# Formal Languages and Linguistics

Pascal Amsili

Sorbonne Nouvelle, Lattice (CNRS/ENS-PSL/SN)

Cogmaster, september 2022

Sorbonne
Nouvelle

## General introduction

1. Mathematicians (incl. Chomsky) have formalized the notion of **language**                              oversimplification ?
                                                                                                           maybe...

2. It buys us:
    2.1 Tools to think about theoretical issues about language/s
        (expressiveness, complexity, comparability...)
    2.2 Tools to manipulate concretely language (e.g. with computers)
    2.3 A research programme:
        • Represent the syntax of natural language in a fully
          unambiguously specified way

    Now let's get familiar with the mathematical notion of language

Sorbonne ;;;
Nouvelle ;;;

# Overview

Sorbonne
Nouvelle

# Alphabet, word

### Def. 1 (Alphabet)

An *alphabet* $\Sigma$ is a finite set of symbols (letters).
The *size* of the alphabet is the cardinal of the set.

### Def. 2 (Word)

A *word* on the alphabet $\Sigma$ is a finite sequence of letters from $\Sigma$.
Formally, let $[p] = (1, 2, 3, 4, ..., p)$ (ordered integer sequence).
Then a word is a *mapping*

$$u : [p] \longrightarrow \Sigma$$

$p$, the length of $u$, is noted $|u|$.

Sorbonne
Nouvelle

# Examples I

Alphabet   {■, ▬}

Words   ▬ ▬ ▬ ·

■

·■ ▬·

. . .

Alphabet   {▬ , ▬··· , ▬·▬· , ▬·· , · , . . . }

Words   ··· ▬ ▬ ▬ ···

▬·· ·· ▬·· · ·▬· ▬ ▬ ▬ ▬

·▬ ▬··· ▬··· ·▬ ·▬ ▬··· ▬·· ·▬

. . .

# Examples II

| Alphabet | $\{0,1,2,3,4,5,6,7,8,9,\cdot\}$ |
|---|---|
| Words | $235 \cdot 29$ |
| | $007 \cdot 12$ |
| | $\cdot 1 \cdot 1 \cdot 00 \cdot \cdot$ |
| | ~~$3 \cdot 1415962\ldots$~~ $(\pi)$ |
| | $\ldots$ |
| Alphabet | $\{$a, woman, loves, man $\}$ |
| Words | a |
| | a woman loves a woman |
| | man man a loves woman loves a |
| | $\ldots$ |

## Monoid

### Def. 3 ($\Sigma^*$)

Let $\Sigma$ be an alphabet.
The set of all the words that can be formed with any number of letters from $\Sigma$ is noted $\Sigma^*$

$\Sigma^*$ includes a word with no letter, noted $\varepsilon$

Example: $\quad \Sigma \quad = \{a, b, c\}$
$\qquad\quad\; \Sigma^* \; = \{\varepsilon, a, b, c, aa, ab, ac, ba, \ldots, bbb, \ldots\}$

N.B.: $\Sigma^*$ is always infinite, except...

## Monoid

### Def. 3 ($\Sigma^*$)

Let $\Sigma$ be an alphabet.
The set of all the words that can be formed with any number of letters from $\Sigma$ is noted $\Sigma^*$

$\Sigma^*$ includes a word with no letter, noted $\varepsilon$

Example: $\quad \Sigma \quad = \{a, b, c\}$
$\qquad\quad \Sigma^* \quad = \{\varepsilon, a, b, c, aa, ab, ac, ba, \ldots, bbb, \ldots\}$

N.B.: $\Sigma^*$ is always infinite, except...
$\qquad\qquad\qquad$ if $\Sigma = \emptyset$. Then $\Sigma^* = \{\varepsilon\}$.

Sorbonne
Nouvelle

## Structure of $\Sigma^*$

Let $k$ be the size of the alphabet $k = |\Sigma|$.

