

Free Grammars - III

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, i.e., iff $L(G) = L(G')$

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, i.e., iff $L(G) = L(G')$, 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)

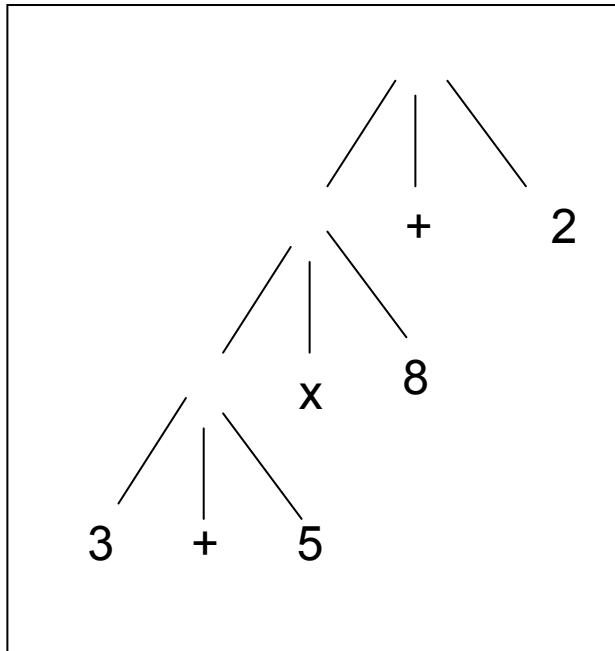
Strong equivalence implies weak equivalence, but the opposite implication does not hold

Strong equivalence is a *decidable* property

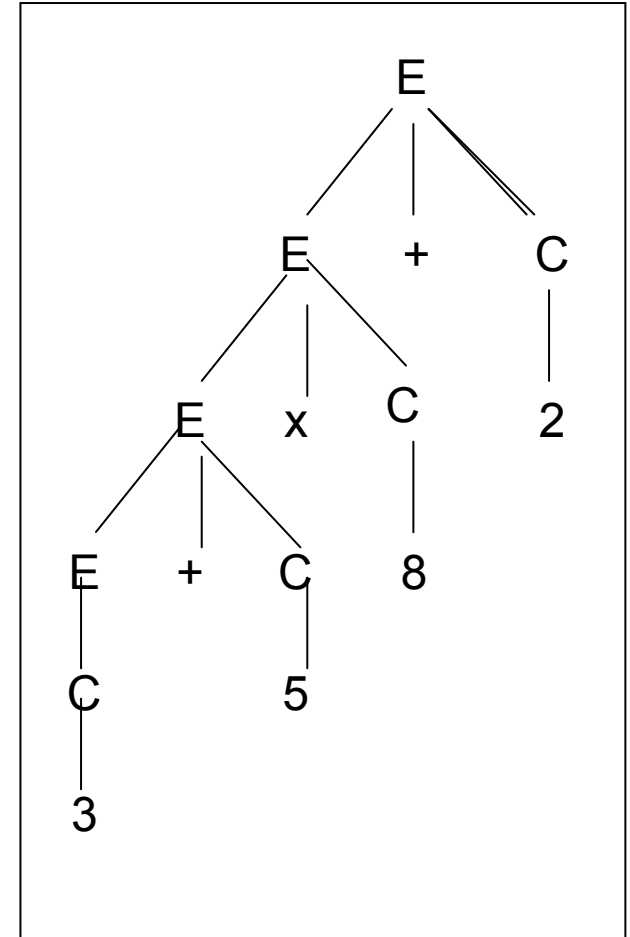
Weak equivalence is an *undecidable* property (but in some special cases)

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

$$G_2: E \rightarrow E + T \quad E \rightarrow T \quad T \rightarrow T \times C \quad T \rightarrow C \\ C \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$$

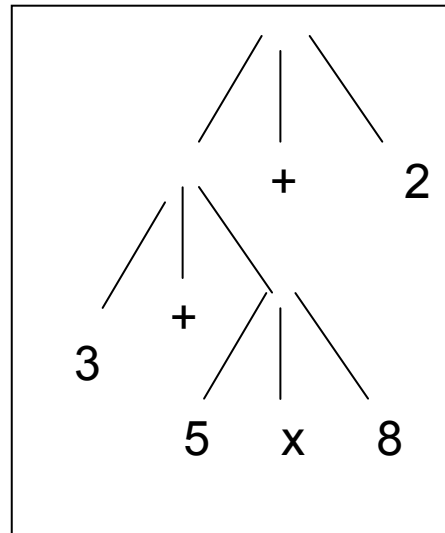
G_1 and G_2 are not equivalent in the structural sense

