

Free Grammars - B

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)

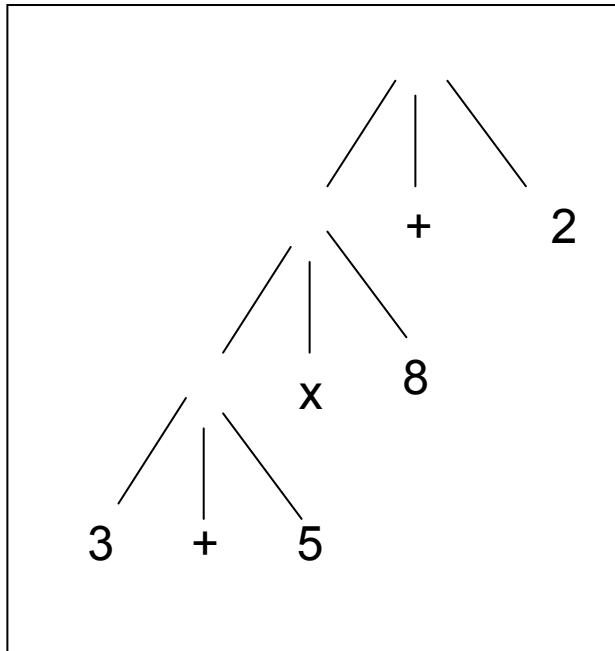
Strong equivalence implies weak equivalence, but the opposite implication does not hold

Strong equivalence is a *decidable* property

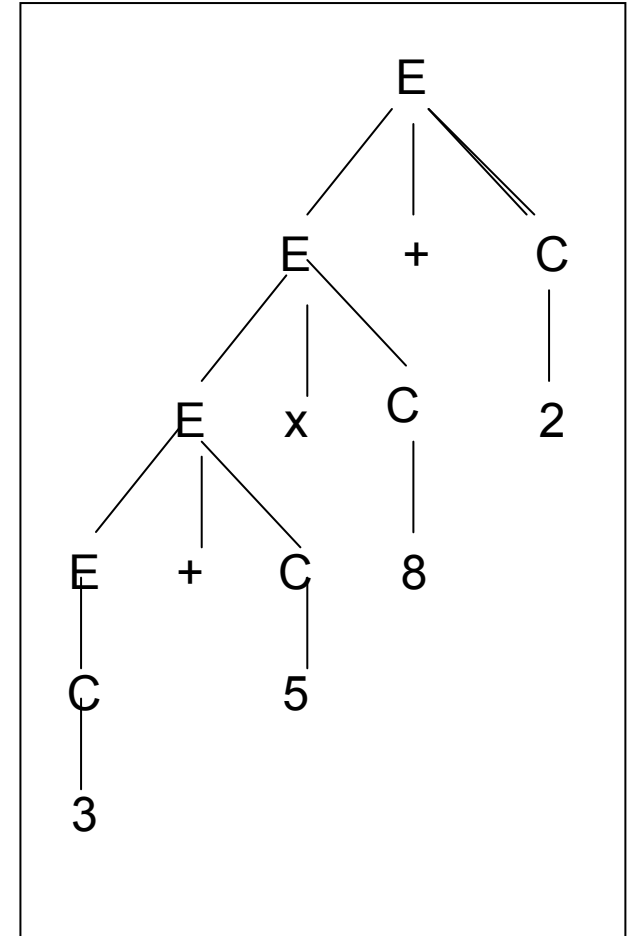
Weak equivalence is an *undecidable* property (but in some special cases)

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

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

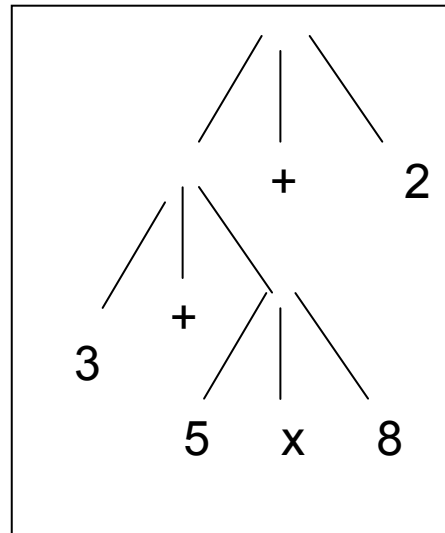
G_1 and G_2 are not equivalent in the structural sense

