# Theory of Formal Languages Basic Notions

Translated and adapted by L. Breveglieri

#### BASICS / 1

ALPHABET: finite set of elements cardinality string, word, phrase

STRING: ordered set of atomic elements, possibly repeated

LANGUAGE: finite or infinite set of strings

The set-theoretic structure of a language has two levels

LANGUAGE: unordered set L of non-atomic elements (strings) that are in turn sequences of atomic elements (characters or terminals) cardinality of language L finite or infinite language

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

$$\sum = \{a_1, a_2, ..., a_k\}$$

$$|\Sigma| = k$$

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

$$abc, aabc, ac, bbb$$

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

$$L_1 = \{ab, ac\}$$

$$L_2 = \{bc, bbc\}$$

 $L_1 = \{abc, aabbcc, abcabc, ...\}$ 

$$|bbc|_b = 2, |bbc|_a = 0$$

#### BASICS / 2

LENGTH OF A STRING x:  $|x| \ge 0$  is the number of elements (letters)

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

EQUALITY OF TWO STRINGS: two strings are equal if and only if (iff) they have the same length their elements orderly coincide, from left to right

$$x = a_1 a_2 ... a_h, y = b_1 b_2 ... b_k$$

$$x = y \quad \text{if} \quad h = k$$

$$a_i = b_i \quad \text{with} \quad i = 1, ... h;$$

$$bbc \neq bcb \neq bc$$

#### **OPERATIONS ON STRINGS / 1**

## is a basic operation is associative changes the length

EMPTY STRING (or NULL string)  $\epsilon$  is the neutral element with respect to concatenation: chaining  $\epsilon$  on the left or right does not change the string Pay attention:  $\epsilon$  is NOT the same as  $\Phi$  (the empty set)!

$$x = a_1 a_2 ... a_h, y = b_1 b_2 ... b_k$$

$$x. y = a_1 a_2 ... a_h b_1 b_2 ... b_k = xy$$

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

$$|xyz| = |x| + |y| + |z|$$

$$x\varepsilon = \varepsilon x = x$$

$$|\varepsilon| = 0$$

Proper substring: y if u,  $v \neq \epsilon$ 

Start 
$$_{k}(x) = k : x$$

$$x = abccbc$$
  
 $p$  prefix  $a, ab, abc, abcc, abccb, abccbc$   
 $si$  suffix  $c, bc, cbc, ccbc, bccbc, abccbc$   
 $si$  substring ....,  $bc, cc, cb, ...$ 

#### **OPERATIONS ON STRINGS / 2**

MIRRORING or REFLECTION

$$x = atri$$
  $x^{R} = irta$   
 $x = bon$   $y = ton$   
 $xy = bonton$   
 $(xy)^{R} = y^{R}x^{R} = notnob$ 

$$x = a_1 a_2 ... a_h$$

$$x^R = a_h a_{h-1} ... a_2 a_1$$

$$(x^R)^R = x$$

$$(xy)^R = y^R x^R$$

$$\varepsilon^R = \varepsilon$$

REPETITION (or ITERATION): the m-th power of a string (with  $m \ge 1$ ) concatenates the string to itself for m-1 times

$$x = ab \quad x^{0} = \varepsilon \quad x^{1} = x = ab \quad x^{2} = (ab)^{2} = abab$$

$$y = a^{3} = aaa \quad y^{3} = a^{3}a^{3}a^{3} = a^{9}$$

$$\varepsilon^{0} = \varepsilon \quad \varepsilon^{2} = \varepsilon$$

$$x^{m} = \underset{123 \dots m}{xxx...x}$$

$$x^{m} = x^{m-1}x, \quad m > 0$$

$$x^{0} = \varepsilon$$

PRECEDENCE BETWEEN OPERATORS power precedes concatenation mirroring precedes concatenation

$$ab^{2} = abb$$
  $(ab)^{2} = abab$   
 $ab^{R} = ab$   $(ab)^{R} = ba$ 

An operation defined on a language applies to each string in the language (and needs to be definable over any string)

$$L^R = \left\{ x \mid x = y^R \land y \in L \right\}$$
 characteristic predicate prefix  $L(L) = \left\{ y \mid x = yz \land x \in L \land y, z \neq \varepsilon \right\}$ 

PREFIX-FREE LANGUAGE: in the language there is not any string that is a prefix of another string of the language