semantic interpretations

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

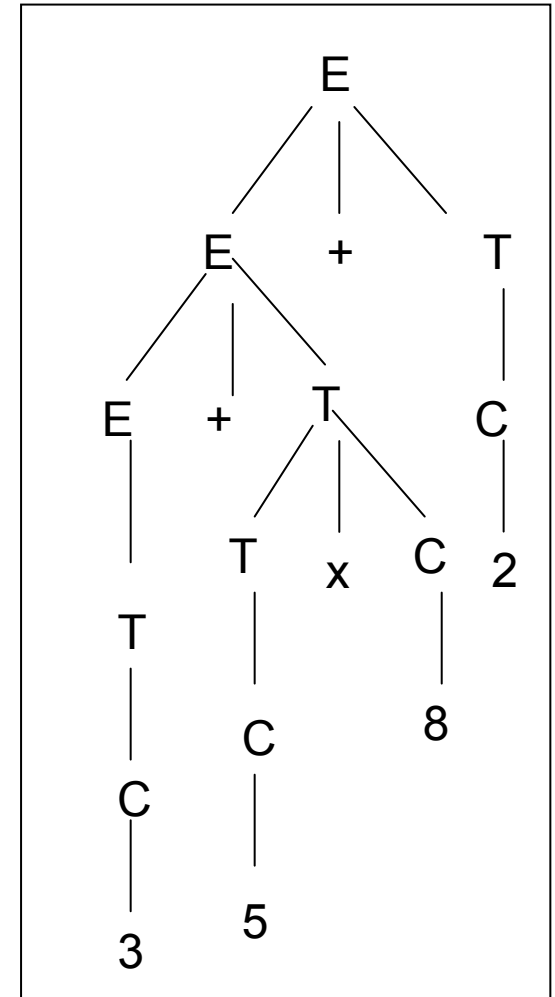
$G_2: ((3 + (5 \times 8)) + 2)$

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

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

G_3 contains more 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

$$\{S \rightarrow Sa \mid a\} \quad \{X \rightarrow aX \mid a\}$$
$$L = a^+$$
$$\underbrace{a \ a} \qquad \qquad a \ \underbrace{a}$$
$$\underbrace{\qquad} \qquad \qquad \underbrace{\qquad}$$

In the strict sense, these two grammars are not structurally equivalent, as 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 right linear 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*, i.e., by abstracting from some minor structural detail

GRAMMAR NORMAL FORM

Normal forms impose restrictions to the rules, but without reducing the family of generated languages. Such forms are useful for instance to prove theorems, rather than to design 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

$$A \rightarrow \alpha B \gamma \quad B \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

$$A \rightarrow \alpha \beta_1 \gamma \mid \alpha \beta_2 \gamma \mid \dots \mid \alpha \beta_n \gamma$$

$$A \Rightarrow \alpha B \gamma \Rightarrow \alpha \beta_i \gamma$$

$$A \Rightarrow \alpha \beta_i \gamma$$

REMOVING THE AXIOM FROM THE RIGHT MEMBER OF THE RULES

It is always possible to restrict the right member of a rule so as it is a string over the alphabet $(\Sigma \cup (V - \{S\}))$, i.e., it does not contain the axiom S

It suffices to introduce the new axiom S_0 and the new 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 generates 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

$A \in Null$ if $A \rightarrow \varepsilon \in P$

$A \in Null$ 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 rule P, add the alternative rules that can be obtained from rule P by deleting each nullable non-terminal that occurs in the right member of the rule, in all possible combinations
- 3) remove all the empty rules of the type $A \rightarrow \varepsilon$, except when $A = S$
- 4) remove $S \rightarrow \varepsilon$ if and only if the language does not contain ε
- 5) put the grammar into reduced form and remove the circular derivations

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 with syntactic trees that are less deep

$\text{Copy}(A) \subseteq V$: the set of the non-terminal symbols into which A can be copied, possibly in a transitive way

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

$$\textit{iterative_phrase} \rightarrow \textit{while_phrase} \mid \textit{for_phrase} \mid \textit{repeat_phrase}$$

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

Copy rules are not totally useless, as sometimes they allow us to share some parts of the grammar. For this reason copy rules may be used in the grammars of technical languages, like for instance in the 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 until a fixpoint is 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 rules of the grammar G' , equivalent to G but without copy rules

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

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

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: every internal node has arity 2,
and every node that is parent of a leaf has arity 1

transformation

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

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

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

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

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 a left immediate recursion to the right

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

$E \quad T$ 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 always work

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

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 recursions that occur in the alternatives of  $A_i$ ,  
  and introduce the new non-terminal  $A'_i$   
end do
```

The same algorithm can be adapted to transform right recursion into left recursion
This 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 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 reads and consumes exactly one terminal symbol at every step. Therefore, the number of steps needed to complete the analysis is equal to the length of the string to analyze

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 rules
- 3) introduce new non-terminal symbols to replace the terminal symbols possibly occurring inside of the right member of the rules (or at the right end)

If the last step of the above algorithm were skipped, the 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

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

$$A_1 \rightarrow A_2 a \quad A_2 \rightarrow bA_1 A'_2 \mid dA'_2 \mid d \mid bA_1 \quad A'_2 \rightarrow acA'_2 \mid ac$$

2) iteratively replace every leftmost non-terminal symbol, until a terminal symbol takes its place

$$A_1 \rightarrow bA_1 A'_2 a \mid dA'_2 a \mid da \mid bA_1 a$$

$$A_2 \rightarrow bA_1 A'_2 \mid dA'_2 \mid d \mid bA_1$$

$$A'_2 \rightarrow acA'_2 \mid ac$$

3) replace the terminal symbols that occur in the middle of the rule (or at the right end) with non-terminal symbols (and with the related rules)

$$A_1 \rightarrow bA_1 A'_2 <a> \mid dA'_2 <a> \mid d <a> \mid bA_1 <a>$$

$$A_2 \rightarrow bA_1 A'_2 \mid dA'_2 \mid d \mid bA_1$$

$$A'_2 \rightarrow a <c> A'_2 \mid a <c>$$

$$<a> \rightarrow a \quad <c> \rightarrow c$$