# Theory of Formal Languages
# Introduction

*Prof. Licia Sbattella*

*aa 2007-08*

*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 of non-atomic elements that are in turn ordered sets of atomic elements:
- cardinality of L
- finite or infinite language

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

$$|\Sigma| = k$$

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

$$abc, aabc, ac, bbb$$

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

$$L_1 = \{a\,b, a\,c\}$$

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

$$L_1 = \{a\,bc, aabbcc, abcabc, ...\}$$

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

# BASICS / 2

LENGTH OF A STRING x: |x|
is the number of elements (letters).

$$\left|bbc\right| = 3$$

$$\left|abbc\right| = 4$$

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

CONCATENATION (product of strings):
- •is a basic opertion
- •is associative
- •changes the length

EMPTY STRING (or NULL string):
ε is the neutral element with respect
to concatenation: chaining ε  on the left
or right does not change the string

$P$ay attention: ε is NOT the same as $\Phi$
(the empty set) !

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

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

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

$$\left| xyz \right| = \left| x \right| + \left| y \right| + \left| z \right|$$

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

$$\left| \varepsilon \right| = 0$$

SUBSTRING: x = u y v
      y is a substring
      u is a prefix
      v is a suffix

$P$roper substring: y if u, v $\neq$ ε

$S$tart $_k$ (x) = k : x

$$x = abccbc$$

$p$ prefix $\quad \imath, ab, abc, abcc, abccb, abccbc$

$s\iota$ suffix $\quad c, bc, cbc, ccbc, bccbc, abccbc$

$s\iota$ substring $\quad ....., bc, cc, cb, ...$

# OPERATIONS ON STRINGS / 2

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

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

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

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

$$\varepsilon^R = \varepsilon$$

MIRRORING or
REFLECTION

$$x = atri \qquad x^R = irta$$

$$x = bon \qquad y = ton$$

$$xy = bonton$$

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

REPETITION (or ITERATION): the $m$-th power of a string (where $m$ is greater than or equal to 1) consists of concatenating the string to itself for $m - 1$ times.

$$x^m = \underset{1\,2\,3\,...\quad m}{xxx...x}$$

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

$$x^0 = \varepsilon$$

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

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

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

PRECEDENCE AMONG OPERATORS:
- power precedes concatenation
- mirroring precedes concatenation

$$ab^2 = abb \qquad (ab)^2 = abab$$

$$ab^R = ab \qquad (ab)^R = ba$$

# OPERATIONS ON LANGUAGES / 1

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

$$L^R = \left\{ x \mid x = y^R \wedge y \in L \right\}$$

characteristic predicate

$$\text{prefix } (L) = \left\{ y \mid x = yz \wedge x \in L \wedge y, z \neq \varepsilon \right\}$$

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

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

$$L_1 = \left\{ x \mid x = a^n b^n \wedge n \geq 1 \right\} \quad a^2 b^2 \in L_1 \quad a^2 b \notin L_1$$

$L_1$ is prefix free    prefixes are $a^n b^m$   where $n > m \geq 0$

$$L_2 = \left\{ a^m b^n \mid m > n \geq 1 \right\} \quad a^4 b^3 \in L_2 \quad a^4 b^2 \in L_2$$

$L_2$ is not prefix-free

Caution: $\varepsilon$ is prefix (or suffix, or substing) to any other string, including itself.

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

CONCATENATION:

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

*m*-th POWER ($m \geq 0$)

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

$$L^0 = \left\{ \varepsilon \right\}$$

Pay attention to the following consequences:

$$\varnothing^0 = \left\{ \varepsilon \right\} \quad L.\varnothing = \varnothing.L = \varnothing \quad L.\left\{ \varepsilon \right\} = \left\{ \varepsilon \right\}.L = L$$

# OPERATIONS ON LANGUAGES / 3

## EXAMPLES:

$$L_1 = \left\{ a^i \mid i \geq 0, \; even \right\} = \left\{ \varepsilon, a^2, a^4, a^6, ... \right\}$$

$$L_2 = \left\{ b^j a \mid j \geq 1, \; odd \right\} = \left\{ ba, b^3 a, b^5 a, ... \right\}$$

$$L_1 L_2 = \left\{ a^i b^j a \mid (i \geq 0, \; even) \wedge (j \geq 1, \; odd) \right\}$$

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

$$(L_1)^2 = \left\{ \varepsilon, a^2, a^4, a^6, ... \right\} \left\{ \varepsilon, a^2, a^4, a^6, ... \right\} =$$

$$= \left\{ \varepsilon, \varepsilon a^2, \varepsilon a^4, ..., a^2 \varepsilon, a^4, ..., a^4 \varepsilon, a^6 ... \right\} = L_1$$

For every pair of even integers $h$ and $k$, $h + k$ is even and $a^{h+k}$ belongs to $L_1$.

