

Free Grammars - III

Prof. Licia Sbattella

aa 2007-08

Translated and adapted by L. Breveglieri

WEAK AND STRONG (or STRUCTURAL) EQUIVALENCE

WEAK EQUIVALENCE: two grammars G and G' are weakly equivalent if and only if they generate the same language, that is if $L(G) = L(G')$.

Such grammars might give completely different structures to the same phrase, one of which structures might be inadequate as for the desired semantic interpretation.

STRONG or STRUCTURAL EQUIVALENCE: two grammars G and G' are strongly equivalent if and only if they generate the same language, that is if $L(G) = L(G')$, and moreover if G and G' assign to every phrase syntactic trees that can be considered structurally equivalent (this means that for every phrase the condensed skeleton trees should be identical).

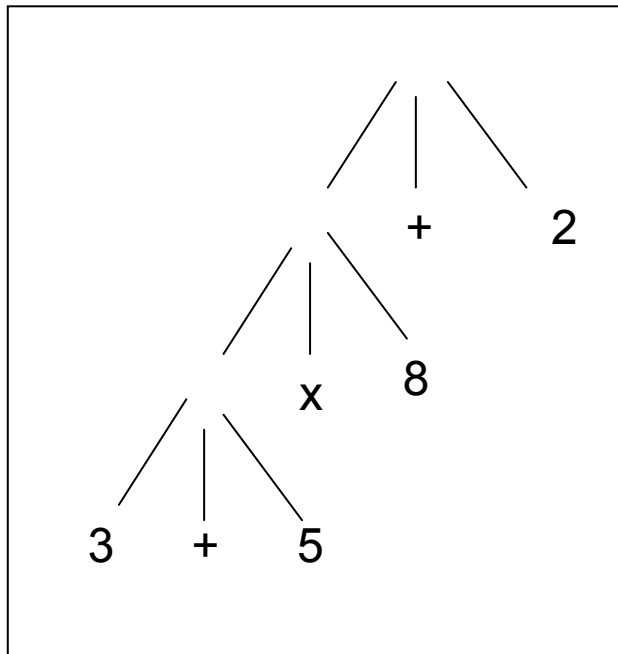
Strong equivalence implies weak equivalence, but the opposite implication does not hold.

Strong equivalence is a *decidable* property.

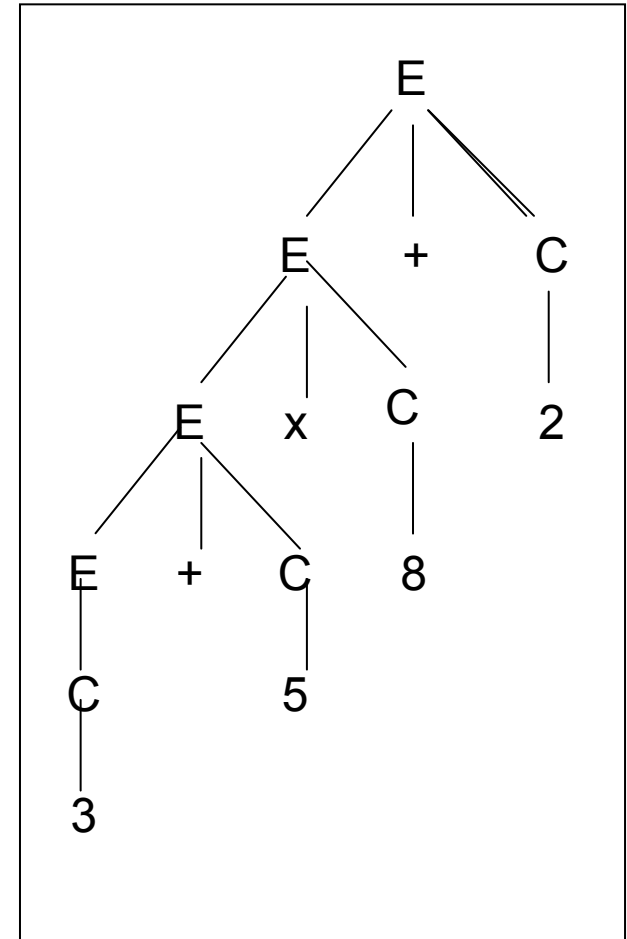
Weak equivalence is a *undecidable* property (but in some special cases).

