...

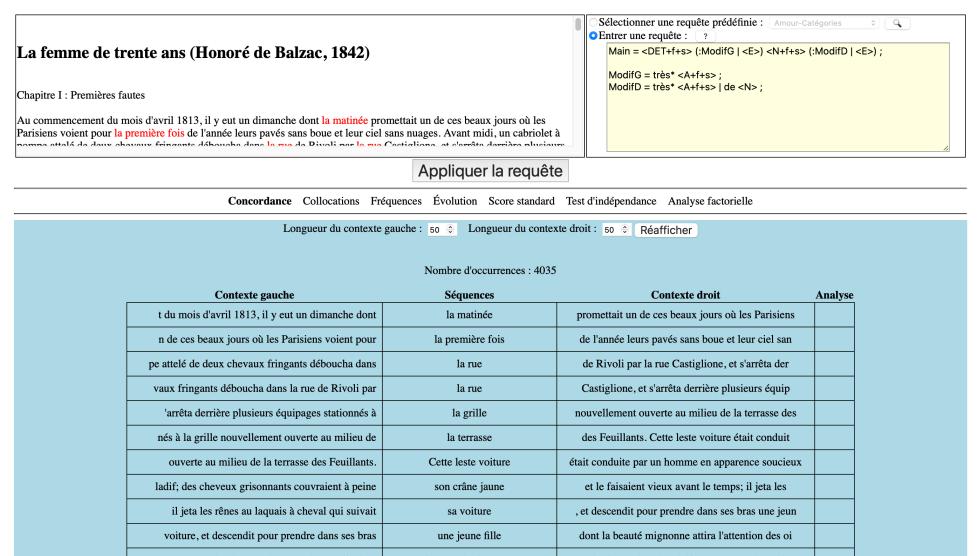
DET = cette | la | ma | ta | sa | notre | votre | leur | une | ...

NOM = chaise | fleur | maison | table | ...
```

<u>GNFémininSingulier</u> = DET (ε | ModifG) NOM (ε | ModifD)

ModifG = très* ADJ

Construire une grammaire des groupes nominaux féminins singuliers que l'on peut construire avec des déterminants, des noms, des adjectifs.



Construire une grammaire qui permet de reconnaître les noms de personne, (soit les hommes, soit les femmes), comme par ex. :

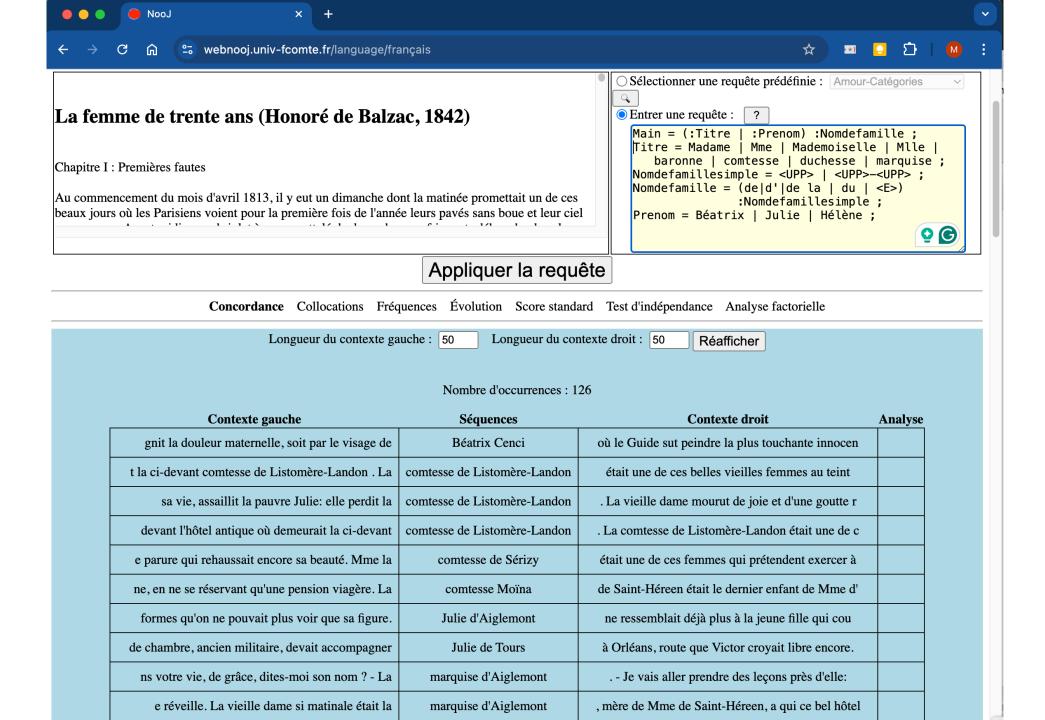
*Mme Jeanne Dupont-Durand, Mlle de la Castellerie, Jeanne d'Artagnan*On devra utiliser les symboles auxiliaires Prénom, Titre, Nomdefamille

```
Personne = ...

Titre = ...

Prénom = ...

Nomdefamille = ...
```



Construire une grammaire qui permet de reconnaître les dates en anglais, comme par ex. :

```
Monday; Tuesday, June 14th; Wednesday, July 1st; on January the 3rd
NameOfDay = Monday | ...
Number 29 = (2|\epsilon)(1st|2nd|3rd) | 20st | (10|11|12|13)th | ((1|2|\epsilon)(4|5|6|7|8|9)th
Number 30 = ...
Number31 = ...
Month29 = February
Month30 = April | ...
Month31 = January | ...
```

Attention à ne pas reconnaître : Monday, June 0th ; February 31st ; 21st

```
Day = Monday|Tuesday|Wednesday|Thursday|Friday|Saturday|Sunday
```

```
Month29 = February
```

Month30 = April | June | September | November

Month31 = January|March|May|July|August|October|December

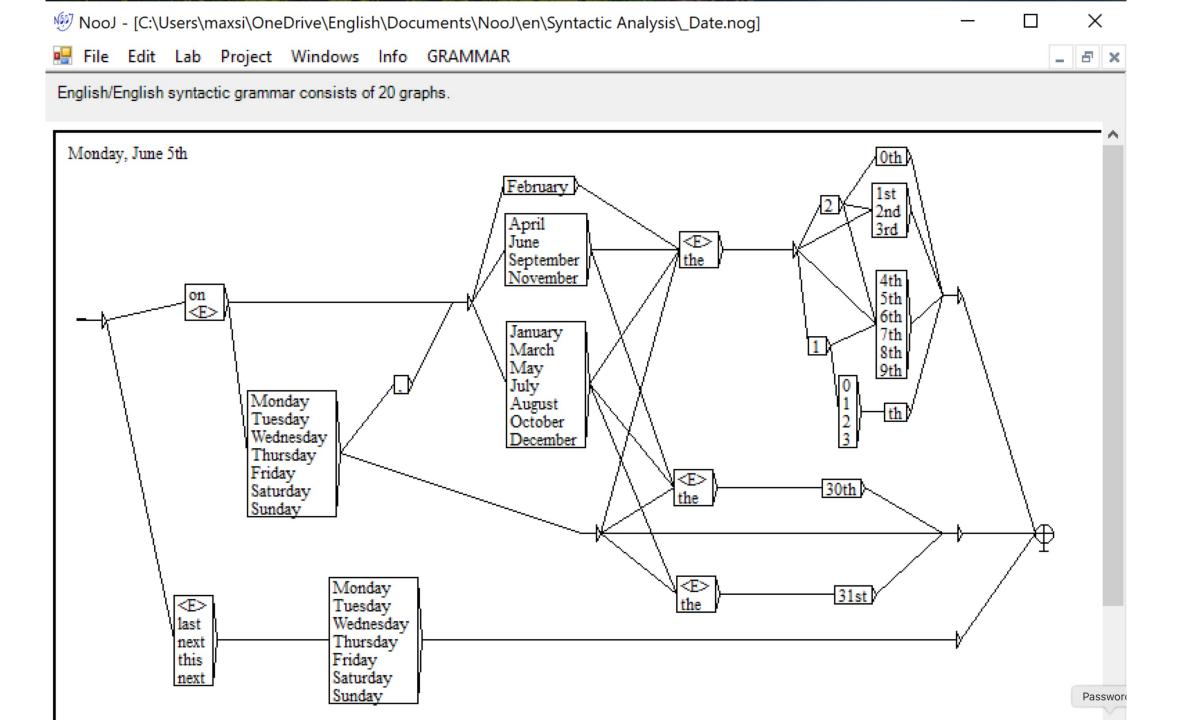
Number 29 = $(2|\epsilon)(1st|2nd|3rd) | 20st | (10|11|12|13)th | ((1|2|\epsilon)(4|5|6|7|8|9)th)$

Number30 = Number29 | 30st

Number31 = Number30 | 31st

<u>Date</u> = ...

```
Day = Monday|Tuesday|Wednesday|Thursday|Friday|Saturday|Sunday
Month29 = February
Month30 = April | June | September | November
Month31 = January | March | May | July | August | October | December
Number 29 = (2|\epsilon)(1st|2nd|3rd) | 20st | (10|11|12|13)th | ((1|2|\epsilon)(4|5|6|7|8|9)th
Number30 = Number29 | 30st
Number31 = Number30 | 31st
\underline{\text{Date}} = (\text{ on } | (\text{on } | \epsilon) \underline{\text{Day }}'','')
       ( Month 29 (the |\epsilon| Number 29 |
          Month 30 (the |\epsilon| Number 30 |
          Month31 (the |\epsilon| Number31)
       (last|next|this) Day
```



Construire une grammaire qui permet de reconnaître les heures en anglais, comme par ex. :

a quarter to three; four a.m.; half past noon; three o'clock in the afternoon

```
Hours = ...

HoursAfternoon = ...