Then $\Sigma^*$ contains :
$$k^0 = 1 \quad \text{word(s) of 0 letters } (\varepsilon)$$
$$k^1 = k \quad \text{word(s) of 1 letters}$$
$$k^2 \qquad \text{word(s) of 2 letters}$$
$$\ldots$$
$$k^n \qquad \text{words of } n \text{ letters, } \forall n \geq 0$$

# Representation of $\Sigma^*$

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



▶ Words can be enumerated according to different orders
▶ $\Sigma^*$ is a *countable* set

## Concatenation

$\Sigma^*$ can be equipped with a binary operation: *concatenation*

### Def. 4 (Concatenation)

Let $[p] \xrightarrow{u} \Sigma$, $[q] \xrightarrow{w} \Sigma$. The concatenation of $u$ and $w$, noted $uw$ $(u.w)$ is thus defined:

$$uw : \quad [p+q] \longrightarrow \Sigma$$
$$uw_i = \begin{cases} u_i & \text{for} \quad i \in [1, p] \\ w_{i-p} & \text{for} \quad i \in [p+1, p+q] \end{cases}$$

Formal Languages   Regular Languages   Formal Grammars   Formal complexity of Natural Languages   References
○○○○○○○●○○○            ○○                ○○○○○○○○○○○○○○○○○○○○○○
○○○○○                 ○○○○○○○○           ○○○○○○○○○
                      ○○○○○○○○○○○○        ○○○○○○○○○○○○
                                        ○○○○○○○○○○

Basic concepts

## Concatenation

$\Sigma^*$ can be equipped with a binary operation: *concatenation*

### Def. 4 (Concatenation)

Let $[p] \xrightarrow{u} \Sigma$, $[q] \xrightarrow{w} \Sigma$. The concatenation of $u$ and $w$, noted $uw$ $(u.w)$ is thus defined:

$$uw : \quad [p+q] \longrightarrow \Sigma$$
$$uw_i = \begin{cases} u_i & \text{for} \quad i \in [1, p] \\ w_{i-p} & \text{for} \quad i \in [p+1, p+q] \end{cases}$$

Example :   $u$   bacba
            $v$   cca

## Concatenation

$\Sigma^*$ can be equipped with a binary operation: *concatenation*

### Def. 4 (Concatenation)

Let $[p] \xrightarrow{u} \Sigma$, $[q] \xrightarrow{w} \Sigma$. The concatenation of $u$ and $w$, noted $uw$ $(u.w)$ is thus defined:

$$uw : \quad [p+q] \longrightarrow \Sigma$$
$$uw_i = \begin{cases} u_i & \text{for} \quad i \in [1, p] \\ w_{i-p} & \text{for} \quad i \in [p+1, p+q] \end{cases}$$

Example :   $u$    bacba
            $v$    cca
            $uv$   bacbacca

# Factor

### Def. 5 (Factor)

A *factor* $w$ of $u$ is a subset of adjascent letters in $u$.
$-w$ is a factor of $u$ $\qquad \Leftrightarrow \quad \exists u_1, u_2$ s.t. $u = u_1 w u_2$
$-w$ is a left factor (*prefix*) of $u$ $\quad \Leftrightarrow \quad \exists u_2$ s.t. $u = w u_2$
$-w$ is a right factor (*suffix*) of $u$ $\quad \Leftrightarrow \quad \exists u_1$ s.t. $u = u_1 w$

### Def. 6 (Factorization)

We call *factorization* the decomposition of a word into factors.

## Role of concatenation

1. Words have been defined on $\Sigma$.
   If one takes two such words, it's always possible to form a new word by concatenating them.

2. Any word can be factorised in many different ways:
   *a b a c c a b*

## Role of concatenation

1. Words have been defined on $\Sigma$.
   If one takes two such words, it's always possible to form a new word by concatenating them.

2. Any word can be factorised in many different ways:
   *a b a c c a b*
   (*a b a*)(*c c a b*)

## Role of concatenation

1. Words have been defined on Σ.
   If one takes two such words, it's always possible to form a new word by concatenating them.

2. Any word can be factorised in many different ways:
   *a b a c c a b*
   (*a b*)(*a c c*)(*a b*)

## Role of concatenation

1. Words have been defined on Σ.
   If one takes two such words, it's always possible to form a new word by concatenating them.

