

# Chapter 1

## Languages

While we may think that computers are just machines, there is an incredible work behind that aims to make some kind of language that computers can understand. It is well renewed that computers act depending on streams of binary code, but that same stream of binary code can be considered as a language, something that the computers "speak" and "think" with. Let's define more in detail this "computers' language".

### Alphabet and Strings

DEFINITION

An **alphabet**  $\Sigma$  is a **non-empty, finite set**, which contains elements called **symbols** (or **characters**). For instance, the following two sets are considered alphabets:

- $\Sigma = \{0, 1\}$
- $\Sigma = \{a, b, c, \dots, x, y, z\}$

A **string**  $w$  over an alphabet  $\Sigma$  is a **sequence of symbols**, all belonging to  $\Sigma$ , which are all written one after the other and aren't separated by other symbols

Strings also have different properties, such as a **length**, a **reverse** and the possibility to include one or more **substrings**:

- **Length**: defined as  $|w|$ , it denotes the **number of symbols** contained within  $w$ . If a string has length 0, then such string is called **empty string**, and is denoted with  $\epsilon$ ;
- **Reverse**: defined as  $w^R$ , the reverse is a string which contains all the symbols of  $w$  in the reverse order;
- **Substring**: we say that a string  $z$  is a substring of  $w$  if  $z$  appears consecutively within  $w$ .

💡 1.0.0

EXAMPLE

For instance, the string  $w = 00101$  has a length of 5, its reverse is  $w^R = 10100$ , and has, as possible substrings the following:  $001, 10, \epsilon$ . Its alphabet is  $\Sigma = \{0, 1\}$

When we have two strings  $x$  and  $y$ , we can **concatenate** them by appending  $y$  at the end of  $x$ . The result is denoted as  $xy$ . Clearly, the length of the concatenated string is equal to the sum of the length of the two strings ( $|xy| = |x| + |y|$ ). A specific notation  $x^k$  denotes the concatenation of string  $x$  with itself for  $k$  times. So, if  $x$  was for instance "01", then we would have that

$$x^0 = \epsilon \quad x^1 = 01 \quad x^2 = 0101 \quad x^3 = 010101 \quad \text{and so on...}$$

By combining the definitions of concatenation and substring, we can define properly what a **prefix** and **suffix** is:

- **Prefix**: we say that a string  $x$  is a prefix of a string  $z$  if there exists a string  $y$  such that  $z = xy$ . Moreover, we say that  $x$  is a **proper prefix** of  $z$  if, additionally,  $x \neq z$ , so if  $y \neq \epsilon$ ;

- **Suffix:** we say that a string  $x$  is a suffix of a string  $z$  if there exists a string  $y$  such that  $z = yx$ . Moreover, we say that  $x$  is a **proper suffix** of  $z$  if, additionally,  $x \neq z$ , so if  $y \neq \epsilon$ .

Given these tools, we can now define what a language is:

### Language

A **language** is a **set of strings**. It is also defined as **prefix-free** if no member is a proper prefix of any other member

Languages can follow an order. There are different types of orders, one of which is the **lexicographic order**, which is defined by the order of the alphabet; we can think of it as the familiar dictionary order. Another order is the **shortlex order** (or **string order**), which is the same of the lexicographic order with the exception that shorter strings precede longer strings.

## 1.1 Deterministic Finite Automata (DFA)

We know that sometimes circuits can use the output from a given combination of inputs as input for the next output. In that case, we talk about circuits with states, or more specifically, we talk about **automata**.

### Deterministic Finite Automaton (DFA)

We define a **finite automaton** as a tuple of 5 elements  $(Q, \Sigma, \delta, q_0, F)$ , where:

- $Q$  is a finite set, denoting the **states** that the automaton can reach;
- $\Sigma$  is a finite set, denoting the **alphabet** of the automaton;
- $\delta$  is a function, called **transition function**, which given  $Q$  and  $\Sigma$  returns the set of next states  $Q'$ ;

$$\delta : Q \times \Sigma \longrightarrow Q'$$

- $q_0 \in Q$  is a state, and denotes the **starting state** of the automaton;
- $F \subseteq Q$  is a finite set, denoting the set of **accepted states** (or **final states**). An accepted state tells the automata if a string can be accepted or not.