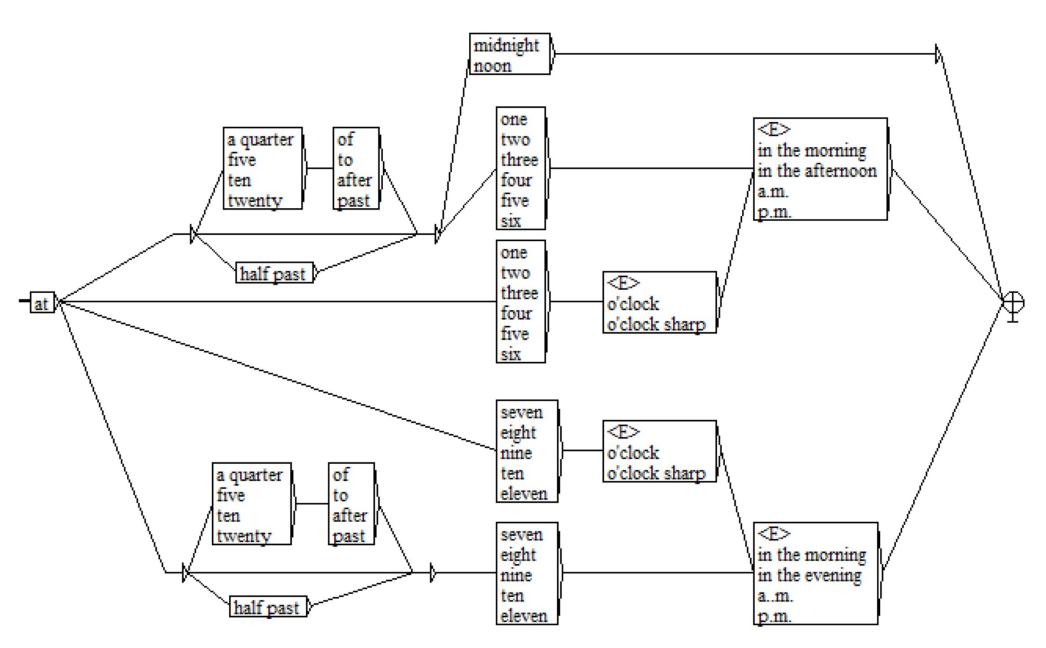
HoursEvening = ...

Minutes = ...

Time = ...
```

Attention à ne pas reconnaître : Two in the evening ; Midnight p.m. ; A quarter past two o'clock

```
Minutes = (five | ten | a quarter | twenty | twenty-five) (past | to) | half past
HoursAfternoon = one | two | three | four | five
HoursEvening = six | seven | eight | nine | ten | eleven
```



Construire une grammaire qui combine les dates et les heures, par ex. :

Monday at a quarter to three; On Sunday, January 4th at four a.m.; at half past noon; The 4th of February at three o'clock in the afternoon

Construire une grammaire qui combine les dates et les heures, par ex. :

Monday at a quarter to three; On Sunday, January 4th at four a.m.; half past noon; The 4th of February at three o'clock in the afternoon

Complement = Date (at Hours $| \varepsilon |$ at Hours

Récursivité

• Dans une grammaire hors-contexte, l'expression (à droite) peut contenir des symboles auxiliaires

• Que se passe-t-il lorsqu'une règle est définie à partir d'elle-même, ex. :

```
Phrase = Sujet Verbe (ε | que Phrase)
Sujet = Jean | Marie | Paul | Hélène
Verbe = dit | espère | pense | dort
```

Récursivité

• Dans une grammaire hors-contexte, l'expression (à gauche) peut contenir des symboles auxiliaires

• Que se passe-t-il lorsqu'une règle est définie à partir d'elle-même, ex. :

```
Phrase + Sujet Verbe (ε | que Phrase)
Sujet = Jean | Marie | Paul | Hélène
Verbe = dit | espère | pense | dort
```

Jean dit que Marie espère que Paul pense que Hélène dort

Trois sortes de récursivité

- Récursivité à droite
 Phrase = GN (pense | dit | espère | voit) GN
 GN = son voisin | que Phrase
- Récursivité à gauche GN = GN (rouge | de Paul | de campagne | en Bretagne) | la maison
- Récursivité générale
 <u>Phrase</u> = <u>DET NOM</u> que <u>Phrase</u> <u>VERBE</u>

Supprimer la récursivité à gauche

• Récursivité à gauche :

GN = GN (rouge | de Paul | de campagne | en Bretagne) | la maison

• On peut supprimer les récursivités à gauche et construire la grammaire régulière équivalente :

GN = la maison (rouge | de Paul | de campagne | en Bretagne)*

GN = GN (rouge | de Paul | de campagne | en Bretagne) | la maison

GN => GN en Bretagne => GN de campagne en Bretagne => GN rouge de campagne en Bretagne => la maison rouge de campagne en Bretagne

Supprimer la récursivité à droite

• Récursivité à droite :

```
Phrase = GN (pense | dit | espère | voit) GN
GN = son voisin | que Phrase
```

• On peut supprimer les récursivités à droite et construire la grammaire régulière équivalente :

```
Phrase = (son voisin (pense | dit | espère | voit) que)*
son voisin (pense | dit | espère | voit) son voisin
```

Supprimer la récursivité générale

• Récursivité générale (ni à gauche, ni à droite) :

```
Phrase = DET NOM (& | Relative) Verbe
Relative = que Phrase
Verbe = dort | aime | a vu | a embauchée
```

• Impossible de supprimer cette récursivité...

Supprimer la récursivité générale

• Récursivité générale :

```
Phrase = DET NOM (ε | Relative) Verbe
Relative = que Phrase
Verbe = dort | aime | a vu | a embauchée
```

Mais y a-t-il vraiment un besoin linguistique ?

Le chat dort. Le chat que la cousine a vu dort

? Le chat que la voisine que le voisin aime a vu dort

* Le chat que la voisine que le voisin que la police recherche aime a vu dort

Supprimer la récursivité générale

• Récursivité générale :

```
Phrase = <DET> <N> (ε | Relative) Verbe
Relative = que Phrase
Verbe = dort | aime | a vu | a embauchée
```

• Y a-t-il vraiment un besoin linguistique?

Le chat dort. Le chat que la cousine a vu dort

? Le chat que la voisine que le voisin aime a vu dort

* Le chat que la voisine que le voisin que la police recherche aime a vu dort

• On ne peut pas supprimer les récursivités générales, mais ce n'est pas trop grave... Et on peut toujours fixer une limite aux imbrications de relatives, ex. 3 maximum.

Supprimer les récursivités

- Il vaut mieux développer des grammaires hors-contexte que des grammaires régulières pour nommer, gérer et combiner des grands nombres de règles
- Les algorithmes traitant des grammaires régulières sont bien plus rapides que les algorithmes traitant des grammaires hors-contexte
- Certains outils informatiques (ex. NooJ) peuvent supprimer les récursivités à gauche et à droite, et transformer les grammaires hors-contexte en grammaires régulières, et ainsi analyser les textes avec de bonnes performances.

4. Grammaires contextuelles

Grammaires génératives contextuelles (type 1)

• Les règles de réécriture peuvent contenir 2 symboles à gauche ; le premier doit être présent aussi dans le second membre, ex. :

```
SINGULIER PHRASE \rightarrow SINGULIER GN voit GN
PLURIFI PHRASE → PLURIFI GN voient GN
GN → PLURIEL GN | SINGULIER GN
SINGULIER GN → SINGULIER IN NOMSINGULIER
PLURIEL GN → PLURIEL les NOMPLURIEL
SINGULIER \rightarrow \epsilon
PLURIEL \rightarrow \varepsilon
NOMSINGULIER \rightarrow chat
NOMSINGULIER \rightarrow chien
NOMPLURIEL \rightarrow chats
NOMPLURIEL \rightarrow chiens
```

Grammaires génératives contextuelles (Type 1)

• Les règles de réécriture peuvent contenir 2 symboles à gauche ; le premier doit être présent aussi dans le second membre, ex. :

pour marquer les contextes.

Langage décrit : { le chat voit le chat, le chat voit le chien, le chat voit les chats, le chat voit les chiens, le chien voit le chien voit le chien voit le chien, le chien voit les chats, le chien voit les chiens, les chats voient le chat, les chats voient le chien, les chats voient les chiens voient le chat, les chiens voient le chien, les chiens voient les chiens les chiens voient les chiens les chiens voient les chiens les chiens voient les chiens voient les chiens voient les chiens l

Les séquences *les chats voit le chien* et *les chat voit les chien* ne sont pas reconnues

```
SINGULIER PHRASE \rightarrow SINGULIER ON voit GN

PLURIEL PHRASE \rightarrow PLURIEL GN voient GN

GN \rightarrow PLURIEL GN | SINGULIER GN
```

SINGULIER GN → SINGULIER le NOMSINGULIER

PLURIEL GN - PLURIEL les NOMPLURIEL

SINGULIER $\rightarrow \epsilon$

PLURIEL $\rightarrow \varepsilon$

NOMSINGULIER \rightarrow chat

NOMSINGULIER → chien

NOMPLURIEL \rightarrow chats

NOMPLURIEL → *chiens*

Les symboles SINGULIER et PLURIEL sont utilisés

Grammaires contextuelles (Type 1), notation alternative

(ici: notation Lexical-Functional Grammar)

• Exemple de règle de grammaire LFG :

```
VP --> V: \uparrow = \downarrow;

(NP: (\uparrowOBJ)=\downarrow

(\downarrowCASE)=ACC)

PP*: \downarrow \in (\uparrowADJUNCT).
```

- Le constituant VP contient un verbe V, un objet optionnel (NP) et un nombre quelconque de compléments prépositionnels PP*
- « ↑=↓ » signifie que le constituant VP hérite des propriétés du verbe V
- « (↑OBJ)=↓ » signifie que le constituant NP hérite des propriétés de OBJ
- « (↓CASE)=ACC » signifie que le NP doit être à l'accusatif
- « ↓∈(↑ADJUNCT) » signifie que les compléments PP appartiennent à l'ensemble des modifieurs (« ADJUNCT » en anglais) du VP

Grammaires contextuelles (Type 1), notation alternative (ici : notation simplifiée)

• Exemple de grammaire hors contexte :