EXAMPLE: structural equivalence of arithmetic expressions $- 3 + 5 \times 8 + 2$

$G_1: E \rightarrow E + C \quad E \rightarrow E \times C \quad E \rightarrow C$
 $C \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$



skeleton tree



full tree

$$G_2: E \rightarrow E + T \quad E \rightarrow T \quad T \rightarrow T \times C \quad T \rightarrow C$$

$$C \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$$

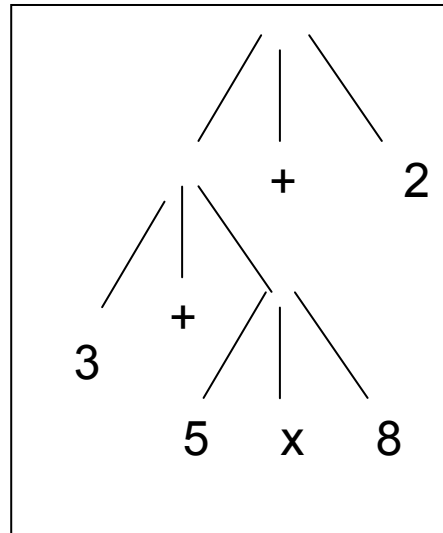
G_1 and G_2 are not equivalent in the structural sense.

semantic interpretations:

$G_1: (((3 + 5) \times 8) + 2)$

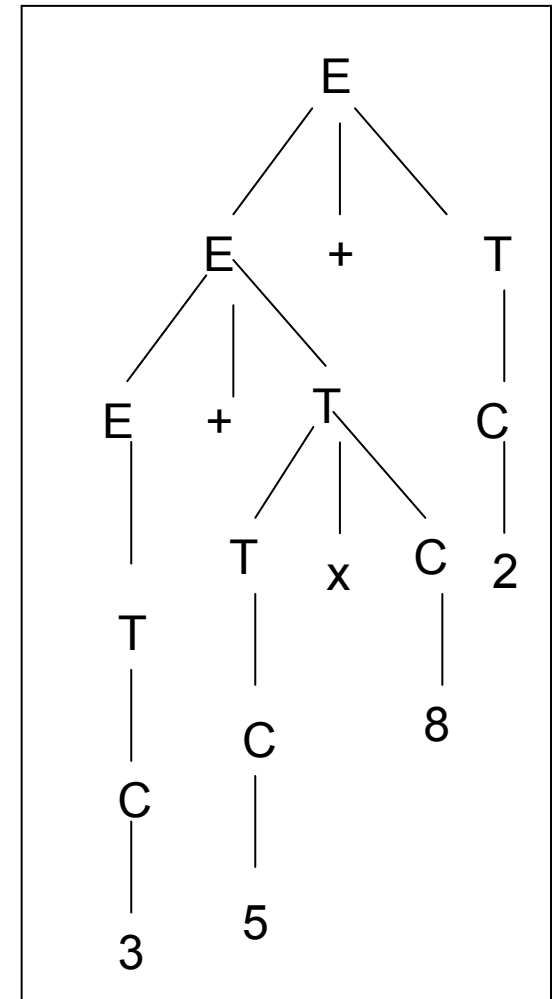
$G_2: ((3 + (5 \times 8)) + 2)$

Only G_2 is structurally adequate with reference to the usual precedence rules between operators (but G_2 is more complex).



skeleton tree

full tree



When defining formally a language, STRUCTURAL ADEQUACY is an issue to be considered carefully, as usually the grammar generating the language is the syntactic basis to define the semantic interpretation or to design a syntax-driven transduction.

G_3 is structurally equivalent to the previous grammar:

$$\begin{aligned} E &\rightarrow E + T \mid T + T \mid C + T \mid E + C \mid T + C \mid C + C \mid T \times C \mid C \times C \mid C \\ T &\rightarrow T \times C \mid C \times C \mid C \\ C &\rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

G_3 contains more production rules than the previous grammar, as it is not designed by sharing or grouping substructures according to taxonomy. Categories and taxonomy may be helpful to reduce grammar complexity.

