

# *Alphabets, words, and languages*

Chapitre 01

# Alphabet

- A non-empty finite set of symbols (letters)
- An alphabet will, for example, be denoted  $X$  or  $\Sigma$ .
- For example
  - $\Sigma = \{0,1\}$  : alphabet of binary numbers.
  - $\Sigma = \{a,b,\dots,z\}$  : the set of lowercase letters.

# Word

- A word (finite)  $\omega$  over the alphabet  $\Sigma$  is a (finite) sequence of letters and is denoted by simple juxtaposition:

$$\omega = x_1 x_2 \dots x_n \text{ where } \omega_i \in \{1, \dots, n\}, x_i \in \Sigma$$

Example:

- `abbac` and `bccca` are two words on the alphabet  $\{a, b, c\}$ .
- `01101` is a word on the alphabet  $\{0, 1\}$ .

# Length of a word

- The number of characters (letters, digits, or other symbols) it contains. For instance, in the word sequence  $\omega = x_1x_2\dots x_n$ , the length of the word is  $n$ , representing the total number of symbols in the sequence.

Thus,  $|abbac| = 5$  and  $|ba| = 2$ .

- We also define the number of occurrences of a letter  $d$  from  $\Sigma$  in  $\omega$ ;  $|\omega|_d$

Example :

- let  $\omega = |00011001|$  be a word in  $\Sigma = \{0,1\}$  hence :  $|\omega|_0=5$  and  $|\omega|_1=3$ .

# The empty word

- The empty word is a word without symbols and therefore of length 0. This word is represented by the symbol  $\varepsilon$  ( $|\varepsilon| = 0$ ).

# Concatenation of words

- Let  $X$  be an alphabet,  $x \in \Sigma^*$  is a word of length  $m$ , and  $y \in \Sigma^*$  is a word of length  $n$ . The concatenation of  $x$  and  $y$ , denoted  $xy$ , is the word of length  $m+n$  whose first  $m$  symbols represent a word equal to  $x$ , and the last  $n$  symbols represent a word equal to  $y$ .
- More specifically, if  $x = a_1 a_2 \dots a_m$  and  $y = b_1 b_2 \dots b_n$  then
$$xy = a_1 a_2 \dots a_m b_1 b_2 \dots b_n .$$

# Concatenation of words

- Example :

Let  $x=01101$  and  $y=001$ , then

$xy=01101001$  and  $yx=00101101$ .

- Concatenation is associative ( $(xy)z = x(yz)$ ) but generally not commutative.
- The empty word is the neutral element for concatenation :  $\varepsilon x = x \varepsilon = x$ .
- Concatenation is regular on both the right and the left
  - $wu = wv \Rightarrow u = v$
  - $uw = vw \Rightarrow u = v$
- $|uv| = |u| + |v|$

# Power of an alphabet :

- Let  $\Sigma$  un alphabet, be an alphabet, we denote by  $\Sigma^k$  the set of all words of a given length  $k$  over this alphabet.
- Examples :
  - $\Sigma^0 = \{\varepsilon\}$  whatever the alphabet  $\Sigma$ .
  - if  $\Sigma = \{a,b\}$  then  $\Sigma^1 = \{a,b\}$ ,  $\Sigma^2 = \{aa,ab,ba,bb\}$ ,  
 $\Sigma^3 = \{aaa,aab,aba,abb,baa,bab,bba,bbb\}$  , ....



# Power of an alphabet :

- The set of all words over  $\Sigma$  is denoted by  $\Sigma^*$ .
- For instance  $\{a, b, c\}^* = \{\epsilon, a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc, aaa, aab, \dots\}$ .