```
GNSINGULIER = DETSINGULIER (\varepsilon | ADJSINGULIER) NOMSINGULIER

GNPLURIEL = DETPLURIEL (\varepsilon | ADJPLURIEL) NOMPLURIEL

GN = GNSINGULIER | GNPLURIEL
```

• La grammaire contextuelle correspondante est beaucoup plus simple si elle ne contient pas les contraintes d'accord en nombre :

$$GN = DET (\epsilon \mid ADJ) NOM$$

• On va donc ajouter à cette grammaire hors-contexte les contraintes d'accord

Grammaires contextuelles (Type 1), notation alternative (ici : notation simplifiée)

Grammaire contextuelle :

```
\underline{GN} = DET
(\varepsilon \mid ADJ)
NOM
```

• Grammaire contextuelle = Grammaire hors contexte + on ajoute les contraintes d'accord

Grammaires contextuelles (Type 1), notation alternative (ici : notation simplifiée)

Grammaire contextuelle :

```
GN = DET/<THIS$Nombre=N$Nombre>
  (ε | ADJ/<THIS$Nombre=N$Nombre>)
  (N NOM)
```

- Grammaire contextuelle = Grammaire hors contexte + on ajoute les contraintes d'accord en production (outputs) :
- → des variables contiennent des ALU : THIS (ALU courante), N
- des propriétés morphologiques et lexicales : \$Nombre, \$Genre, \$Cas, \$Temps, \$Classedistributionnelle, etc.
- ➤ des contraintes en production (outputs) : <THIS\$Nombre=N\$Nombre>

Exemples de grammaires contextuelles morphologiques (ici : notation simplifiée)

• Reconnaître des verbres dérivés en adjectifs en *-able* :

$$L = a \mid b \mid c \mid ... \mid z$$

$$ADJECTIF = \begin{pmatrix} Prefixe \mid L \mid L^* \end{pmatrix} able / \langle Prefixe \# er = : VERBE \rangle$$

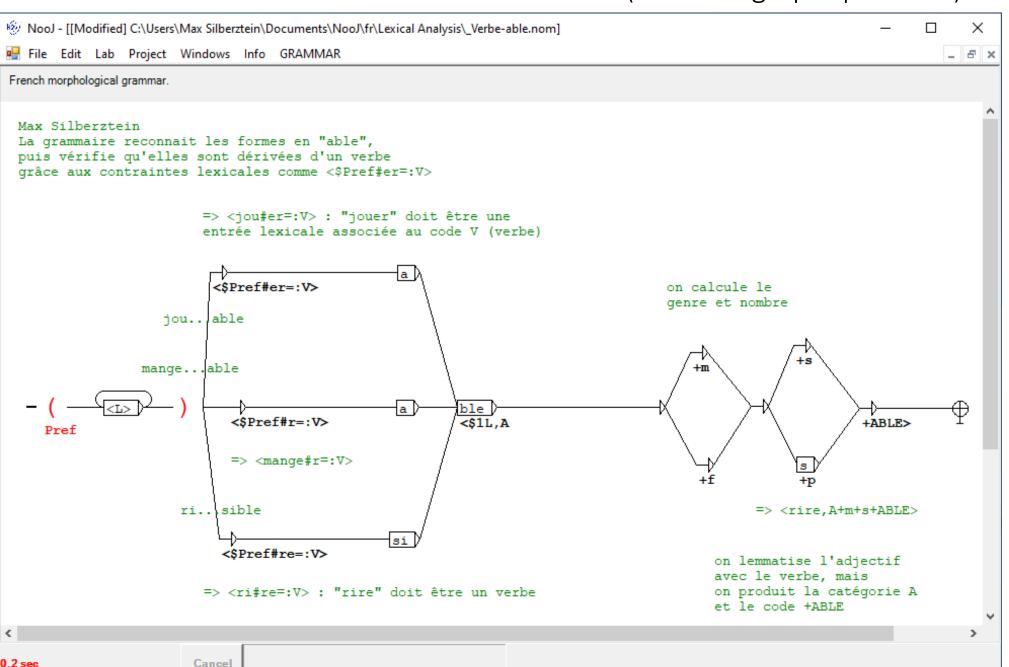
$$VERBE = aimer \mid brancher \mid casser \mid donner \mid écouter \mid voler \mid ...$$

Exemples de formes reconnues : aimable, branchable, cassable La forme table n'est pas reconnue car « ter » n'est pas un verbe

Exemples de grammaires contextuelles morphologiques (notation NooJ)

- Le symbole Main est le symbole auxiliaire initial
- Le symbole <L> représente n'importe quelle lettre ; <L><L>* représente toutes les séquences de lettres ; ces séquences sont rangées dans la variable \$Prefixe.
- La contrainte <\$Prefixe#er=:V> vérifie que si l'on concatène la suite de lettres rangée dans la variable \$Prefixe avec le « er », on obtient un verbe.

Verbes dérivés en -able (notation graphique NooJ)



- On range toutes les lettres (<L>) du préfixe avant « able » dans la variable \$Pref
- contraintes, ex. :
 <\$Pref#er=:V>
 (on ajoute « er » après
 \$Pref ; on vérifie que
 la forme obtenue est

- On vérifie les

un verbe)

- On produit le résultat d'analyse, par exemple, les deux annotations pour jouable :
- <jouer,A+m+s+ABLE>
 <jouer,A+f+s+ABLE>

Exemples de grammaires contextuelles morphologiques

• Reconnaître des verbes préfixés en re- et dé- :

```
VERBEPREF = (r | re | ré | (\varepsilon|re) (dé|dés)) (suffixe L L*)/<Suffixe=:VERBE>
```

Exemples de formes reconnues : rappeler, remanger, réapprendre, redémonter Le verbe rater n'est pas reconnu car « ater » n'est pas un verbe

Réduplications

- Mots comme bébé, chouchou, dodo, mémé, papa
- Quelques expressions, ex.: C'est pas joli joli
- En ukrainien : *Elle est belle belle* [très belle]. *Il marcha marcha marcha marcha pour arriver à destination* [marcha longtemps]
- Les réduplications marquent un pluriel (indonésien), une emphase (tagalog), un ensemble d'objets identiques (arbre arbre [forêt] en quechua), un pluriel distributif (ex. en japonais : tous les mois = つきづき), etc.

Reconnaître les réduplications

L'accord en français (grammaire régulière)

```
<DET+m+s> (<A+m+s> | <E>) <N+m+s> (<A+m+s> | <E>) |
<DET+f+s> (<A+f+s> | <E>) <N+f+s> (<A+f+s> | <E>) |
<DET+m+p> (<A+m+p> | <E>) <N+m+p> (<A+m+p> | <E>) |
<DET+f+p> (<A+f+p> | <E>) <N+f+p> (<A+f+p> | <E>)
```

L'accord en français (grammaire hors-contexte)

GN = GNMS | GNFS | GNMP | GNFP

$$GNMS = \langle DET+m+s \rangle (\langle A+m+s \rangle | \langle E \rangle) \langle N+m+s \rangle (\langle A+m+s \rangle | \langle E \rangle)$$

$$GNFS = \langle DET + f + s \rangle (\langle A + f + s \rangle | \langle E \rangle) \langle N + f + s \rangle (\langle A + f + s \rangle | \langle E \rangle)$$

$$GNMP = \langle DET + m + p \rangle (\langle A + m + p \rangle | \langle E \rangle) \langle N + m + p \rangle (\langle A + m + p \rangle | \langle E \rangle)$$

$$GNFP = \langle DET+f+p \rangle (\langle A+f+p \rangle | \langle E \rangle) \langle N+f+p \rangle (\langle A+f+p \rangle | \langle E \rangle)$$

Problème

Langues slaves, par ex. ukrainien:

- -- 2 nombres : singulier et pluriel,
- -- 7 genres : masculin animé, masculin inanimé, féminin, neutre, double genre, genre commun, sans genre,
- -- 7 cas : nominatif, accusatif, génitif, instrumental, datif, locatif, vocatif.

2 x 7 x 7 = 98 fois la même règle de grammaire!

L'accord en français (grammaire contextuelle)

```
GN = \langle DET \rangle
(\langle A \rangle
(\langle A \rangle)
(\langle A \rangle)
(\langle A \rangle)
```

L'accord en français (grammaire contextuelle)

```
GN = <DET>/<THIS$Nombre=N$Nombre><THIS$Genre=$N$Genre> (<A>/<THIS$Nombre=N$Nombre><THIS$Genre=$N$Genre> | <E>) (N <N>) (<A>/<THIS$Nombre=N$Nombre><THIS$Genre=$N$Genre> | <E>)
```

Grammaire contextuelle qui reconnait les phrases sémantiquement correctes :

• Propriétés des noms Nom :

Classe = Hum | Conc | Abst

• Propriétés des verbes Verbe :

ClasseSujet = Hum | Conc | Abst

ClasseObjet = Hum | Conc | Abst

Ajouter les contraintes d'accord, de temps et distributionnelle à la grammaire suivante :

Phrase = Sujet Verbe Objet

Grammaire contextuelle qui reconnait les phrases sémantiquement correctes :

Ajouter les contraintes d'accord, de temps et distributionnelle à la grammaire suivante :