CAUTION:

$$\left\{ x \mid x = y^m \wedge y \in L \right\} \subset L^m$$

$$m = 2 \quad L_1 = \left\{ a, b \right\}$$

$$\left\{ a^2, b^2 \right\} \subset L_1^2 = \left\{ a^2, ab, ba, b^2 \right\}$$

# OPERATIONS ON LANGUAGES / 4

STRINGS OF FINITE LENGTH: the power operator allows to define
in an expressive way the language of the strings having length not greater
than (= less than or equal to) a given fixed integer K.

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

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

Notice the role of ε,
that allows to obtain
all the strings of length
0, 1, 2.

$$\{\varepsilon, a, b\}$$

$$\{\varepsilon, a, b\}$$

$$\{\varepsilon, a, b\}$$

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

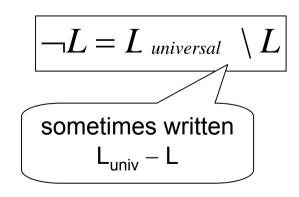$$L = \{a, b\}\{\varepsilon, a, b\}^2$$

# OPERATIONS ON LANGUAGES / 5

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

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

$$L_{universal} = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup ...$$
$$L_{universal} = \neg\varnothing$$

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.

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

sometimes written
$L_{univ} - L$

# OPERATIONS ON LANGUAGES / 6

## EXAMPLES

The complement of a finite language is always an infinite language.

The complement of an infinite language may be infinite,
but need not be always such
(sometimes happens to be finite).

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

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

Set-theoretic difference:

sometimes written $L_1 - L_2$

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

$$L_1 = \{x \mid |x|_a = |x|_b = |x|_c \geq 0\}$$

$$L_2 = \{x \mid |x|_a = |x|_b \wedge |x|_c = 1\}$$

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

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

# OPERATIONS ON LANGUAGES / 7

In both natural and artificial languages, the phrases can be of any length.

But only formulae 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.

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

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

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.

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

# OPERATIONS ON LANGUAGES / 8

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

PROPERTIES OF THE STAR OPERATOR:
- monotonic
- closed w.r.t. concatenation
- idempotent
- commutes with mirroring

Moreover:

$$L \subseteq L^*$$

$$\text{if } \left( x \in L^* \wedge y \in L^* \right) \text{ then } xy \in L^*$$

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

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

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

Example (idempotence):

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

$$\text{idempotence and } L^* = \{aa\}^*$$

# OPERATIONS ON LANGUAGES / 9

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).

$$\Sigma_A = \{A, B, ..., Z\} \quad \Sigma_N = \{0, 1, 2, ..., 9\}$$

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

$$\text{if} \quad \Sigma = \Sigma_A \cup \Sigma_N$$

$$I_5 = \Sigma_A (\Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \Sigma^3 \cup \Sigma^4)$$

$$I_5 = \Sigma_A (\Sigma \cup \varepsilon)^4$$

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

# OPERATIONS ON LANGUAGES / 10

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.

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

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

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

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:

$$\Sigma^4 \Sigma^*$$

$$(\Sigma^+)^4$$

# OPERATIONS ON LANGUAGES / 11

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

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

Example of quotienting:

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

$$L' / L'' = \{ a^r b^s \mid (r \geq 2 \ \textit{even}) \wedge (1 \leq s < r, s \quad \textit{odd}) \ \}$$

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

$$L'' / L' = \varnothing$$

Question: what happens if x ∈ L' does not admit any suffix z ∈ L'' ?

# Bibliography

- S. Crespi Reghizzi, *Linguaggi Formali e Compilazione*, Città Studi, UTET

- Hopcroft, Ullman, *Formal Languages and their Relation to Automata*, Addison Wesley, 1969

- A. Salomaa – *Formal Languages*, Academic Press, 1973

- D. Mandrioli, C. Ghezzi – *Theoretical Foundations of Computer Science*, John Wiley & Sons, 1987

- L. Breveglieri, S. Crespi Reghizzi, *Linguaggi Formali e Compilatori: Temi d'Esame Risolti,* web site