STRUCTURAL EQUIVALENCE IN THE LARGE SENSE

$$\begin{array}{c} \{S \rightarrow Sa \mid a\} \quad \{X \rightarrow aX \mid a\} \\ L = a^+ \\ \begin{array}{cc} \underbrace{a \ a} & \underbrace{a \ a} \\ \underbrace{\hspace{1.5cm}} & \underbrace{\hspace{1.5cm}} \end{array} \end{array}$$

In the strict sense, these two grammars are not structurally equivalent, since the syntactic trees of the sample string aa below are not identical.

However, one can observe that every left linear tree of the first grammar corresponds to a linear right tree of the second grammar (that is, for the same string the two grammars generate two trees that are mirror images of each other).

The two grammars are strongly equivalent *in the large sense*, that is by abstracting from some minor structural detail.

GRAMMAR NORMAL FORM

Normal forms impose restrictions to the production rules, but without reducing the family of generated languages. Such forms are useful for instance in proving theorems, rather than in designing a grammar for a specific language.

There are some transformations to put a general grammar into an equivalent normal form. Such transformations are helpful to design syntactic analyzers.

GRAMMAR to start from:

$$G = \{\Sigma, V, P, S\}$$

NON-TERMINAL EXPANSION

Expanding a non-terminal does not change the generative power of the grammar.

$$\begin{aligned} A &\rightarrow \alpha B \gamma & B &\rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n \\ A &\rightarrow \alpha \beta_1 \gamma \mid \alpha \beta_2 \gamma \mid \dots \mid \alpha \beta_n \gamma \\ A &\Rightarrow \alpha B \gamma \Rightarrow \alpha \beta_i \gamma \\ A &\Rightarrow \alpha \beta_i \gamma \end{aligned}$$

REMOVING THE AXIOM FROM THE RIGHT MEMBER OF THE PRODUCTION RULES

It is always possible to restrict the right member of a production rules to being a string over the alphabet $(\Sigma \cup (V - \{S\}))$, that is to not containing the axiom S . It suffices to introduce the new axiom S_0 and the new production rule $S_0 \rightarrow S$.

NULLABLE NON-TERMINAL SYMBOL and NON-NULLABLE NORMAL FORM

A non-terminal A is said to be *nullable* if and only if there exists a derivation as aside:

$$A \xRightarrow{+} \varepsilon$$

CAUTION: this does not exclude that A can generate also a string different from ε .

$Null \subseteq V$ is the set of the nullable non-terminal symbols

Computation of the set of the nullable non-terminal symbols:

$$A \in Null \text{ if } A \rightarrow \varepsilon \in P$$

$$A \in Null \text{ if } (A \rightarrow A_1 A_2 \dots A_n \in P \text{ with } A_i \in V \setminus \{A\}) \wedge \forall A_i (A_i \in Null)$$

EXAMPLE – how to identify the nullable non-terminal symbols

$$S \rightarrow SAB \mid AC \quad A \rightarrow aA \mid \varepsilon \quad B \rightarrow bB \mid \varepsilon \quad C \rightarrow cC \mid c$$
$$Null = \{A, B\}$$

NON-NULLABLE NORMAL FORM (and not containing empty rules): every non-terminal symbol different from the axiom is not nullable, and the axiom itself is nullable if and only if the language contains the empty string ε .

HOW TO OBTAIN THE NON-NULLABLE NORMAL FORM OF A GRAMMAR:

- 1) Compute the *Null* set.
- 2) For every production rule P , add the alternative rules that can be obtained from the rule P by deleting each nullable non-terminal occurring in the right member of the rule, in all possible combinations.
- 3) Remove all the empty production rules of the type $A \rightarrow \varepsilon$, but when $A = S$.
- 4) Remove $S \rightarrow \varepsilon$ if and only if the language does not contain ε .
- 5) Put the grammar in reduced form and remove circular derivations, if any.