In another way :

- $\Sigma^* = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \dots$
- Sometimes we want to exclude the empty word from the set of words. The set of non-empty words over the alphabet  $X$  is denoted by  $X^+$ .
- $\Sigma^+ = \Sigma^1 \cup \Sigma^2 \cup \Sigma^3 \cup \dots$
- $\Sigma^* = \Sigma^+ \cup \{\epsilon\}$

# Mirror of a word

- Let  $\Sigma$  be an alphabet and let  $\omega \in \Sigma^*$ ,  $|\omega| = n$  and  $\omega = \omega_0 \omega_1 \dots \omega_n$ ,  $n \geq 0$
- The mirror (reverse) of  $\omega$  is  $\omega^R = \omega_n \omega_{n-1} \dots \omega_1 \omega_0$
- Property:
  - $\forall u, v \in \Sigma^* : (uv)^R = v^R u^R$
- Palindrome :
- Let  $\Sigma$  be an alphabet and let  $\omega \in \Sigma^*$ , a word is called a palindrome if :  $\omega = \omega^R$

Example :  $\omega = 0110 \rightarrow \omega^R = 0110 \rightarrow \omega = \omega^R$

# Factors

- We call a left factor of  $\omega$  a word  $u$  such that  $uv = \omega$ .
- We call a right factor a word  $v$  such that  $uv = \omega$ .
- We call a factor of  $w$  a word  $u$  such that there exist  $v$  and  $v'$  such that  $vuv' = \omega$ .

Example:

- "bon" is a left factor of the word "bonjour".
- "jour" is a right factor of the word "bonjour".
- "jo" is a factor of the word "bonjour".

# *Languages*

---

# Langages

We call a **language over  $\Sigma$**  any set of words over  $\Sigma$

*Definition :*

- A **(formal) language** is any subset  $L$  of  $\Sigma^*$ , that is  $L \subset \Sigma^*$ .

*Examples :*

- $L_1 = \Sigma^*$  ,  $L_2 = \emptyset$  ,  $L_3 = \{\varepsilon\}$  ,  $L_4 = \{ \omega \in \Sigma^* , = \omega_1 ab \omega_2 \}$
- A language can be finite or infinite

Let  $\Sigma = \{0,1\}$

- $L_1 = \{ \omega \in \Sigma^* , \omega_1 \equiv [3] \}$  ,  $L_1$  is infinite.
- $L_2 = \{ \omega \in \Sigma^* , |\omega| < 5 \}$  ,  $L_2$  est finite

# Languages

- Note :

Among languages, it is important to distinguish:

- The set  $\emptyset$  (the empty set, which contains no words).
- The language  $\{\epsilon\}$  (the language that contains only the empty word as its sole element).

# Notes :

- A **finite language** is a language that contains a finite number of words.
- The **empty language** contains no words.
- A language is said to be **proper** if it does not contain the empty word.
- A language is **infinite** if it is neither empty nor finite.
- Some infinite languages (**semi-decidable languages**) can be described by a set of rules called a **formal grammar**. There are other infinite languages for which no description method exists; these are called **undecidable languages**.

# Operations on languages

- Let  $X$  be an alphabet, a certain number of operations can be performed on languages:
- **union** :  $L1 \cup L2 = \omega \in X^* \mid \omega \in L1 \text{ or } \omega \in L2$
- **intersection** :  $L1 \cap L2 = \omega \in X^* \mid \omega \in L1 \text{ and } \omega \in L2$
- **complement** with respect to  $X^*$  :  $L = \omega \in X^* \mid \omega \notin L$
- **difference** :  $L1 - L2 = L1 \cap \bar{L2} = \omega \mid \omega \in X^* \mid \omega \in L1 \text{ and } \omega \notin L2$
- **concatenation** :  $L1.L2 = u.v \mid u \in L1 \text{ et } v \in L2$



# Operations on languages

- Power of a language:  $L^2 = L.L$  et  $\forall n \in \mathbb{N}$ ,  
 $L^{n+1} = L^n . avec L^0 = \varepsilon$
- (The transition to) the Kleene Star
- The language  $L^*$  (Kleene star of  $L$ ) is defined by:
  - $L^* = L^0 \cup L \cup L^2 \dots \cup L^n \dots =$   
 $\{u \mid \exists n \in \mathbb{N}, u_1, \dots, u_n \in L \text{ tel que } u = u_1 \dots u_n \}$
- The plus operation :  $L^+ = L^* . L = L \cup L^2 \dots \cup L^n \dots$