```
Phrase = Sujet/<THIS$Classe=V$ClasseSujet>
   (v Verbe)
   Objet/<THIS$Classe=V$ClasseObjet>
```

5. Grammaires non restreintes

Grammaires génératives non restreintes

 Aucune restriction : les membres gauches et les membres droits des règles de réécriture peuvent contenir n'importe quelle séquence d'ALUs, de symboles auxiliaires et de mots vides.

• Les grammaires non restreintes permettent de décrire n'importe quel langage récursivement énumérable (donc, n'importe quelle langue naturelle)

Grammaires génératives non restreintes

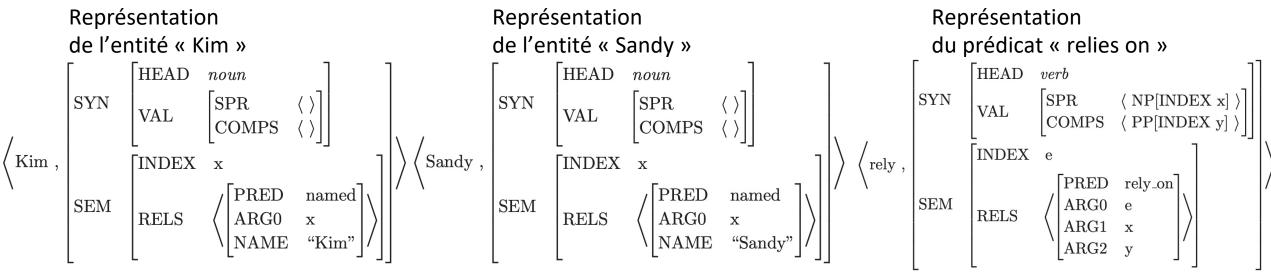
• Exemple de grammaire générative :

```
PHRASE → GN0 mange GN1
    GNO mange GN1 \rightarrow GNO ne mange pas GN1
                  GNO → Luc
              GN1 \rightarrow une pomme
  GN0 mange GN1 → C'est GN0 qui mange GN1
        GN0 mange GN1 → il mange GN1
        GN0 mange GN1 → GN0 la mange
GNO mange GN1 → GN1 est mangée par GNO
```

• •

Grammaires génératives non restreintes alternatives (notation HPSG)

• Exemple de grammaire HPSG qui représente la phrase « Kim relies on Sandy »



Représentation du mot « on »

$$\left\langle \text{on ,} \begin{bmatrix} \text{SYN} & \begin{bmatrix} \text{HEAD} & prep \\ \text{VAL} & \begin{bmatrix} \text{SPR} & \langle \ \rangle \\ \text{COMPS} & \langle \ \text{NP[INDEX x]} \ \rangle \end{bmatrix} \end{bmatrix} \right\rangle$$

$$\left\langle \text{SEM} & \begin{bmatrix} \text{INDEX} & \mathbf{x} \\ \text{RELS} & \langle \ \rangle \end{bmatrix} \right\rangle$$

Représentation de la structure de phrase

Head-Complement Rule
$$\begin{bmatrix} phrase \\ COMPS & \langle \ \rangle \end{bmatrix} \rightarrow \mathbf{H} \begin{bmatrix} word \\ COMPS & \langle \ 1 \ , ..., \ n \ \rangle \end{bmatrix} 1 ... n$$

Règles de combination des matrices

Semantic Compositionality Principle $\begin{bmatrix} \operatorname{RELS} & \boxed{A_1} \oplus \ldots \oplus \boxed{A_n} \end{bmatrix} \to \begin{bmatrix} \operatorname{RELS} & \boxed{A_1} \end{bmatrix} \ldots \begin{bmatrix} \operatorname{RELS} & \boxed{A_n} \end{bmatrix}$ Semantic Inheritance Principle $\begin{bmatrix} \operatorname{INDEX} & \boxed{1} \end{bmatrix} \to \ldots \ \mathbf{H} \begin{bmatrix} \operatorname{INDEX} & \boxed{1} \end{bmatrix} \ldots$

Grammaires génératives non restreintes alternatives (notation simplifiée)

• Exemple (sans les contraintes d'accord) :

NEGATION:

GNO mange GN1 = GNO/THIS <E>/ne mange/THIS <E>/pas GN1/THIS

Le voisin mange la pomme \rightarrow Le voisin ne mange pas la pomme

EXTRACTION_0:

GNO mange GN1 = $\langle E \rangle / C'$ est GNO/THIS $\langle E \rangle / qui$ mange/THIS GN1/THIS

Le voisin mange la pomme \rightarrow C'est le voisin qui mange la pomme

PASSIF_FS:

GN0 mange GN1 = GN0/GN1 <E>/est <V>/THIS\$Participepassé+f+s GN1FS/par GN0 Le voisin mange la pomme \rightarrow la pomme est mangée par le voisin

•••

Analyses sémantiques

• Notation prédicative utilisée en traitement des connaissances :

PRED:

```
(Acteur GN0) (Prédicat mange) (Agent GN1) / Prédicat (Acteur, Agent)

