

# Formal Language Theory

## an Introduction

*Prof. A. Morzenti*

ALPHABET  $\Sigma$  : any ***finite*** set of symbols  $\Sigma = \{a_1, a_2, \dots a_k\}$

cardinality of the alphabet  $|\Sigma| = k$

String: a sequence ( $\Rightarrow$ ordered) of alphabet elements (possibly repeated)

Language: any set of strings

$$\Sigma = \{a, b, c\} \quad L_1 = \{ab, ac\} \quad L_2 = \{bc, bbc\} \quad L_3 = \{abc, aabbcc, aaabbbccc, \dots\}$$

The strings of a language are called its ***sentences*** or ***phrases***

Language ***cardinality***: the number of its sentences

$$|L_2| = |\{bc, bbc\}| = 2 \qquad |\emptyset| = 0$$

Number of occurrences of a symbol in a string  $|bbc|_b = 2, \quad |bbc|_a = 0$

With a slight *abuse of notation* sometimes we denote with  $\Sigma$   
both the alphabet and  
the language of all strings of length 1

*length* of a string  $x$ :  $|x|$   
number of its elements

$$\begin{array}{l} |bbc| = 3 \\ |abbc| = 4 \end{array}$$

*string equality* : two strings are equal if and only if (*iff*, for short)

- have the same length
- their elements, from left to right, coincide

$$x = a_1 a_2 \dots a_h \quad y = b_1 b_2 \dots b_k$$

$$x = y \text{ iff } h = k \text{ and } a_i = b_i \text{ for all } i = 1 \dots h$$

$$bbc \neq bcb \neq bc$$

# OPERATIONS ON STRINGS /1

CONCATENATION (product):  $x \cdot y$  or  $xy$  for short

$$x = a_1 a_2 \dots a_h \quad y = b_1 b_2 \dots b_k \quad x \cdot y = xy = a_1 a_2 \dots a_h b_1 b_2 \dots b_k$$

- associative  $(xy)z = x(yz)$

- length  $|xy| = |x| + |y|$

EMPTY STRING (or ***null string***)  $\varepsilon$  is the neutral element for concatenation

for any  $x$ ,  $x\varepsilon = \varepsilon x = x$

length of  $\varepsilon$ :  $|\varepsilon| = 0$

NOTICE:  $\varepsilon$  is ***NOT*** the empty set:  $\varepsilon \neq \emptyset$

SUBSTRINGS: if  $x = uyv$  (NB: both  $u$  and  $v$  can be  $\varepsilon$ ) then

- $y$  is a substring of  $x$
- $y$  is a ***proper substring*** iff  $u \neq \varepsilon$  or  $v \neq \varepsilon$
- $u$  is a ***prefix*** of  $x$
- $v$  is a ***suffix*** of  $x$

## EXAMPLES

if  $x = abccbc$  then

prefixes:  $a, ab, abc, abcc, abccb, abccbc$

suffixes:  $c, bc, cbc, cc bc, bcc bc, abcc bc$

substrings:  $\dots, bc, cc, cb, abc, bcc, \dots$

## OPERATIONS ON STRINGS /2

REFLECTION  $x^R$

$$\begin{aligned} x &= a_1 a_2 \dots a_h \\ x^R &= a_h a_{h-1} \dots a_2 a_1 \\ (x^R)^R &= x \\ (xy)^R &= y^R x^R \\ \varepsilon^R &= \varepsilon \end{aligned}$$

$$\begin{aligned} x &= atri & x^R &= irta \\ x &= bon & y &= ton \\ xy &= bonton \\ (xy)^R &= y^R x^R = notnob \end{aligned}$$

REPETITION:  $m$ -th power ( $m \geq 1$ ) of string  $x$ : concatenation of  $x$  with itself  $m-1$  times

$$\begin{aligned} x^m &= \underset{1 \ 2 \ 3 \ \dots \ m}{xxxx \dots x} \\ \left( \begin{array}{l} \text{inductive} \\ \text{definition} \end{array} \right) & \left\{ \begin{array}{l} x^m = x^{m-1} x, \quad m > 0 \\ x^0 = \varepsilon \end{array} \right. \end{aligned}$$

$$\begin{aligned} x &= ab & x^0 &= \varepsilon & x^1 &= x = ab & x^2 &= (ab)^2 = abab \\ y &= a^3 = aaa & y^3 &= a^3 a^3 a^3 = a^9 \\ \varepsilon^0 &= \varepsilon & \varepsilon^2 &= \varepsilon \end{aligned}$$