Equivalently, prefix(L) and L are disjoint sets (i.e., prefix(L)  $\cap$  L =  $\Phi$ )

$$\begin{aligned} L_1 &= \left\{x \mid x = a^n b^n \land n \geq 1\right\} &\quad a^2 b^2 \in L_1 \quad a^2 b \not\in L_1 \\ \text{L}_1 &\text{ is prefix free} & \text{prefixes are } a^n b^m \quad \text{where } n > m \geq 0 \end{aligned}$$
 
$$L_2 &= \left\{a^m b^n \mid m > n \geq 1\right\} \quad a^4 b^3 \in L_2 \quad a^4 b^2 \in L_2 \\ \text{L}_2 &\text{ is not prefix-free} \end{aligned}$$

Caution:  $\epsilon$  is prefix (or suffix, or substing) of any other string, including itself

### OPERATIONS ON LANGUAGES / 2 binary (two arguments) operations

CONCATENATION

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

m-th POWER ( $m \ge 0$ )

$$egin{aligned} L^m &= L^{m-1}L, m > 0 \ L^0 &= \left\{ \varepsilon \right\} \end{aligned}$$

Pay attention to the following consequences

$$\varnothing^0 = \{\varepsilon\}$$
  $L.\varnothing = \varnothing.L = \varnothing$   $L.\{\varepsilon\} = \{\varepsilon\}.L = L$ 

#### **EXAMPLES**

$$L_{1} = \{a^{i} \mid i \geq 0, even \} = \{\varepsilon, a^{2}, a^{4}, a^{6}, ...\}$$

$$L_{2} = \{b^{j}a \mid j \geq 1, odd \mid j = \{ba, b^{3}a, b^{5}a, ...\}$$

$$L_{1}L_{2} = \{a^{i}b^{j}a \mid (i \geq 0, even) \land (j \geq 1, odd)\}$$

$$= \{\varepsilon ba, a^{2}ba, a^{4}ba, ...\varepsilon b^{3}a, a^{2}b^{3}a, ...\}$$

$$(L_{1})^{2} = \{\varepsilon, a^{2}, a^{4}, a^{6}, ...\} \{\varepsilon, a^{2}, a^{4}, a^{6}, ...\} =$$

$$= \{\varepsilon, \varepsilon a^{2}, \varepsilon a^{4}, ..., a^{2} \varepsilon, a^{4}, ..., a^{4} \varepsilon, a^{6} ...\} = L_{1}$$

For every pair of even integers h and k, h + k is even and  $a^{h+k}$  belongs to L<sub>1</sub>

CAUTION

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

STRINGS OF FINITE LENGTH: the power operator allows us to expressively define the language of the strings that have a length not greater than (= less than or equal to) a given fixed integer k

$$L = \{\varepsilon, a, b\}^3 \quad k = 3$$
$$L = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, ...bb\}$$

Notice the role of  $\epsilon$ , which allows us to obtain all the strings of length 0, 1, 2

$$egin{array}{l} \{arepsilon,a,b\} \ \{arepsilon,a,b\} \ \ \{arepsilon,a,b\} \end{array}$$

And in order to exclude the empty string ε, do as follows

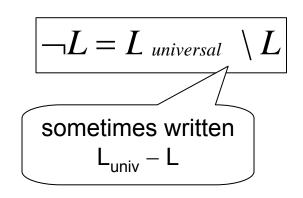
$$L = \{a,b\}\{\varepsilon,a,b\}^2$$

SET-THEORETIC OPERATIONS: these are the traditional operations of elementary set theory: union  $\cup$ , intersection  $\cap$ , complement  $\neg$  (or overlining  $\overline{\phantom{a}}$ ) and the traditional relational operators between sets: strict inclusion  $\subset$ , inclusion  $\subseteq$ , equality  $\neq$ , etc

UNIVERSAL LANGUAGE: the set of ALL the strings defined over the alphabet Σ, of any length (including also length 0)
Also sometimes called the FREE MONOID

$$\begin{bmatrix} L_{universal} &= \sum^{0} \cup \sum^{1} \cup \sum^{2} \cup ... \\ L_{universal} &= \neg \varnothing$$

COMPLEMENT of a language L over the alphabet  $\Sigma$ : it si defined as the set-theoretic difference with respect to the universal language over  $\Sigma$  Equivalently, it is the set of all the strings over the alphabet  $\Sigma$  that do not belong to L



#### **EXAMPLES**

The complement of a finite language is always an infinite language.
The complement of an infinite language may be infinite, but not necessarily (sometimes it happens to be finite)

Set-theoretic difference

$$\sum = \{a, b, c\} 
L_1 = \{x \mid |x|_a = |x|_b = |x|_c \ge 0\} 
L_2 = \{x \mid |x|_a = |x|_b \land |x|_c = 1\}$$