Let's make an example:

#### 1.1.1

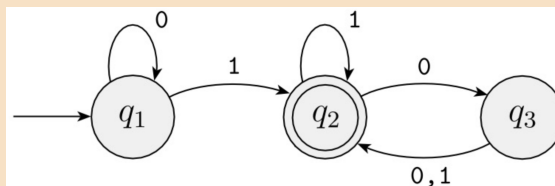
Consider the following tuple:

$$(\{q_1, q_2, q_3\}, \{0, 1\}, \delta, q_1, q_2)$$

It denotes that we have 3 possible states  $\{q_1, q_2, q_3\}$ , the alphabet is made of 2 symbols,  $\{0, 1\}$ , the starting state is  $q_1$ , the final state is  $q_2$  and the transition function is  $\delta$ . Such transition function could be the following:

	0	1
$q_1$	$q_1$	$q_2$
$q_2$	$q_3$	$q_2$
$q_3$	$q_2$	$q_2$

We can represent it visually as follows:



Automata change states depending on the string that they receive in input. So when a string  $x$  is given to the automaton, it goes through the string symbol by symbol, and changes states depending on the input characters read and the result of the transition function. If, after reading the whole string, the automaton arrives at a final input, then we say that it **accepts** the string, otherwise it **rejects** it.

### Language of an automaton and Regular Language

DEFINITION

The **language of an automaton**  $M$ , denoted as  $L(M)$ , is the **set** of all those **input strings** which make the automaton **go** from the starting state **to** a **final state**.

More specifically: given an automaton  $M = (Q, \Sigma, \delta, q_0, F)$  and let  $w$  be a string of  $n$  symbols where  $w_i$  denotes the  $i^{\text{th}}$  symbol of the string; let also that

$$\forall w_i \in w, \text{ then } w_i \in \Sigma$$

Then, we say that  $M$  accepts  $w$  if there exists a set of states  $r = \{r_0, r_1, \dots, r_n\} \in Q$  which respects all the following three conditions:

- $r_0 = q_0$ : the first state of  $r$  must coincide with the initial state;
- $\delta(r_i, w_{i+1}) = r_{i+1}, \forall i \in [0, 1, \dots, n-1]$ : each character must lead to a transition of state;
- $r_n \in F$ : the final state in  $r$  should be in  $F$ .

$M$  recognizes a language  $A$  If

$$A = \{w \mid M \text{ accepts } w\}$$

A **language**  $L$  is said to be **regular** if there exists a finite automaton  $M$  which accepts  $L$

Note that it's possible for an automaton to have that  $q_0 \in F$ .

## 1.2 Non-Deterministic Finite Automata (NFA)

With DFA, we always knew that, given a current state  $q_i$  and a character  $w_i \in \Sigma$ , then we would **always know** what the **next state** would be. There is another type of automata which breaks this deterministic rule, and is called **non-deterministic finite automata**.

### Non-Deterministic Finite Automaton (NFA)

DEFINITION

A **non-deterministic finite automaton (NFA)** is an automaton denoted by a tuple of 5 elements  $(Q, \Sigma, \delta, q_0, F)$  where:

- $Q$  is a finite set of **states** that the automaton can reach;
- $\Sigma$  denotes a finite **alphabet** recognized by the automaton;
- $\delta$  is the **transition function**, which given  $Q$  and  $\Sigma_c$  returns the **power set** of  $Q$  (denoted by  $\mathcal{P}(Q)$ );

$$\delta : Q \times \Sigma_c \longrightarrow \mathcal{P}(Q)$$

- $q_0 \in Q$  is the **starting state** of the automaton;
- $F \subseteq Q$  is the set of **accepted states**.

The power set  $\mathcal{P}(Q)$  of a set  $Q$  is the set of all subsets of  $Q$ , for any  $Q$ ; the alphabet denoted with  $\Sigma_c$  is equal to  $\Sigma \cup \{c\}$ .