2. Any word can be factorised in many different ways:
   *a b a c c a b*
   (*a b a c c*)(*a b*)

# Role of concatenation

1. Words have been defined on Σ.
   If one takes two such words, it's always possible to form a new
   word by concatenating them.

2. Any word can be factorised in many different ways:
   *a b a c c a b*
   (*a*)(*b*)(*a*)(*c*)(*c*)(*a*)(*b*)

# Role of concatenation

1. Words have been defined on $\Sigma$.
   If one takes two such words, it's always possible to form a new word by concatenating them.

2. Any word can be factorised in many different ways:
   *a b a c c a b*
   (*a*)(*b*)(*a*)(*c*)(*c*)(*a*)(*b*)

3. Since all letters of $\Sigma$ form a word of length 1
   (this set of words is called the *base*),

4. any word of $\Sigma^*$ can be seen as a (unique) sequence of concatenations of length 1 words :
   *a b a c c a b*
   $(((((( ab)a)c)c)a)b)$
   $(((((( a.b).a).c).c).a).b)$

## Properties of concatenation

1. Concatenation is non commutative
2. Concatenation is associative
3. Concatenation has an identity (neutral) element: $\varepsilon$

1. $uv.w \neq w.uv$
2. $(u.v).w = u.(v.w)$
3. $u.\varepsilon = \varepsilon.u = u$

Notation : $a.a.a = a^3$

# Overview

# Language

### Def. 7 (Formal Language)

Let $\Sigma$ be an alphabet.
A language on $\Sigma$ is a set of words on $\Sigma$.

## Language

### Def. 7 (Formal Language)

Let $\Sigma$ be an alphabet.
A language on $\Sigma$ is a set of words on $\Sigma$.
or, equivalently,
A language on $\Sigma$ is a subset of $\Sigma^*$

# Examples I

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

## Examples I

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

$$L_1 = \{aa, ab, bac\} \qquad \text{finite language}$$

# Examples I

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

| $L_1 = \{aa, ab, bac\}$ | finite language |
|---|---|
| $L_2 = \{a, aa, aaa, aaaa \ldots\}$ | |

## Examples I

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

| | |
|---|---|
| $L_1 = \{aa, ab, bac\}$ | finite language |
| $L_2 = \{a, aa, aaa, aaaa \ldots\}$ | |
| or $L_2 = \{a^i \ / \ i \geq 1\}$ | infinite language |

## Examples I

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

| | |
|---|---|
| $L_1 = \{aa, ab, bac\}$ | finite language |
| $L_2 = \{a, aa, aaa, aaaa \ldots\}$ | |
| or $L_2 = \{a^i \; / \; i \geq 1\}$ | infinite language |
| $L_3 = \{\varepsilon\}$ | finite language, |
| | reduced to a singleton |

## Examples I

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

| | |
|---|---|
| $L_1 = \{aa, ab, bac\}$ | finite language |
| $L_2 = \{a, aa, aaa, aaaa \ldots\}$ | |
| or $L_2 = \{a^i \ / \ i \geq 1\}$ | infinite language |
| $L_3 = \{\varepsilon\}$ | finite language, |
| | reduced to a singleton |
| $\neq$ | |

## Examples I

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

| | |
|---|---|
| $L_1 = \{aa, ab, bac\}$ | finite language |
| $L_2 = \{a, aa, aaa, aaaa \ldots\}$ | |
|    or $L_2 = \{a^i \ / \ i \geq 1\}$ | infinite language |
| $L_3 = \{\varepsilon\}$ | finite language, |
| | reduced to a singleton |
| | $\neq$ |
| $L_4 = \emptyset$ | "empty" language |

## Examples I

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

| | |
|---|---|
| $L_1 = \{aa, ab, bac\}$ | finite language |
| $L_2 = \{a, aa, aaa, aaaa \dots\}$ | |
|    or $L_2 = \{a^i \ / \ i \geq 1\}$ | infinite language |
| $L_3 = \{\varepsilon\}$ | finite language, |
| | reduced to a singleton |
| | $\neq$ |
| $L_4 = \emptyset$ | "empty" language |
| $L_5 = \Sigma^*$ | |

## Examples II

Let $\Sigma = \{a, man, loves, woman\}$.