Le voisin mange la pomme → mange (Le voisin, la pomme)
```

• Notation RDF/Turtle (utilisée dans les applications de Web sémantique) :

TURTLE:

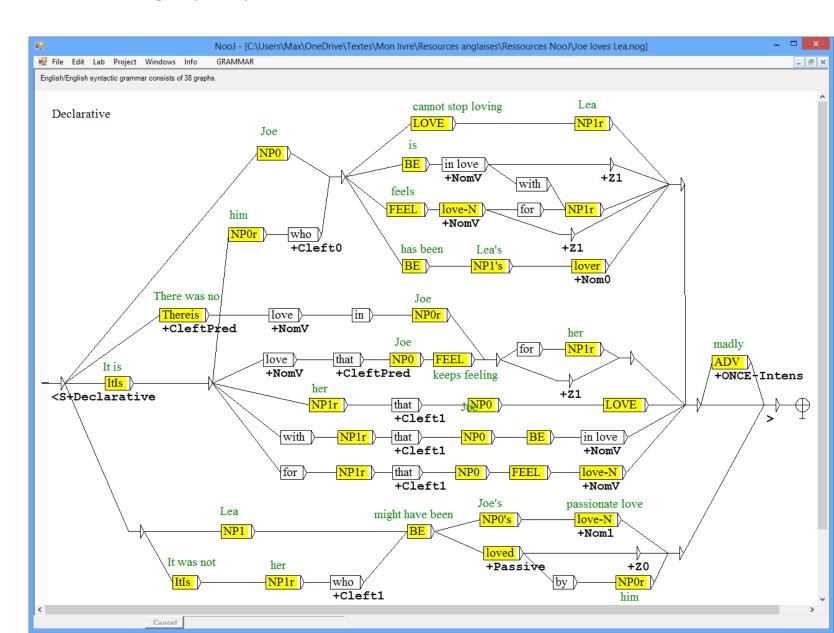
Exemple de grammaire non restreinte

(notation graphique NooJ)

Les noeuds en jaune sont des références à des graphes imbriqués ; ils correspondent aux symboles auxiliaires des grammaires génératives.

Cette grammaire contient 38 graphes

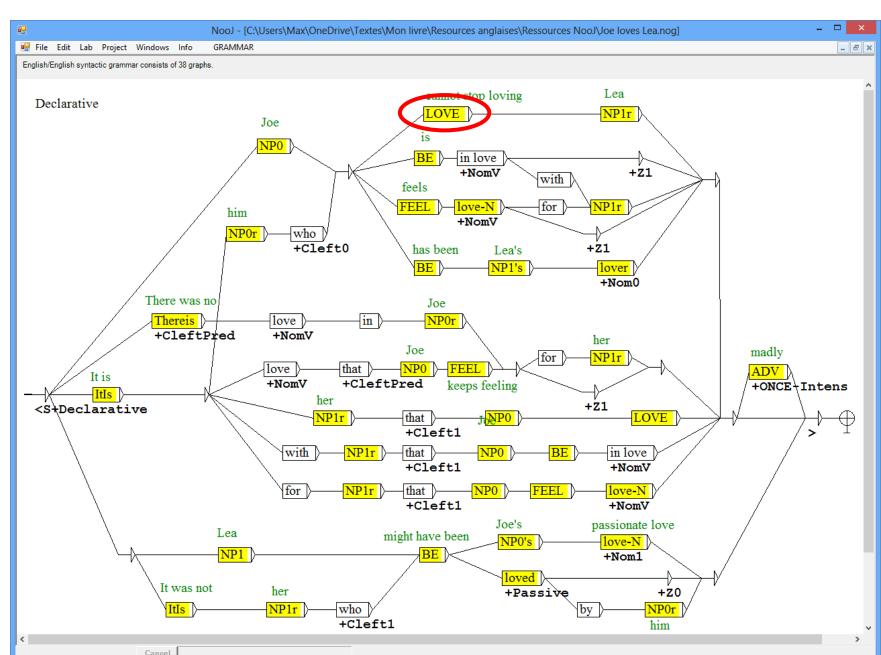
En général, les linguistes préfèrent utiliser la notation graphique lorsque les grammaires deviennent complexes.



Exemple de grammaire transformationnelle sous NooJ (notation graphique)

Les noeuds en jaune sont des références à des graphes imbriqués ; ils correspondent aux symboles auxiliaires des grammaires génératives

Le graphe LOVE est présenté ciaprès

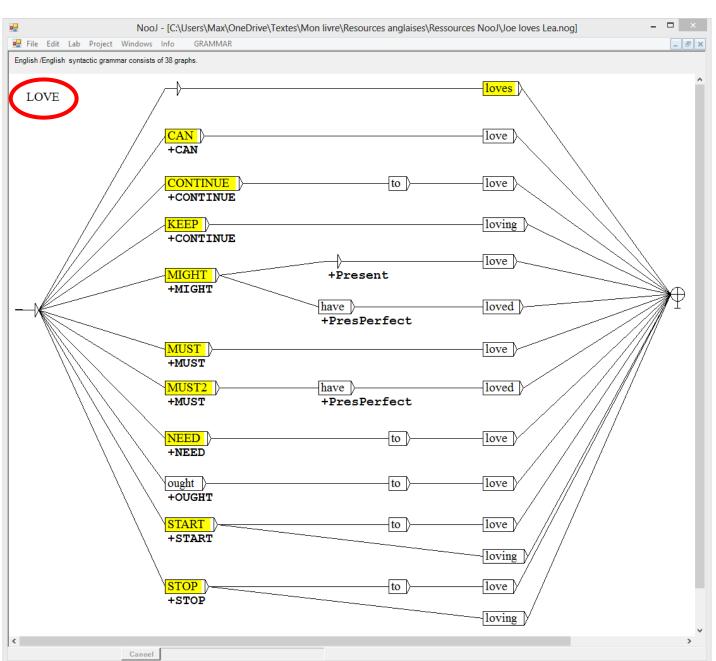


Exemple de grammaire transformationnelle sous NooJ (notation graphique)

Le graphe présenté ici est le graphe LOVE

Cette grammaire décrit les séquences constituées d'une forme conjuguée du verbe to love et des formes avec verbes auxiliaires, aspectuels et modaux, ainsi que les négations éventuelles, ex. :

(The butcher) must not have loved (his wife)



Exemple de grammaire transformationnelle sous NooJ (notation graphique)

On peut appliquer cette grammaire pour produire automatiquement toutes les phrases transformées à partir d'une phrase donnée, ici :

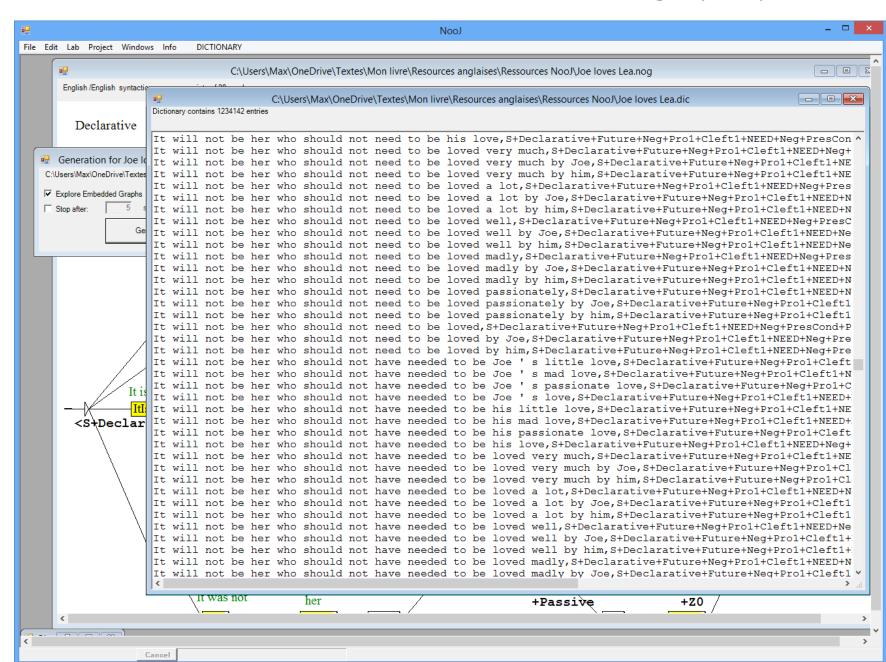
Joe Joves Lea

Chaque phrase transformée est associées à une série de codes sémantiques, ex. : *It is Joe who did not love her* est associée à : +Declarative+Past+Neg+Pron1+Focus0

Certaines phrases sont stylistiquement discutables, ex. :

It will not be her who should not need to be loved by Joe.

On pourra ajouter des contraintes stylistiques sur la taille des phrases, le nombre de négation, limiter le nombre de verbes auxiliaires dans les phrases passives, etc.



Τάλως

Automates





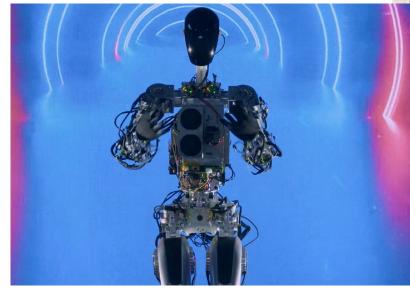
Le Turc



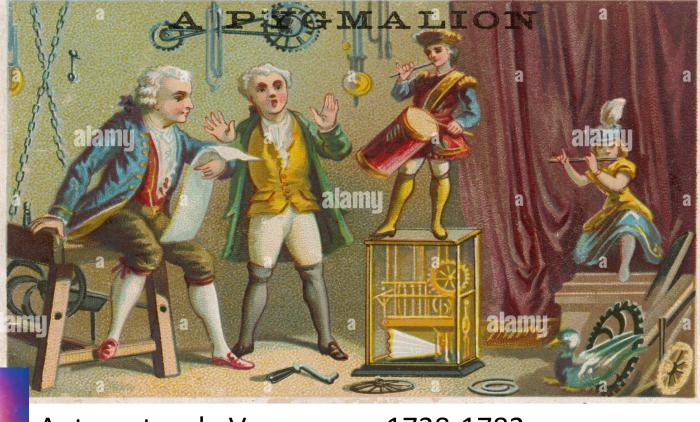


Horloge astronomique de Strasbourg,

XIVe



Robot Tesla, 2022



Automates de Vaucanson, 1738-1782

Automates

Formalisation du comportement des automates

- Un automate est un mécanisme qui
 - est dans un certain état
 - accepte des signaux entrants
 - change d'état en fonction des signaux.

Exemple d'automate :

- une lampe est éteinte ; on presse un bouton => elle s'allume ; on represse le bouton => elle s'éteint
- Poupée interactive : on appuie sur sa main gauche => elle gazouille ; on appuie sur son pied gauche => elle rit ; on effleure son pied droit => elle éternue ; ou elle fait des vocalises ; on lui donne son biberon => elle tète (bruit de succion) ; etc.

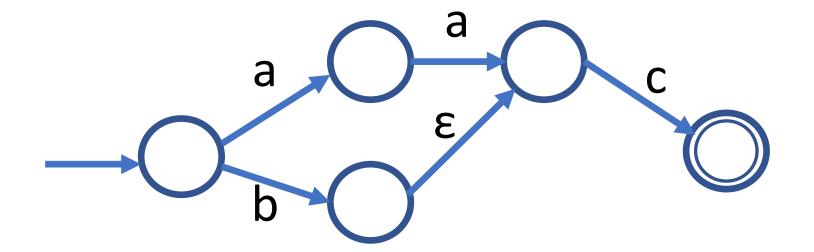




Formalisation du comportement des automates

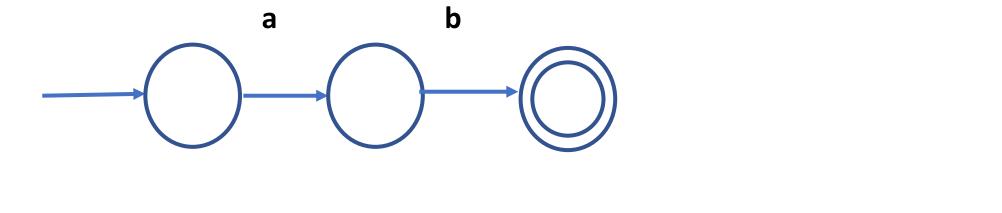
- un alphabet A
- un ensemble fini d'états Q
- un ensemble non vide d'états initiaux I