semantic interpretations

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

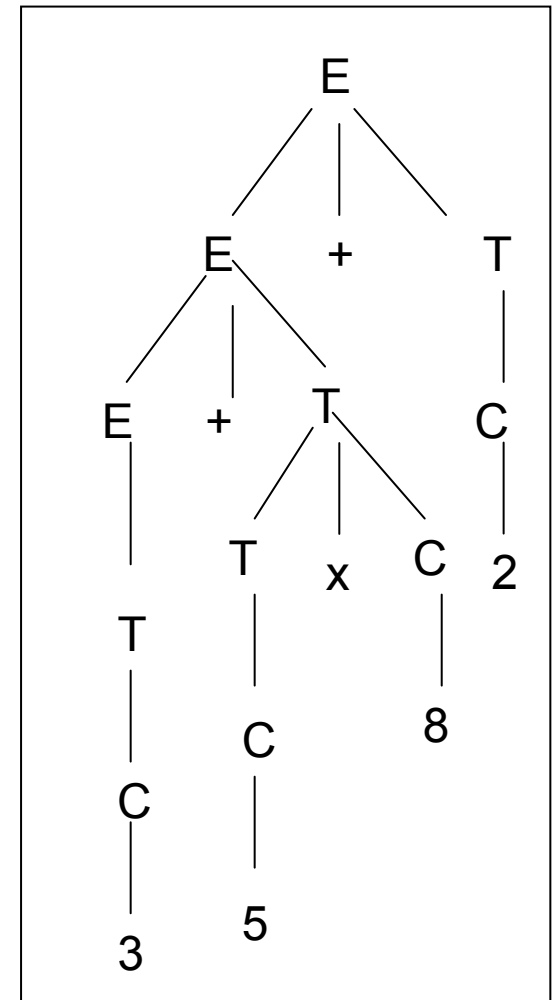
$G_2: ((3 + (5 \times 8)) + 2)$

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



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

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

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

$$\begin{aligned} P' &:= P \setminus \{A \rightarrow B \mid A, B \in V\} \\ P' &:= \{A \rightarrow \alpha \mid \alpha \in ((\Sigma \cup V)^* \setminus V)\} \\ \text{where } B &\rightarrow \alpha \in P \wedge B \in \text{Copy} (A) \\ A &\stackrel{*}{\Rightarrow} B \Rightarrow \alpha \text{ becomes } A \Rightarrow \alpha \end{aligned}$$

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 $j > i$

GRAMMARS FOR REGULAR LANGUAGES

Regular languages are a special case of the free ones

Regular languages can be generated by grammars that exhibit strongly restricted rules

The phrases of a regular language necessarily contain repeated factors as the length of the string grows

It is possible to prove in a rigorous way that some (free) languages cannot be generated by any regular expression

REGULAR LANGUAGES

Recognizing phrases requires only a finite amount of memory

FREE LANGUAGES

Recognizing phrases requires an unbounded amount of memory

FROM THE REGULAR EXPRESSION TO THE FREE GRAMMAR

The iterative regular operators (star and cross) must be replaced with recursive rules

Decompose the regexp into subexpressions and denumerate them progressively. After the very definition of regexp, the following cases are obtained, which allow Us to design the corresponding rules shown aside

1. $r = r_1.r_2....r_k$
2. $r = r_1 \cup r_2 \cup \cup r_k$
3. $r = (r_1)^*$
4. $r = (r_1)^+$
5. $r = b \in \Sigma$
6. $r = \varepsilon$