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

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

sometimes written  $L_1 - L_2$ 

In both natural and artificial languages, the phrases can be of any length But only formulas of finite length can be written to define a language It is necessary to introduce some operators to create infinitely many strings

STAR OPERATOR (also called Kleene star or concatenation closure): it is the limit of the power operator

The union of all the powers of a language, for every positive or null exponent

all the powers of a language, for every positive of null exp 
$$L^* = \bigcup_{h=0...\infty} L^h = L^0 \cup L^1 \cup L^2... = \varepsilon \cup L^1 \cup L^2...$$
 
$$L = \{ab, ba\} \quad L^* = \{\varepsilon, ab, ba, abab, abba, baab, baba,...\}$$
 L is finite but L\* is infinite

Every string in the star language of L can be factored into substrings, each of which belongs to the language L

Sometimes, the star language happens to be identical to the base language

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

If one takes the alphabet  $\Sigma$  as the base language,  $\Sigma^*$  contains all strings  $(\Sigma^*)$  is the universal language over the alphabet  $\Sigma$ ). One may signify that L is a language over the alphabet  $\Sigma$  by writing as follows:

#### PROPERTIES OF THE STAR OPERATOR monotonic

closed w.r.t. concatenation idempotent

commutes with mirroring

Moreover

$$L \subseteq L^*$$
if  $\left(x \in L^* \land y \in L^*\right)$  then  $xy \in L^*$ 

$$\left(L^*\right)^* = L^*$$

$$\left(L^*\right)^R = \left(L^R\right)^*$$

$$\varnothing^* = \{\varepsilon\}$$
  $\{\varepsilon\}^* = \{\varepsilon\}$ 

$$L_1 = \left\{ a^{2n} \mid n \ge 0 \right\} \qquad L_1^* = L_{1_*}$$

$$\text{moreover } L = \left\{ aa \right\}^*$$

Example of star operator: an identifier, modeled as a string of letters and digits (alphanumeric), of arbitrary length (not null), but starting with a letter (not with a digit)

$$\begin{split} & \sum_{A} = \left\{A, B, ..., Z\right\} \quad \sum_{N} = \left\{0, 1, 2, ..., 9\right\} \\ & I = \sum_{A} \left(\sum_{A} \bigcup \sum_{N}\right)^{*} \\ & \text{if} \quad \sum = \sum_{A} \bigcup \sum_{N} \\ & I_{5} = \sum_{A} \left(\sum^{0} \bigcup \sum^{1} \bigcup \sum^{2} \bigcup \sum^{3} \bigcup \sum^{4}\right) \\ & I_{5} = \sum_{A} \left(\sum \bigcup \varepsilon\right)^{4} \end{split}$$

The C language would admit the underscore "\_" as well, but not as the initial symbol. Extend the definition (do it yourself)

CROSS OPERATOR (also called Kleene cross or  $\varepsilon$ -free concatenation closure): is the non-reflexive closure with respect to concatenation (see below)

The unitory does not contain the null power

Sometimes very useful, but not indispensable

$$L^{+} = \bigcup_{h=1...\infty} L^{h} = L^{1} \bigcup L^{2} \bigcup ...$$

$$\left\{ab, bb\right\}^{+} = \left\{ab, bb, ab^{3}, b^{2}ab, abab, b^{4}, ...\right\}$$

$$\left\{\varepsilon, aa\right\}^{+} = \left\{\varepsilon, a^{2}, a^{4}, ...\right\} = \left\{a^{2n} \mid n \ge 0\right\}$$

The same language can be defined in different ways by different combinations of the same or other operators

Example: the strings of length greater than or equal to 4

$$\left| \sum^{4} \sum^{*} (\sum^{+})^{4} \right|$$

QUOTIENT OPERATOR: it shortens the phrases of a language L', by stripping off a suffix out of another language L''

$$L = L'/L'' = \{ y \mid (x = yz \in L') \land z \in L'' \}$$

#### Example of quotienting

$$L' = \{a^{2n}b^{2n} \mid n > 0\}, \quad L'' = \{b^{2n+1} \mid n \ge 0\}$$

$$L'/L'' = \{a^{r}b^{s} \mid (r \ge 2 \text{ even}) \land (1 \le s < r, s \text{ odd}) \}$$

$$= \{a^{2}b, a^{4}b, a^{4}b^{3}, ...\}$$

$$L''/L' = \emptyset$$

Question: what happens if  $x \in L'$  does not admit any suffix  $z \in L''$ ?