- un ensemble non vide d'états terminaux T
- un ensemble de transitions, c'est-à-dire de triplets (q1, l, q2) tels que q1 et q2 sont des états de Q, et l est une lettre de A

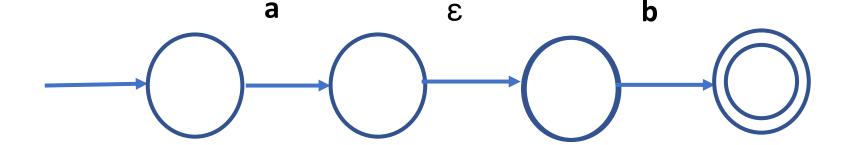
Automate fini: 5 états, 5 transitions



Un automate reconnaît tous les mots épelés par un chemin qui part d'un état initial et qui rejoint un état terminal.

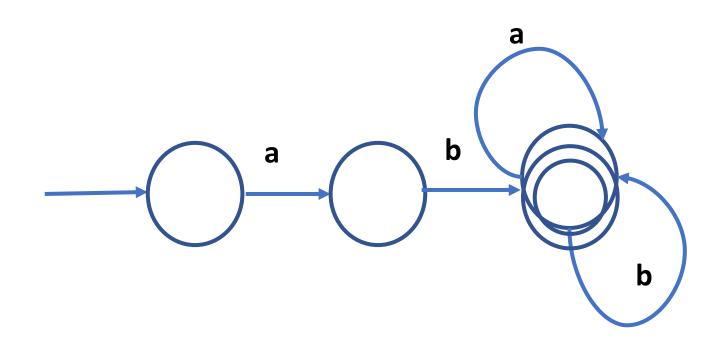
Sur l'alphabet {a, b}, construire l'automate fini qui reconnaît le mot « ab »





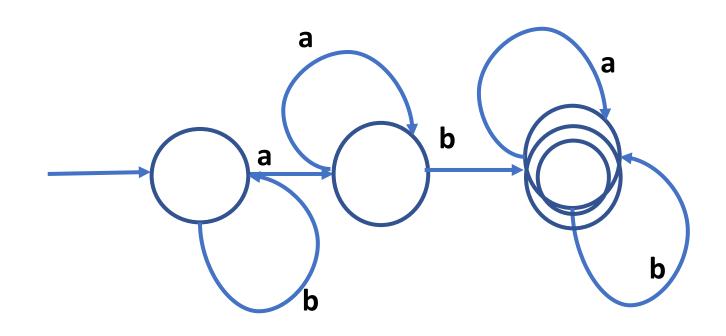
Sur l'alphabet {a, b}, construire l'automate fini qui reconnaît tous les mots qui commencent par « ab »

Sur l'alphabet {a, b}, construire l'automate fini qui reconnaît tous les mots qui commencent par « ab » (mais qui peuvent être suivis par n'importe quelle séquence de lettres)



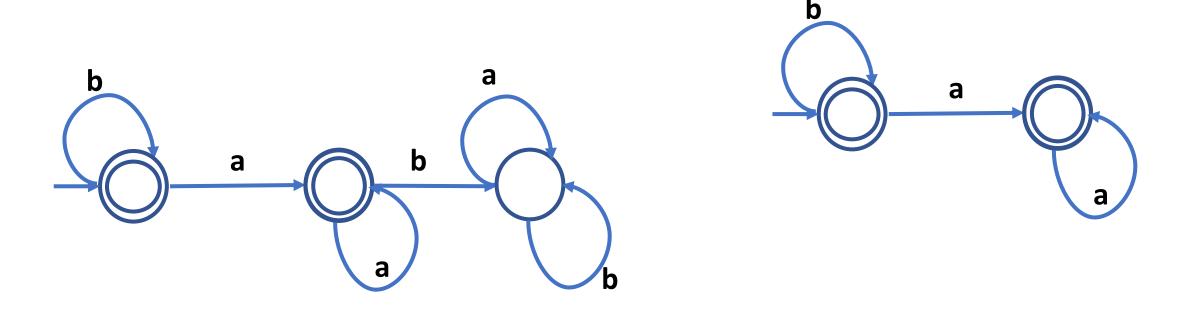
Sur l'alphabet {a, b}, construire l'automate fini qui reconnaît tous les mots qui contiennent l'affixe « ab »

Sur l'alphabet {a, b}, construire l'automate fini qui reconnaît tous les mots qui contiennent l'affixe « ab »



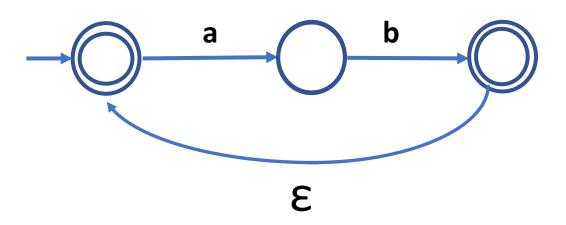
Sur l'alphabet {a, b}, construire l'automate fini qui reconnaît tous les mots qui ne contiennent pas l'affixe « ab »

Sur l'alphabet {a, b}, construire l'automate fini qui reconnaît tous les mots qui ne contiennent pas l'affixe « ab »



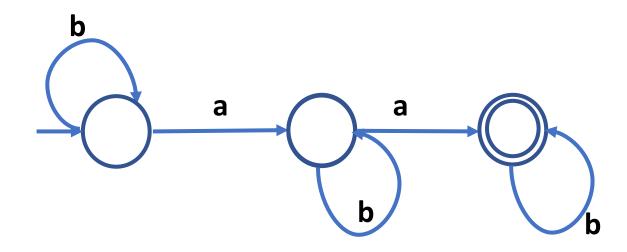
Sur l'alphabet {a, b}, construire l'automate fini qui reconnaît tous les mots qui contiennent un nombre quelconque de suites « ab »

Sur l'alphabet {a, b}, construire l'automate fini qui reconnaît tous les mots qui ne contiennent pas l'affixe « ab »



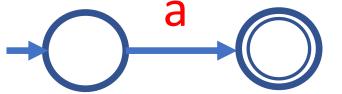
Sur l'alphabet {a, b}, construire l'automate fini qui reconnaît tous les mots qui contiennent exactement deux occurrences de « a »

Sur l'alphabet {a, b}, construire l'automate fini qui reconnaît tous les mots qui contiennent exactement deux occurrences de « a »

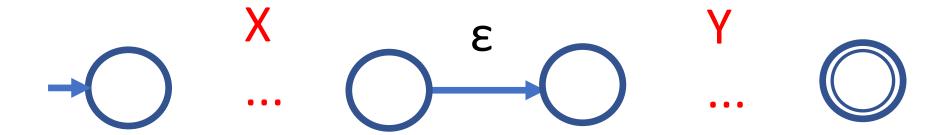


Construction de Thompson (1955)

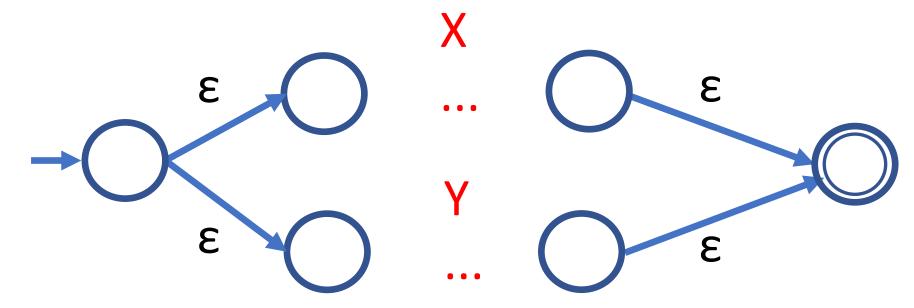
• Pour chaque lettre a, on construit l'automate :



• Pour chaque concaténation X Y, on construit l'automate :



• Pour chaque disjonction X | Y, on construit l'automate :



Pour chaque opération de Kleene X*, on construit l'automate : χ
 Χ

Construction de Thompson

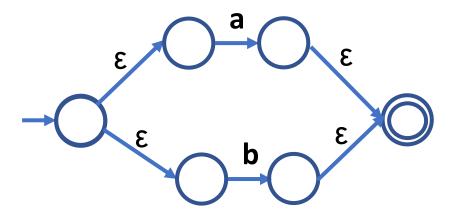
• Sur l'alphabet {a, b, c}, construire en suivant la méthode de Thompson l'automate fini pour l'expression régulière :

a | b

Construction de Thompson

• Sur l'alphabet {a, b, c}, construire en suivant la méthode de Thompson l'automate fini pour l'expression régulière :

a | b

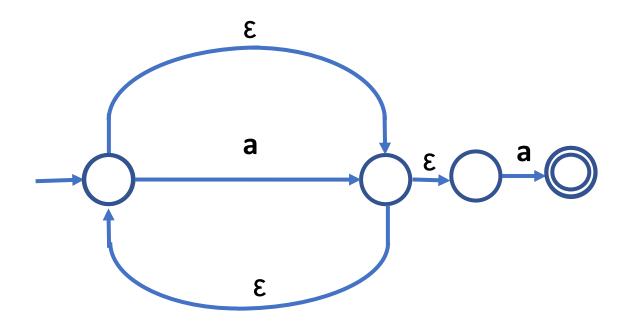


• Sur l'alphabet {a, b, c}, construire en suivant la méthode de Thompson l'automate fini pour l'expression régulière :

a* a

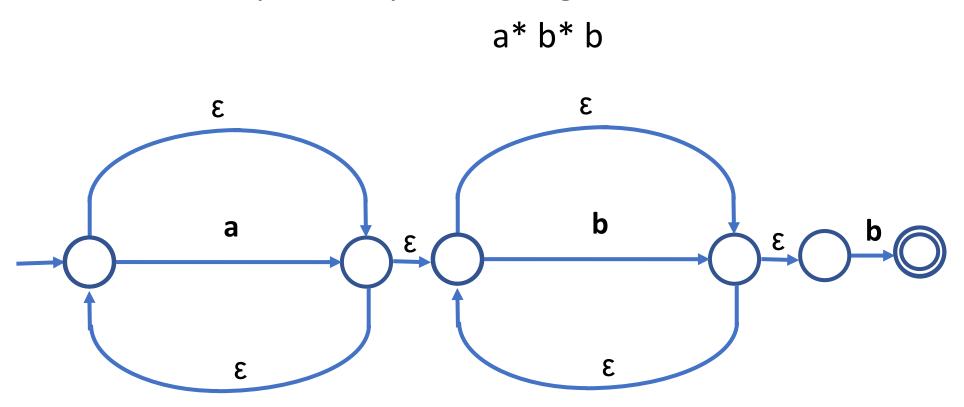
• Sur l'alphabet {a, b, c}, construire en suivant la méthode de Thompson l'automate fini pour l'expression régulière :