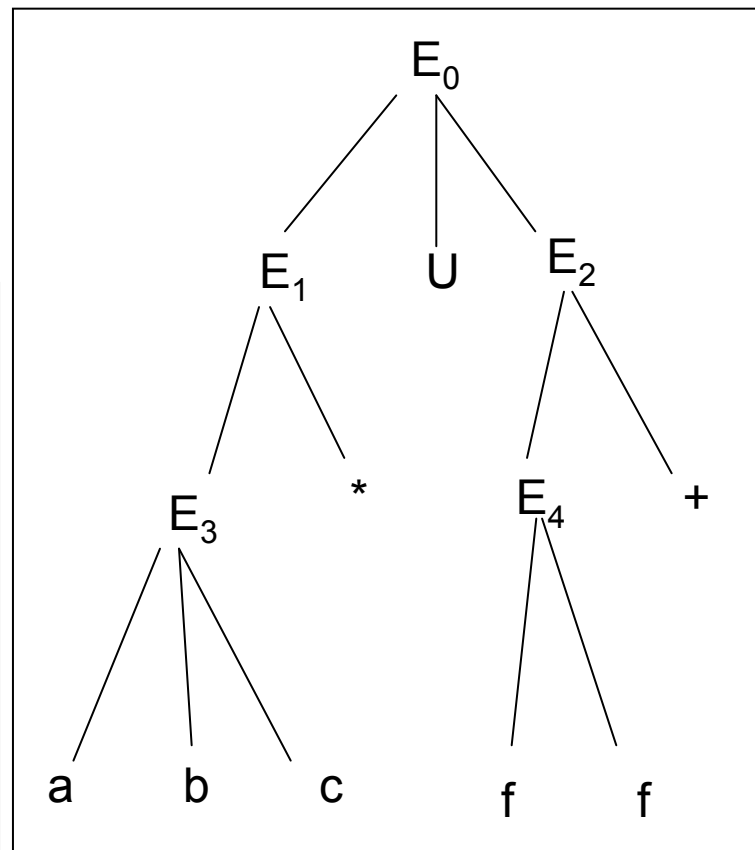
1. $E = E_1E_2...E_k$
2. $E = E_1 \cup E_2 \cup ... \cup E_k$
3. $E = EE_1 \mid \varepsilon$ or $E = E_1E \mid \varepsilon$
4. $E = EE_1 \mid E_1$ or $E = E_1E \mid E_1$
5. $E = b$
6. $E = \varepsilon$

Uppercase letters indicate the non-terminal symbols. If E_i is a terminal symbol $a_i \in \Sigma$, the symbol a_i is written instead of E_i

EXAMPLE

$$E = (abc)^* \cup (ff)^+$$

- | | | |
|----|----------------|--|
| 2. | $E_1 \cup E_2$ | $E_0 \rightarrow E_1 \mid E_2$ |
| 3. | E_3^* | $E_1 \rightarrow E_1 E_3 \mid \varepsilon$ |
| 4. | E_4^+ | $E_2 \rightarrow E_2 E_4 \mid E_4$ |
| 1. | abc | $E_3 \rightarrow abc$ |
| 1. | ff | $E_4 \rightarrow ff$ |



If the regexp to start from is ambiguous, the derived equivalent free grammar is ambiguous as well

Every regular language is free, but there are many free languages that are not regular, like for instance the palindromes and the arithmetic expressions with parentheses

$$REG \subset LIB$$

UNI-LINEAR GRAMMAR (or GRAMMAR of TYPE 3)

RIGHT UNI-LINEAR RULE (or simply RIGHT LINEAR)

$$A \rightarrow uB \quad \text{where} \quad u \in \Sigma^*, B \in (V \cup \varepsilon)$$

LEFT UNI-LINEAR RULE (or simply LEFT LINEAR)

$$A \rightarrow Bv \quad \text{where} \quad v \in \Sigma^*, B \in (V \cup \varepsilon)$$

The corresponding syntax trees grow in a heavily unbalanced way

If a left or right uni-linear grammar contains a recursive derivation, then such a derivation is recursive on the left or right, respectively