OPERATOR PRECEDENCE: repetition and reflection take precedence over concatenation

$$\begin{aligned} ab^2 &= abb & (ab)^2 &= abab \\ ab^R &= ab & (ab)^R &= ba \end{aligned}$$

# OPERATIONS ON LANGUAGES /1

OPERATIONS ARE TYPICALLY DEFINED ON A LANGUAGE  
BY EXTENDING THE STRING OPERATION TO ALL ITS PHRASES

REFLECTION  $L^R$ :  $L^R = \{ x \mid \exists y (y \in L \wedge x = y^R) \}$  def. by *characteristic predicate*

$\text{Prefixes}(L) = \{ y \mid y \neq \varepsilon \wedge \exists x \exists z (x \in L \wedge x = yz \wedge z \neq \varepsilon) \}$  NB: *proper* prefixes

*Prefix-free language*  $L$ : no proper prefix of its phrases  $\in L$ :  $\text{Prefixes}(L) \cap L = \emptyset$

EXAMPLE:  $L_1 = \{ x \mid x = a^n b^n \wedge n \geq 1 \}$  is prefix-free:  $a^2 b^2 \in L_1$   $a^2 b \notin L_1$

EXAMPLE:  $L_2 = \{ x \mid x = a^m b^n \wedge m > n \geq 1 \}$  is not prefix-free:  $a^4 b^3 \in L_2$   $a^4 b^2 \in L_2$

# OPERATIONS ON LANGUAGES / 2

## Operations defined over two arguments

### CONCATENATION

$$L' L'' = \{xy \mid x \in L' \wedge y \in L''\}$$

$m$ -th POWER  
(inductive definition)

$$L^m = L^{m-1} L, m > 0$$
$$L^0 = \{\varepsilon\}$$

NB:  $\{\varepsilon\} \neq \emptyset$

NB: consequences

$$\emptyset^0 = \{\varepsilon\} \quad L.\emptyset = \emptyset.L = \emptyset \quad L.\{\varepsilon\} = \{\varepsilon\}.L = L$$

# OPERATIONS ON LANGUAGES / 3

## EXAMPLES

$$L_1 = \{ a^i \mid i \geq 0, i \text{ even} \} = \{ \varepsilon, a^2, a^4, \dots \}$$

$$L_2 = \{ b^j a \mid j \geq 1, j \text{ odd} \} = \{ ba, b^3 a, b^5 a, \dots \}$$

$$\begin{aligned} L_1 L_2 &= \{ a^i b^j a \mid (i \geq 0, i \text{ even}) \wedge (j \geq 1, j \text{ odd}) \} = \\ &= \{ \varepsilon ba, a^2 ba, a^4 ba, \dots \varepsilon b^3 a, a^2 b^3 a, \dots \} \end{aligned}$$

$$\begin{aligned} (L_1)^2 &= \{ \varepsilon, a^2, a^4, a^6, \dots \} \{ \varepsilon, a^2, a^4, a^6, \dots \} = \\ &= \{ \varepsilon, \varepsilon a^2, \varepsilon a^4, \dots, a^2 \varepsilon, a^4, \dots, a^4 \varepsilon, a^6 \dots \} = L_1 \end{aligned}$$

for each pair of even numbers  
 $h$  and  $k$ ,  $h+k$  is even, hence  
 $a^{h+k} \in L_1$

PAY ATTENTION: the language  $L^m$   
in general does **not** contain **only**  
phrases of  $L$  repeated  $m$  times

$$\begin{aligned} \{ x \mid x = y^m \wedge y \in L \} &\subset L^m \\ m = 2 \quad L_1 &= \{ a, b \} \\ \{ a^2, b^2 \} &\subset L_1^2 = \{ a^2, ab, ba, b^2 \} \end{aligned}$$