# Example

$$L_1 = \{\varepsilon, aa\}, L_2 = \{a^i b^j / i, j \geq 0\} \text{ et } L_3 = \{ab, b\}.$$

$$L_1.L_2, L_1.L_3, L_1 \cup L_2, L_2 \cap L_3, L_1^{10}, L_1^*, L_1^+, L_2^R.$$

**Solutions :**

- $L_1.L_2 = L_2$ ;
- $L_1.L_3 = \{ab, b, aaab, aab\}$ ;
- $L_1 \cup L_2 = L_2$ ;
- $L_2 \cap L_3 = L_3$ ;
- $L_1^{10} = \{a^{2n} / 10 \geq n \geq 0\}$ ;
- $L_1^* = L_1^+ = \{a^{2n} n \geq 0\}$ ;
- $L_2^R = \{b^i a^j / i, j \geq 0\}$ .

# Examples of languages

$$\Sigma = \{a\} \quad L_1 = \{\varepsilon, a, aa, aaa, \dots\}$$

$$\Sigma = \{a, b\} \quad L_2 = \{\varepsilon, ab, aabb, aaabbb, aaaabbbb, \dots\}$$

$$\Sigma = \{a, b\} \quad L_2 = \{\varepsilon, ab, aabb, aaabbb, aaaabbbb, \dots\}$$

$$\Sigma = \{a, b\} \quad L_3 = \{\varepsilon, aa, bb, aaaa, abba, baab, bbbb, \dots\}$$

$$\Sigma = \{a, b, c\} \quad L_4 = \{\varepsilon, abc, aabbcc, aaabbbccc, \dots\}$$

# Description of languages

- Description in natural language : .
- Language L1 over the alphabet  $\{0,1\}$ : set of words whose interpretation as integers are multiples of three. It includes, for example, the word: 1001 but not 1000.
- Language L2 over the alphabet  $\{a,b\}$ : formed of all palindrome words. Thus, language L2 contains abbabba but not abbabab.
- L1 over the alphabet  $\{0,1\}$ : set of words whose interpretation as integers are multiples of three.

# Description of language

- **Enumerative descriptions :**
- They are clearly used for finite languages, but also for certain infinite languages such as: :
- $L3 = \{a^n b^n / n \geq 1\}$  : containing all words formed exactly of a sequence of  $n$  occurrences of the letter  $a$  followed by a sequence containing the same number of occurrences  $n$  of the letter  $b$ .

# Description of languages

- **Definition by expression:**
- for example, the expression  $ab^*cab$  represents the language  $L_4$ , whose words start with an occurrence of the letter a, followed by any number (possibly zero) of occurrences of the letter b, followed by the right factor cab

# Description of languages

- **Generative mechanisms:**

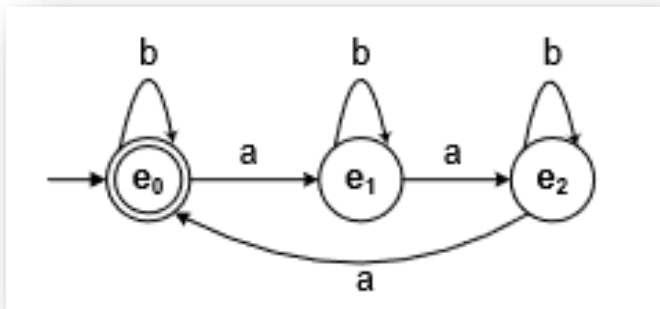
- called grammar or rewriting systems: they define a mechanism for generating words in the form of inductive construction rules.

- The language  $L_3$  produced by the rules of the following grammar:  $S \rightarrow aSb \mid S \rightarrow ab$

# Description of language

## Recognition mechanisms:

- also called automata or machines: they allow determining whether a word belongs to the considered language or not
- Example :



$\# S_0 a \rightarrow \# a S_0$

$\# S_0 b \rightarrow \# b S_0$

$a S_0 a \rightarrow a a S_0$

$a S_0 b \rightarrow S_0$

$b S_0 a \rightarrow S_0$

$b S_0 b \rightarrow b b S_0$

$\# S_0 \rightarrow \# S_f$



*Tank you*

*Any Questions?*