EXAMPLE (continued)

nullable	G original	G'' to be reduced	G'' without empty rules
F	$S \rightarrow SAB \mid$ $\quad \mid AC$	$S \rightarrow SAB \mid SA \mid SB \mid$ $\quad \mid S \mid AC \mid C$	$S \rightarrow SAB \mid SA \mid SB \mid$ $\quad \mid AC \mid C$
V	$A \rightarrow aA \mid \varepsilon$	$A \rightarrow aA \mid a$	$A \rightarrow aA \mid a$
V	$B \rightarrow bB \mid \varepsilon$	$B \rightarrow bB \mid b$	$B \rightarrow bB \mid b$
F	$C \rightarrow cC \mid c$	$C \rightarrow cC \mid c$	$C \rightarrow cC \mid c$

COPY RULES (categorization) AND HOW TO ELIMINATE THEM

Copy rule (or categorization rule): in the example, the syntactic class B is included in the syntactic class A. By removing the copy rules, one obtains an equivalent grammar that has less deep syntactic trees.

Copy (A) $\subseteq V$: the set of the non-terminal symbols A may be copied into, possibly in a transitive way.

$$A \rightarrow B \quad B \in V$$

$$\textit{frase_iterativa} \rightarrow \textit{frase_while} \mid \textit{frase_for} \mid \textit{frase_repeat}$$

$$\text{Copy}(A) = \left\{ B \in V \mid \text{there exists the derivation } A \overset{*}{\Rightarrow} B \right\}$$

Copy rules are not totally unuseful, as sometimes they allow to share some parts of the grammar. For this reason copy rules may be used in the grammars of technical languages, like for instance programming languages.

1) Compute the set *Copy* (assume the grammar does not contain empty rules). The computation is expressed by the following logicals clauses, to be iterated as long as a fixpoint has been reached.

$$A \in \text{Copy} (A)$$

$$C \in \text{Copy} (A) \text{ if } (B \in \text{Copy} (A)) \wedge (B \rightarrow C \in P)$$

2) Construct the production rules of the grammar G' , equivalent to G but without copy rules.

$$P' := P \setminus \{A \rightarrow B \mid A, B \in V\}$$

$$P' := \{A \rightarrow \alpha \mid \alpha \in ((\Sigma \cup V)^* \setminus V)\}$$

$$\text{where } B \rightarrow \alpha \in P \wedge B \in \text{Copy} (A)$$

$$A \overset{*}{\Rightarrow} B \Rightarrow \alpha \text{ becomes } A \Rightarrow \alpha$$

EXAMPLE – grammar of the arithmetic expressions without copy rules

standard non-ambiguous
grammar
but with copy rules

$$\begin{aligned} E &\rightarrow E + T \mid T & T &\rightarrow T \times C \mid C \\ C &\rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \\ \text{Copy}(E) &= \{E, T, C\} & \text{Copy}(T) &= \{T, C\} \\ \text{Copy}(C) &= \{C\} \end{aligned}$$

equivalent grammar without copy rules

$$\begin{aligned} E &\rightarrow E + T \mid T \times C \mid 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \\ T &\rightarrow T \times C \mid 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \\ C &\rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

CHOMSKY NORMAL FORM: the production rules are only of two types

1. binary rule (of the homogeneous type)

$$A \rightarrow BC \quad \text{where } B, C \in V$$

2. terminal rule (with one-letter right side)

$$A \rightarrow \alpha \quad \text{where } \alpha \in \Sigma$$

If the language contains the empty string, add the rule:

$$S \rightarrow \varepsilon$$

Structure of the syntax tree: internal nodes all of arity 2, and nodes that are the parent of a leaf all of arity 1.

transformation

$$A_0 \rightarrow A_1 A_2 \dots A_n$$

$$A_0 \rightarrow A_1 \langle A_2 \dots A_n \rangle$$

$$\langle A_2 \dots A_n \rangle \rightarrow A_2 \dots A_n$$

$$\text{if } A_1 \text{ is a terminal} \quad \langle A_1 \rangle \rightarrow A_1$$