## Examples II

Let $\Sigma = \{$a, man, loves, woman$\}$.

$L = \{$ a man loves a woman, a woman loves a man $\}$

## Examples II

Let $\Sigma = \{$a, man, loves, woman$\}$.

$L = \{$ a man loves a woman, a woman loves a man $\}$

Let $\Sigma' = \{$a, man, who, saw, fell$\}$.

## Examples II

Let $\Sigma = \{$a, man, loves, woman$\}$.

$L = \{$ a man loves a woman, a woman loves a man $\}$

Let $\Sigma' = \{$a, man, who, saw, fell$\}$.

$$L' = \left\{ \begin{array}{l} \text{a man fell,} \\ \text{a man who saw a man fell,} \\ \text{a man who saw a man who saw a man fell,} \\ \dots \end{array} \right\}$$

## Set operations

Since a language is a set, usual set operations can be defined:

- ▶ union
- ▶ intersection
- ▶ set difference

## Set operations

Since a language is a set, usual set operations can be defined:

- ▶ union
- ▶ intersection
- ▶ set difference

⇒ One may describe a (complex) language as the result of set operations on (~~simpler~~) languages:
$$\{a^{2k} \mid k \geqslant 1\} = \{a, aa, aaa, aaaa, \ldots\} \cap \{ww \mid w \in \Sigma^*\}$$

## Additional operations

### Def. 8 (product operation on languages)

One can define the *language product* and its closure *the Kleene star* operation:

▶ The *product* of languages is thus defined:
$$L_1.L_2 \;=\; \{uv \; / \; u \in L_1 \; \& \; v \in L_2\}$$

Notation: $\overbrace{L.L.L\ldots L}^{k \text{ times}} = L^k$ ; $L^0 = \{\varepsilon\}$

▶ The Kleene star of a language is thus defined:
$$L^* = \bigcup_{n \geqslant 0} L^n$$

# References I

Bar-Hillel, Yehoshua, Perles, Micha, & Shamir, Eliahu. 1961. On formal properties of simple phrase structure grammars. *STUF-Language Typology and Universals*, 14(1-4), 143–172.

Chomsky, Noam. 1957. *Syntactic Structures*. Den Haag: Mouton & Co.

Chomsky, Noam. 1995. *The Minimalist Program*. Vol. 28. Cambridge, Mass.: MIT Press.

Gazdar, Gerald, & Pullum, Geoffrey K. 1985 (May). *Computationally Relevant Properties of Natural Languages and Their Grammars*. Tech. rept. Center for the Study of Language and Information, Leland Stanford Junior University.

Gibson, Edward, & Thomas, James. 1997. The Complexity of Nested Structures in English: Evidence for the Syntactic Prediction Locality Theory of Linguistic Complexity. *Unpublished manuscript, Massachusetts Institute of Technology*.

Joshi, Aravind K. 1985. *Tree Adjoining Grammars: How Much Context-Sensitivity is Required to Provide Reasonable Structural Descriptions?* Tech. rept. Department of Computer and Information Science, University of Pennsylvania.

Langendoen, D Terence, & Postal, Paul Martin. 1984. *The vastness of natural languages*. Basil Blackwell Oxford.

Mannell, Robert. 1999. *Infinite number of sentences*. part of a set of class notes on the Internet. http://clas.mq.edu.au/speech/infinite_sentences/.

Schieber, Stuart M. 1985. Evidence against the Context-Freeness of Natural Language. *Linguistics and Philosophy*, 8(3), 333–343.

Stabler, Edward P. 2011. Computational perspectives on minimalism. *Oxford handbook of linguistic minimalism*, 617–643.

Sorbonne
Nouvelle

# References II

Steedman, Mark, *et al.* . 2012 (June). *Combinatory Categorial Grammars for Robust Natural Language Processing*. Slides for NASSLLI course
http://homepages.inf.ed.ac.uk/steedman/papers/ccg/nasslli12.pdf.

Vijay-Shanker, K., & Weir, David J. 1994. The Equivalence of Four Extensions of Context–Free Grammars. *Mathematical Systems Theory*, **27**, 511–546.