## OPERATIONS ON LANGUAGES / 4

### Finite length strings:

The power operator allows one to define concisely the language of strings whose length is not greater than a given integer  $K$

$$L = \{\varepsilon, a, b\}^3 \quad K = 3$$

$$L = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, \dots bbb\}$$

Notice the role of  $\varepsilon$

It allows one to obtain

all strings of length  $< K$  (0, 1, 2)

To rule out the empty string:

$$L = \{a, b\} \{ \varepsilon, a, b \}^2$$

## OPERATIONS ON LANGUAGES / 5

**SET THEORETIC OPERATIONS:** the customary ones are defined:  
union, intersection, difference, inclusion, strict inclusion, equality

$$\cup \quad \cap \quad \setminus \quad \subseteq \quad \subset \quad =$$

**UNIVERSAL LANGUAGE:** the set of all  
strings over the alphabet  $\Sigma$ ,  
of any length, including 0 (i.e., string  $\varepsilon$ )

$$L_{universal} = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \dots$$

**COMPLEMENT** of a language  $L$  over alphabet  $\Sigma$   
is the set difference with respect to (w.r.t.) the  
universal language (i.e., the set of strings over  $\Sigma$  that  $\notin L$ )

$$\neg L = L_{universal} \setminus L$$

hence 
$$L_{universal} = \neg \emptyset$$

# OPERATIONS ON LANGUAGE / 6

## EXAMPLES

The complement of a *finite* language  
is *always infinite*

$$\neg(\{a, b\}^2) = \varepsilon \cup \{a, b\} \cup \{a, b\}^3 \cup \dots$$

The complement of an *infinite* one  
is *not necessarily finite*

$$L = \{a^{2n} \mid n \geq 0\} \quad \neg L = \{a^{2n+1} \mid n \geq 0\}$$

Examples of the difference operation among languages

$$\begin{aligned} \Sigma &= \{a, b, c\} \\ L_1 &= \{x \mid |x|_a = |x|_b = |x|_c \geq 0\} \\ L_2 &= \{x \mid |x|_a = |x|_b \wedge |x|_c = 1\} \end{aligned}$$

$$L_1 \setminus L_2 = \varepsilon \cup \{x \mid |x|_a = |x|_b = |x|_c \geq 2\}$$

(same number of  $a$ ,  $b$ ,  $c$ , but not =1)

$$L_2 \setminus L_1 = \{x \mid |x|_a = |x|_b \neq |x|_c = 1\}$$

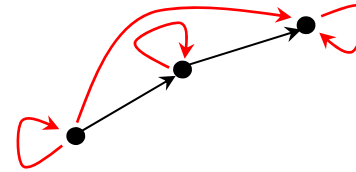
(only one  $c$  and same number of  $a$  and  $b$ , but  $\neq 1$ )

## A frequently used algebraic operation: reflexive and transitive closure $R^*$ of a relation $R$

Given a set  $A$  and a relation  $R \subseteq A \times A$ , the condition  $(a_1, a_2) \in R$  is also denoted as  $a_1 R a_2$

$R^*$  is a *relation* defined by:

$$x R^* y \text{ iff } \exists x_0, x_1, \dots, x_n, n \geq 0, \\ \text{s.t. (s.t.} \equiv \text{such that)} \quad x = x_0, \quad y = x_n, \quad \text{and} \quad \forall i = 1 \dots n \quad x_{i-1} R x_i$$



If we see  $a R b$  as *a step* in relation  $R$ ,  $x R^* y$  can be seen as *a chain of  $n \geq 0$  steps*

Similarly, **transitive closure** (non reflexive)  $R^+$ : idem with  $n \geq 1$   
 **$k$ -th power  $R^k$** : with  $n = k$

Example: if relation  $R$  is the *adjacency* relation on a graph  
 $R^*$  is the *reachability* relation

Similarly, **closure**  $A$  of a *set*  $A$  under an *operation* (function)  
 is obtained from  $A$  by adding all elements obtained  
 by applying the operation any number of times

# OPERATIONS ON LANGUAGES / 7

STAR OPERATOR: reflexive transitive closure under the concatenation operation  
(also called *Kleene star*)

$$L^* = \bigcup_{h=0 \dots \infty} L^h = L^0 \cup L^1 \cup L^2 \dots = \varepsilon \cup L^1 \cup L^2 \dots$$

$$L = \{ab, ba\} \quad L^* = \{\varepsilon, ab, ba, abab, abba, baab, baba, \dots\}$$