EXAMPLE – transformation into Chomsky normal form

$$S \rightarrow dA \mid cB \quad A \rightarrow dAA \mid cS \mid c \quad B \rightarrow cBB \mid dS \mid d$$

$$S \rightarrow \langle d \rangle A \mid \langle c \rangle B$$

$$A \rightarrow \langle d \rangle \langle AA \rangle \mid \langle c \rangle S \mid c$$

$$B \rightarrow \langle c \rangle \langle BB \rangle \mid \langle d \rangle S \mid d$$

$$\langle d \rangle \rightarrow d \quad \langle c \rangle \rightarrow c$$

$$\langle AA \rangle \rightarrow AA \quad \langle BB \rangle \rightarrow BB$$

HOW TO TRANSFORM LEFT RECURSION INTO RIGHT RECURSION

Construction of the NON-LEFT RECURSIVE FORM (indispensable for designing left recursive descent syntactic analyzers).

1) TRANSFORMATION OF IMMEDIATE LEFT RECURSION

$$A \rightarrow A\beta_1 \mid A\beta_2 \mid \dots \mid A\beta_h$$

where no β_i is empty

$$A \rightarrow \gamma_1 \mid \gamma_2 \mid \dots \mid \gamma_k$$

$$A \rightarrow \gamma_1 A' \mid \gamma_2 A' \mid \dots \mid \gamma_k A' \mid \gamma_1 \mid \gamma_2 \mid \dots \mid \gamma_k$$

$$A' \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_h A' \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_h$$

$$A \Rightarrow A\beta_2 \Rightarrow A\beta_3\beta_2 \Rightarrow \gamma_1\beta_3\beta_2$$

$$A \Rightarrow \gamma_1 A' \Rightarrow \gamma_1\beta_3 A' \Rightarrow \gamma_1\beta_3\beta_2$$

EXAMPLE – how to shift to the right a left immediate recursion

$$E \rightarrow E + T \mid T \quad T \rightarrow T * F \mid F \quad F \rightarrow (E) \mid i$$

$E \quad T$ are immediately recursive on the left

$$E \rightarrow TE' \mid T \quad E' \rightarrow +TE' \mid +T$$

$$T \rightarrow FT' \mid F \quad T' \rightarrow *FT' \mid *F \quad F \rightarrow (E) \mid i$$

; simple transformation but does not work always

$$E \rightarrow T + E \mid T \quad T \rightarrow F * T \mid F \quad F \rightarrow (E) \mid i$$

2) TRANSFORMATION OF NON-IMMEDIATE RECURSION

HYPOTHESIS: G is non-homogeneous non-nullable form, with rules of length one (similar to the Chomsky normal form but without having two non-terminals).

The ALGORITHM works iteratively, in two interleaved steps:

- 1) expansion to expose left immediate recursions
- 2) replacement of left recursion with right recursion

It is advisable to denumerate the non-terminal symbols from 1 to m :

$$V = \{A_1, A_2, \dots, A_m\}$$

A_1 is the axiom

IDEA of the algorithm: modify the rules so that, if a rule is $A_i \rightarrow A_j$, it holds $j > i$.

ALGORITHM TO ELIMINATE LEFT RECURSION (possibly non-immediate)

```
for  $i := 1$  to  $m$  do  
  for  $j := 1$  to  $i - 1$  do  
    replace every rule of the form  $A_i \rightarrow A_j \alpha$   
    with the following rules:  
     $A_i \rightarrow Y_1 \alpha \mid Y_2 \alpha \mid \dots \mid Y_k \alpha$   
    (this may introduce new left immediate recursion)  
    where  $A_j \rightarrow Y_1 \mid Y_2 \mid \dots \mid Y_k$  are the alternatives for  $A_j$   
  end do  
  eliminate, by means of the previous algorithm, the possible  
  left immediate recursion that occur in the alternatives of  $A_i$ ,  
  and introduce the new non-terminal  $A'_i$   
end do
```

This same algorithm can be adapted to transform right recursion into left recursion. Such transformation is sometimes required to design bottom-up syntactic analyzers.