a* a



• Sur l'alphabet {a, b, c}, construire en suivant la méthode de Thompson l'automate fini pour l'expression régulière :

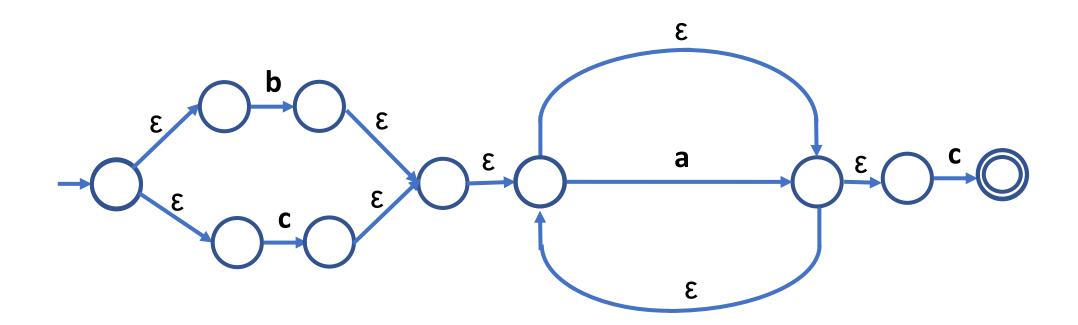
• Sur l'alphabet {a, b, c}, construire en suivant la méthode de Thompson l'automate fini pour l'expression régulière :



ATTENTION : on ne peut pas unifier l'état terminal de l'automate gauche avec l'état initial de l'automate droit car dans ce cas, l'automate reconnaîtrait « abab », ce qui n'est pas décrit par la grammaire

• Sur l'alphabet {a, b, c}, construire en suivant la méthode de Thompson l'automate fini pour l'expression régulière :

• Sur l'alphabet {a, b, c}, construire en suivant la méthode de Thompson l'automate fini pour l'expression régulière :

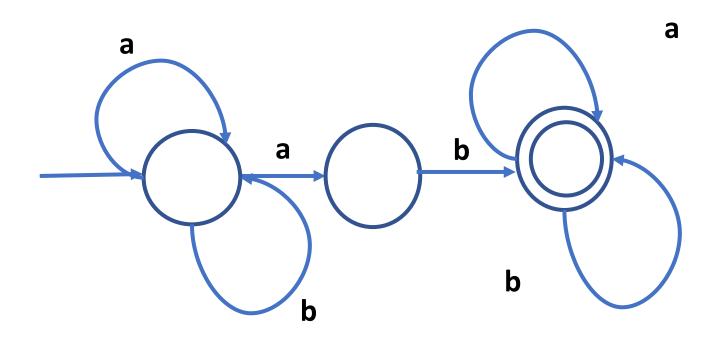


Automates déterministes

- Un automate est déterministe si :
 - il ne contient aucune transition étiquetée par ε
 - pour chaque état de l'automate, il n'existe au plus qu'une seule transition sortante étiquetée par une lettre donnée

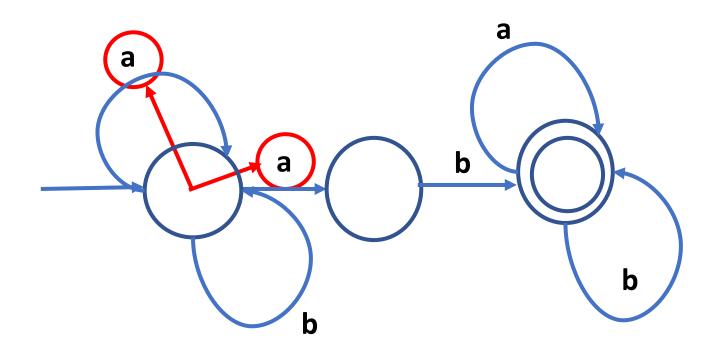
On peut déterminiser tout automate

Automate fini déterministe ou non-déterministe ?

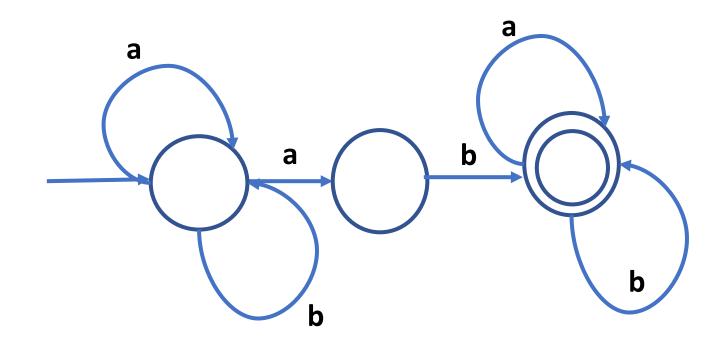


Automate fini non-déterministe :

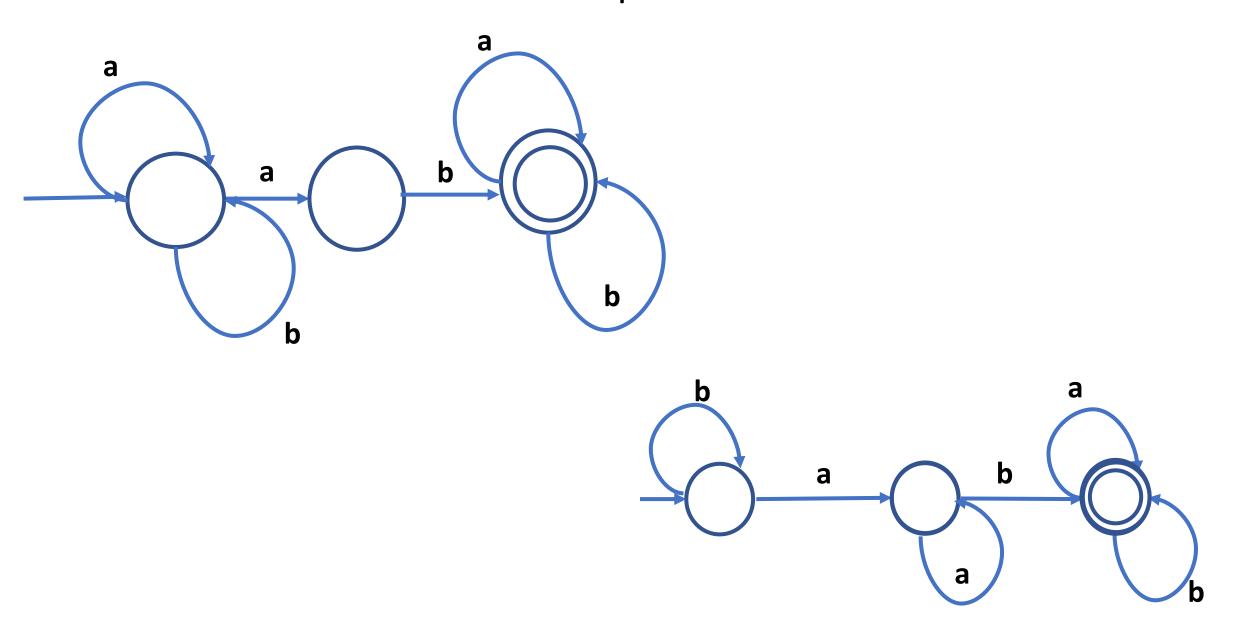
l'état initial contient 2 transitions sortante étiquetée par la même lettre a



Exercice: construire l'automate fini déterministe équivalent, i.e., qui reconnaît tous les mots qui contiennent l'affixe « ab »



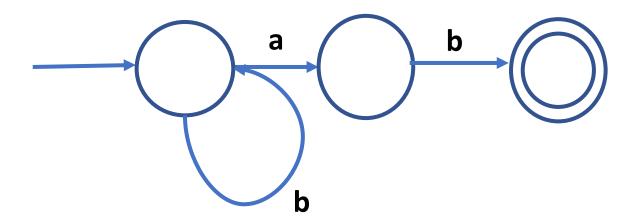
Automate déterministe équivalent



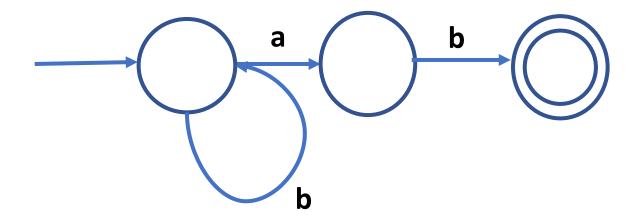
Automates complets

- Un automate est complet si :
 - il est déterministe
 - pour chaque état de l'automate, il existe une transition sortante étiquetée par chaque lettre de l'alphabet
- On peut compléter tout automate

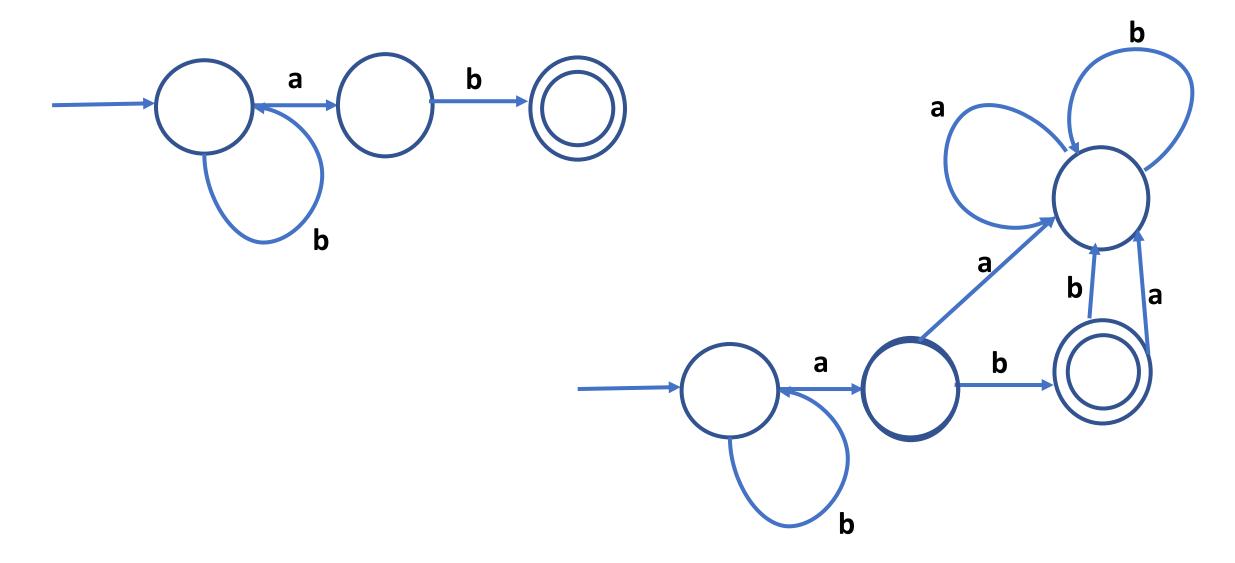
Automate fini complet ou incomplet?



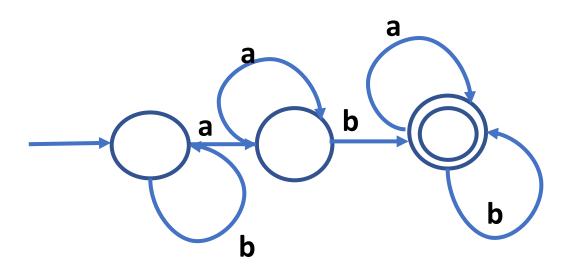
Compléter l'automate



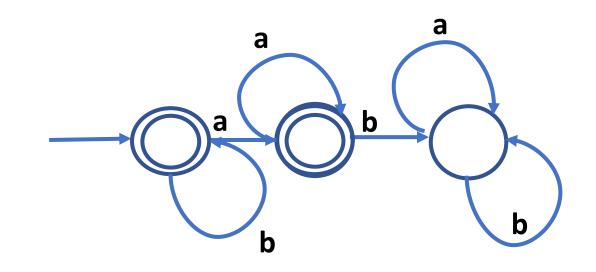
Automate fini complet



Négation : on construit un automate complet, puis on inverse les états terminaux



Cet automate reconnaît tous les mots qui contiennent l'affixe « ab »



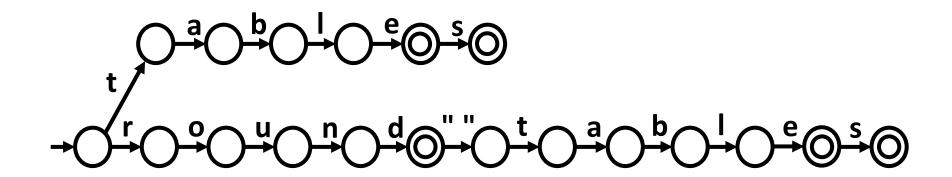
Cet automate reconnaît tous les mots qui ne contiennent pas l'affixe « ab »

• Représenter les entrées d'un dictionnaire de mots simples et de mots composés sous la forme d'un automate fini déterministe :

table round round table

•L'automate doit reconnaître toutes les formes : table, tables, round table, round tables.

• Représenter les entrées d'un dictionnaire de mots simples et de mots composés sous la forme d'un automate fini déterministe :



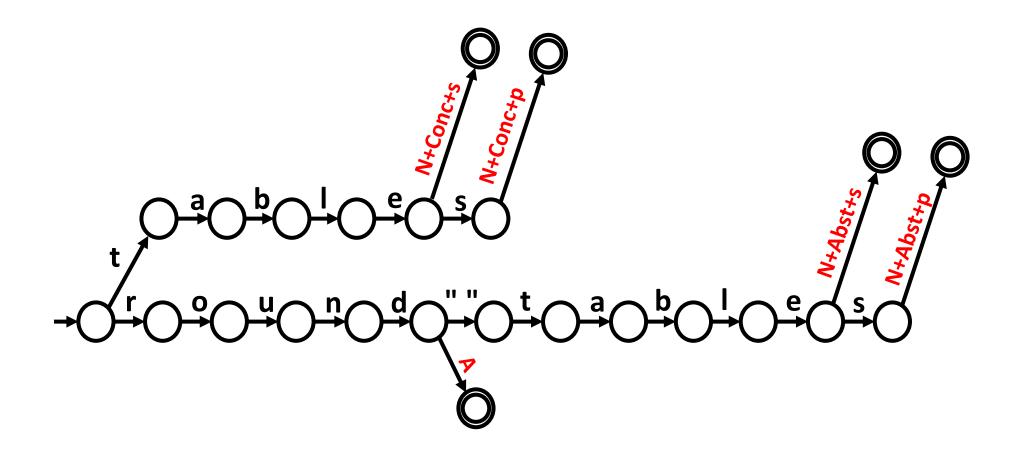
Automates finis avec sorties (outputs): transducteurs finis

• Représenter les entrées d'un dictionnaire de mots simples et de mots composés sous la forme d'un transducteur qui produit leur analyse