EXAMPLE: the phrases that contain the substring *aa* and terminate with *b* are generated by the following regexp (ambiguous)

$$(a \mid b)^* aa(a \mid b)^* b$$

1. right linear grammar G_d

$$S \rightarrow aS \mid bS \mid aaA \quad A \rightarrow aA \mid bA \mid b$$

2. left linear grammar G_s

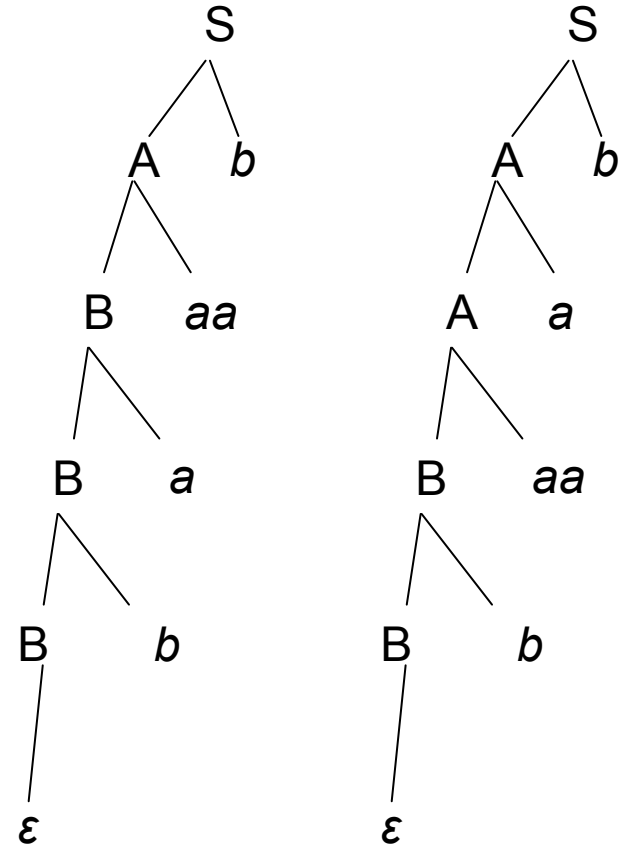
$$S \rightarrow Ab \quad A \rightarrow Aa \mid Ab \mid Baa$$

$$B \rightarrow Ba \mid Bb \mid \varepsilon$$

3. equivalent non-unilinear grammar

$$E_1 \rightarrow E_2 aa E_2 b \quad E_2 \rightarrow E_2 a \mid E_2 b \mid \varepsilon$$

G_s – syntax trees of the ambiguous phrase *baaab*



EXAMPLE : arithmetic expressions without parentheses

$$L = \{a, a + a, a * a, a + a * a, \dots\}$$

$$G_d : S \rightarrow a \mid a + S \mid a * S$$

$$G_s : S \rightarrow a \mid S + a \mid S * a$$

G_s and G_d are not adequate to compilation, as their syntax structure does not reflect the usual precedence rules between arithmetic operators. It is convenient to transform a unilinear grammar into a strictly unilinear grammar, where every rule contains at most a single terminal character

$$A \rightarrow aB \quad \text{or} \quad A \rightarrow Ba \quad \text{where} \quad a \in (\Sigma \cup \varepsilon), B \in (V \cup \varepsilon)$$

Without any loss of generality, one can impose to have only null terminal rules

$$A \rightarrow aB \mid \varepsilon \quad \text{where} \quad a \in \Sigma$$

FROM UNILINEAR GRAMMAR TO REGULAR EXPRESSION: LINEAR EQUATIONS

UNILINEAR GRAMMAR / REGULAR LANGUAGE

First it is proved that the family of the languages generated by unilinear grammars strictly contains that of regular languages. Then it is proved that these two families actually coincide

One can think of the rules of a right (uni)linear grammar as of equations, the unknown terms of which are the languages generated by each non-terminal symbol.

Suppose that G is a strictly right linear grammar and that G contains only null terminal rules $A \rightarrow \varepsilon$

$$G = (V, \Sigma, P, S)$$
$$L_A = \left\{ x \in \Sigma^* \mid A \xRightarrow{+} x \right\} \quad L(G) \equiv L_S$$

replace the (recursive) rule $A \rightarrow a A$ with the equation $L_A = L_a L_A$

replace the rule $A \rightarrow b B$ with the equation $L_A = L_b L_B$

replace the rule $A \rightarrow a$ with the equation $L_A = L_a$

unite terminal rules to non-terminal rules:

$A \rightarrow b B$ and $A \rightarrow c$ so that $L_A = L_c \mid L_b L_B$

and so on when there are more terms

then solve the equations by substitution

in the case of a recursive equation $L_A = L_x \mid L_y L_A$,
solve it by means of the Arden identity: $L_A = L_y^* L_x$