EXAMPLE – apply the algorithm to grammar G_3

$$\begin{array}{l} A_1 \rightarrow A_2 a \mid b \quad A_2 \rightarrow A_2 c \mid A_1 d \mid e \\ A_1 \Rightarrow A_2 a \Rightarrow A_1 da \end{array}$$

$i = 1$ Eliminate the left immediate recursion of A_1 grammar is unchanged
(here it does not exist).

$i = 2 \ j = 1$ Replace $A_2 \rightarrow A_1 d$ by means
of the rules obtained
by expanding A_1

$$\begin{array}{l} A_1 \rightarrow A_2 a \mid b \\ A_2 \rightarrow A_2 c \mid A_2 a d \mid b d \mid e \end{array}$$

Eliminate the left immediate recursion
and eventually obtain G'_3

$$\begin{array}{l} A_1 \rightarrow A_2 a \mid b \\ A_2 \rightarrow b d A' \mid e A' \\ A' \rightarrow c A' \mid a d A' \mid e \end{array}$$

REAL TIME NORMAL FORM (and GREIBACH NORMAL FORM)

In the REAL TIME NORMAL FORM every production rule begins with a terminal symbol.

$$A \rightarrow a\alpha \quad \text{where} \quad a \in \Sigma, \alpha \in \{\Sigma \cup V\}^*$$

The GREIBACH NORMAL FORM is a special case of the above.

Every rule begins with a terminal symbol, followed by zero or more non-terminal symbols.

$$A \rightarrow a\alpha \quad \text{where} \quad a \in \Sigma, \alpha \in V^*$$

“REAL TIME” - this term is justified as a property of the syntax analysis algorithm, which at every step reads and consumes exactly one terminal symbol. Therefore, the number of steps needed to complete the analysis is equal to the length of the string to be analyzed.

ALGORITHM TO TRANSFORM THE GRAMMAR INTO REAL TIME NORMAL FORM AND GREIBACH NORMAL FORM

HYPOTHESIS: the grammar does not contain any nullable non-terminal symbol

- 1) eliminate left recursion (immediate or not)
- 2) by means of elementary transformations, expand the non-terminal symbols possibly on the left end of the right member of the production rules
- 3) introduce new non-terminal symbols to replace the terminal symbols possibly occurring inside of the right member of the production rules (or at the right end)

If the last step of the above algorithm were skipped, the production rules might still contain terminal symbols in the middle or at the right end.

The grammar would then be in real time normal form, though in general not in Greibach normal form.

EXAMPLE – grammar to start from

$$A_1 \rightarrow A_2 a \quad A_2 \rightarrow A_1 c \mid bA_1 \mid d$$

1) eliminate left recursion

$$\begin{aligned} A_1 &\rightarrow A_2 a & A_2 &\rightarrow A_2 a c \mid bA_1 \mid d \\ A_1 &\rightarrow A_2 a & A_2 &\rightarrow bA_1 A'_2 \mid dA'_2 \mid d \mid bA_1 & A'_2 &\rightarrow acA'_2 \mid ac \end{aligned}$$

2) replace iteratively each leftmost non-terminal symbol, as long as a terminal symbol takes its place

$$\begin{aligned} A_1 &\rightarrow bA_1 A'_2 a \mid dA'_2 a \mid da \mid bA_1 a \\ A_2 &\rightarrow bA_1 A'_2 \mid dA'_2 \mid d \mid bA_1 \\ A'_2 &\rightarrow acA'_2 \mid ac \end{aligned}$$

3) replace the terminal symbols occurring in the middle of the production (or at the right end) by means of non-terminal symbols (and of the related productions).

$$\begin{aligned} A_1 &\rightarrow bA_1 A'_2 <a> \mid dA'_2 <a> \mid d <a> \mid bA_1 <a> \\ A_2 &\rightarrow bA_1 A'_2 \mid dA'_2 \mid d \mid bA_1 \\ A'_2 &\rightarrow a <c> A'_2 \mid a <c> \\ <a> &\rightarrow a & <c> &\rightarrow c \end{aligned}$$

Bibliography

- S. Crespi Reghizzi, *Linguaggi Formali e Compilazione*, Pitagora Editrice Bologna, 2006
- 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 (eng + ita)