```
table,N+Conc
rond,A
table ronde,N+Abst
```

• L'automate doit reconnaître toutes les formes : table, tables, round table, rond tables et produire leurs propriétés lexicales (N, Conc, A, Abst) et aussi leur propriétés morphologiques +s ou +p

• Représenter les entrées d'un dictionnaire de mots simples et de mots composés

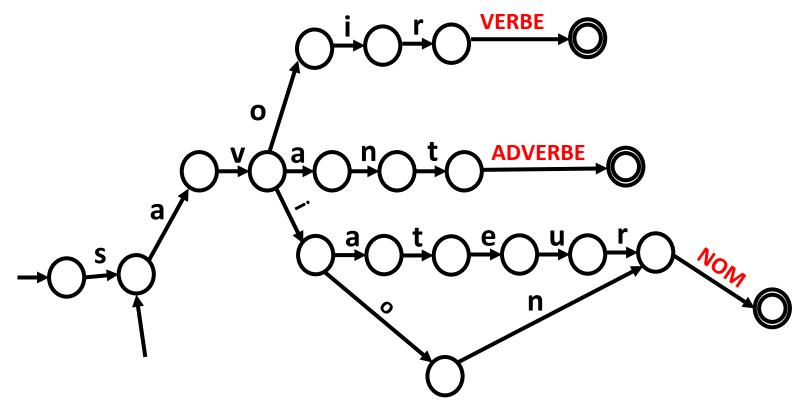


Exercice

•Représenter les entrées d'un dictionnaire sous la forme d'un transducteur déterministe

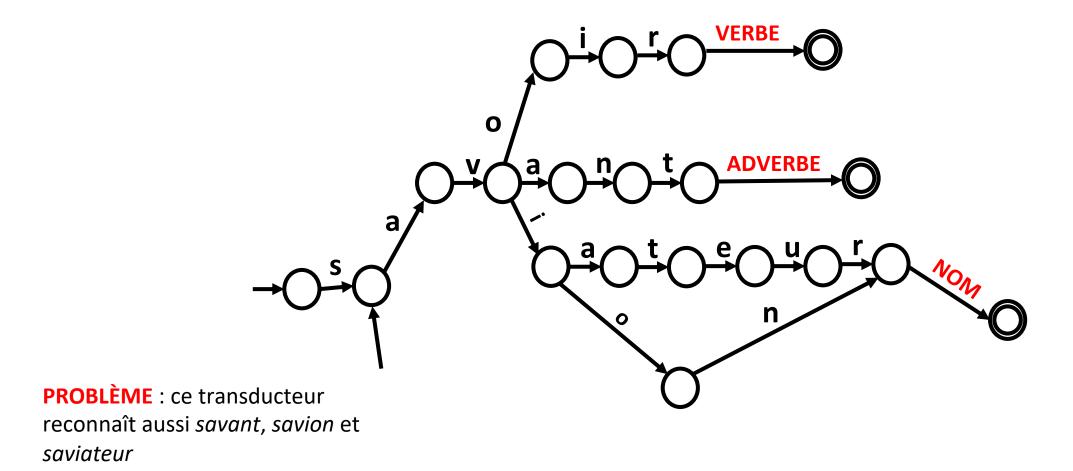
avant, ADVERBE aviateur, NOM avoir, VERBE avion, NOM savoir, VERBE

• Représenter les entrées d'un dictionnaire sous la forme d'un transducteur déterministe

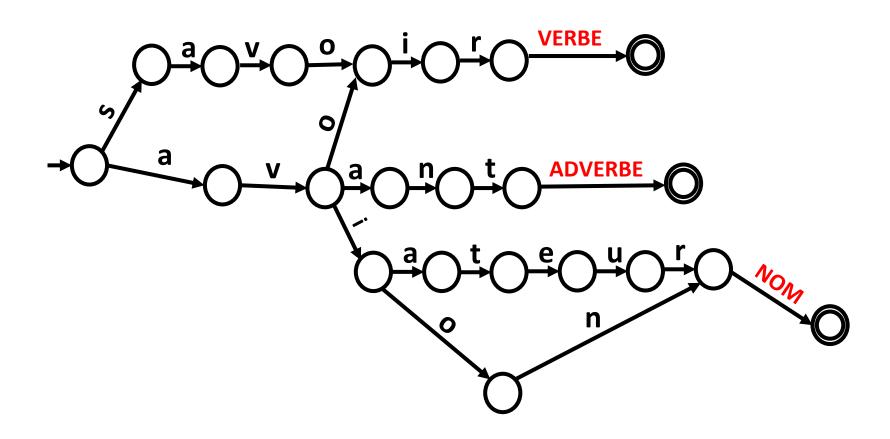


Première tentative : on réutilise le chemin pour *avoir* et *savoir*

• Représenter les entrées d'un dictionnaire sous la forme d'un transducteur déterministe



• Représenter les entrées d'un dictionnaire sous la forme d'un transducteur déterministe



Pour en savoir plus sur les machines

- Grammaires régulières → Automates finis
 - https://fr.wikipedia.org/wiki/Automate_fini
- Grammaires hors-contextes → Automates à pile
 - https://fr.wikipedia.org/wiki/Automate à pile
- Grammaires contextuelles → Automates linéairement bornés
 - https://fr.wikipedia.org/wiki/Automate linéairement borné
- Grammaires non-restreintes → Machines de Turing
 - https://fr.wikipedia.org/wiki/Machine de Turing

Récapitulation

- Equivalence entre grammaires, langages et machines
- Les grammaires génératives sont un mécanisme mathématique pour décrire n'importe quelle langue naturelle
- Hiérarchie de Chomsky-Schützenberger : 4 types de grammaires génératives
- Les grammaires génératives ne sont pas adaptées au développement de grammaires de grande taille ; on a donc inventé des formalismes plus adaptés pour décrire des
 - grammaires régulières (ou rationnelles), ex. XFST, HFST, NooJ
 - grammaires hors-contexte (ou algébriques), ex. YACC, GPSG, NooJ
 - grammaires contextuelles, ex. LFG, TAG, NooJ
 - grammaires non restreintes, ex. HPSG, NooJ
- Chaque type de grammaire peut être traité par des machines/algorithmes plus ou moins efficaces ; plus la grammaire est puissante, plus elle est complexe, et moins l'algorithme est rapide.
- Les grammaires régulières peuvent être traitées par des automates finis : non-déterministes, déterministes, complets, et/ou acycliques.
- Les grammaires régulières augmentées et les dictionnaires peuvent être traitées par des transducteurs finis.