( $L$  is finite       $L^*$  is infinite)

It is the union of all the powers of the language

Every string of the star language  $L^*$  can be chopped into substrings  $\in L$

The star language  $L^*$  can be equal to the base language  $L$

$$L = \{a^{2n} \mid n \geq 0\} \quad L^* = \{a^{2n} \mid n \geq 0\} \equiv L$$

## OPERATIONS ON LANGUAGES / 8

If we take  $\Sigma$  as the base language, then  $\Sigma^*$  contains all the strings built on that alphabet (it is the *universal language* of alphabet  $\Sigma$  )

We often say that  $L$  is a language on alphabet  $\Sigma$  by writing  $L \subseteq \Sigma^*$

### PROPERTIES OF THE STAR OPERATOR

- monotonicity (with  $*$  the set increases):  $L \subseteq L^*$
- closure under concatenation: if  $x \in L^*$  and  $y \in L^*$  then  $xy \in L^*$
- idempotence:  $(L^*)^* = L^*$
- commutativity of star and reflection  $(L^*)^R = (L^R)^*$

Furthermore:  $\emptyset^* = \{ \varepsilon \}$        $\{ \varepsilon \}^* = \{ \varepsilon \}$     NB: these are cases where  $L^*$  is finite

Example of idempotence: We already noticed that, for  $L = \{ a^{2n} \mid n \geq 0 \}$ , it holds  $L^* = L$

This derives from idempotence, because we have  $L = L_0^*$  for  $L_0 = \{ aa \} = \{ a^2 \}$

## OPERATIONS ON LANGUAGES / 9

Example on the STAR OPERATOR

language of identifiers  $I$  as character strings that start with a letter and include any number of letters and digits

$$\Sigma_A = \{ a, b, \dots, z, A, B, \dots, Z \} \quad \Sigma_N = \{ 0, 1, 2, \dots, 9 \}$$

$$I = \Sigma_A (\Sigma_A \cup \Sigma_N)^*$$

if we stipulate  $\Sigma = \Sigma_A \cup \Sigma_N$

language  $I_5$  of identifiers of maximal length 5

$$I_5 = \Sigma_A ( \Sigma \cup \{ \epsilon \} )^4$$

## OPERATIONS ON LANGUAGES / 10

CROSS OPERATOR  $L^+$ : transitive closure (non reflexive) under concatenation

The union does *not* include the first power  $L^0$

Useful but not indispensable, it can be derived from the star operator  $*$ :

$$L^+ = L \cdot L^*$$

$$L^+ = \bigcup_{h=1 \dots \infty} L^h = L^1 \cup L^2 \cup \dots$$

$$\{ab, bb\}^+ = \{ab, bb, ab^3, b^2ab, abab, b^4, \dots\}$$

$$\{\varepsilon, aa\}^+ = \{\varepsilon, a^2, a^4, \dots\} = \{a^{2n} \mid n \geq 0\}$$

if  $\varepsilon \in L$  then  $L^+ = L^*$

Typically, a given language can be defined in different ways using different operators

Example: language  $L$  of strings of length  $\geq 4$ :  $L = \Sigma^4 \Sigma^*$  and also  $L = (\Sigma^+)^4$



## OPERATIONS ON LANGUAGES / 11

**QUOTIENT OPERATOR**  $L_1 / L_2$ : it shortens the phrases of  $L_1$  by cutting off a suffix that belongs to  $L_2$ . NB: **forward** slash (backward slash denotes set difference)

$$L = L_1 / L_2 = \{ y \mid \exists x \in L_1 \exists z \in L_2 (x = yz) \}$$

Example:  $L_1 = \{a^{2n} b^{2n} \mid n > 0\}$        $L_2 = \{b^{2n+1} \mid n \geq 0\}$

$$\begin{aligned} L_1 / L_2 &= \{a^r b^s \mid (r \geq 2, \quad r \text{ even}) \wedge (1 \leq s < r, \quad s \text{ odd})\} \\ &= \{a^2 b, a^4 b, a^4 b^3, \dots\} \end{aligned}$$

$L_2 / L_1 = \emptyset$       because no string in  $L_2$  has a string in  $L_1$  as a suffix