EXAMPLE: set-theoretic equations

Here follows a grammar that generates a list of elements e (possibly empty) separated by the character s

$$S \rightarrow sS \mid eA \quad A \rightarrow sS \mid \varepsilon$$

The grammar can be transformed into a set of equations, as aside

EVERY UNILINEAR LANGUAGE
IS REGULAR

$$\begin{cases} L_s = sL_s \cup eL_A \\ L_A = sL_s \cup \varepsilon \end{cases}$$

$$\begin{cases} L_s = sL_s \cup e(sL_s \cup \varepsilon) \\ L_A = sL_s \cup \varepsilon \end{cases}$$

$$\begin{cases} L_s = (s \cup es)L_s \cup e \\ L_A = sL_s \cup \varepsilon \end{cases}$$

and by using Arden identity one obtains

$$\begin{cases} L_s = (s \cup es)^* e \\ L_A = sL_s \cup \varepsilon \quad L_A = s(s \cup es)^* e \cup \varepsilon \end{cases}$$

THE ROLE OF SELF-INCLUSIVE DERIVATIONS

$$A \overset{+}{\Rightarrow} uAv \quad u \neq \varepsilon \wedge v \neq \varepsilon$$

A self-inclusive derivation may not occur in the unilinear grammar of regular languages. This does not happen in the free grammars that generate free non-regular languages (palindromes, Dyck, etc)

The absence of self-inclusive derivations allows to solve the set-theoretic equations associated with unilinear grammars, and gives rise to the simple structure of the languages that such grammars generate

A GRAMMR THAT DOES NOT ADMIT ANY SELF-INCLUSIVE DERIVATION ALWAYS GENERATES A REGULAR LANGUAGE

THE LIMITATIONS OF FREE LANGUAGES

EXAMPLE: THE LANGUAGE WITH THREE EXPONENTS IS NOT FREE
(the proof exploits the possibility of pumping a long phrase by repeating two independent factors, which is incompatible with the definition of such a language)

$$L = \{a^n b^n c^n \mid n \geq 1\}$$

EXAMPLE: THE COPY LANGUAGE

This sub-language is frequent in the programming languages, whenever two lists need to have identical matching elements, like for instance in the lists of the formal and actual parameters of a procedure or function declaration and call

$$L_{replic} = \{uu \mid u \in \Sigma^+\}$$
$$\Sigma = \{a, b\} \quad x = abbbabbb$$

THE CLOSURE PROPERTIES OF THE REG AND LIB FAMILIES

Here the closure properties of the REG and LIB families of languages are examined, with respect to the most common language operators

Let L and R be two languages belonging to LIB and REG, respectively

$R^R \in REG$	$L^R \in LIB$	denotes both union and concatenation
$R^* \in REG$	$L^* \in LIB$	
$\neg R \in REG$	$\neg L \notin LIB$	denotes both union and concatenation
$R_1 \oplus R_2 \in REG$	$L_1 \oplus L_2 \in LIB$	
$R_1 \cap R_2 \in REG$	$L_1 \cap L_2 \notin LIB$	
	$L \cap R \in LIB$	

4) The intersection of two free languages, in general, is not a free language

$$\boxed{\{a^n b^n c^n \mid n \geq 1\} = \{a^n b^n c^+ \mid n \geq 1\} \cap \{a^+ b^n c^n \mid n \geq 1\}}$$

5) There exist free languages the complements of which (w.r.t. the universal language over the same alphabet) are not free languages

CAUTION: of course, there exist also free languages the complements of which are still free languages (and similarly for intersection)

INTERSECTION OF FREE LANGUAGES WITH REGULAR LANGUAGES

One way to make a free grammar more selective, is to restrict the generated free language through a regular filter

EXAMPLE: regular filter over the Dyck language - $\Sigma = \{ a, c \}$ (a open c close)

$$\begin{aligned} L_1 &= L_D \cap \neg(\Sigma^* cc \Sigma^*) = (ac)^* \\ L_2 &= L_D \cap \neg(\Sigma^* ca \Sigma^*) = \{ a^n c^n \mid n \geq 0 \} \end{aligned}$$

The former intersection filters the Dyck language by removing the phrases that contain the factor cc , that is the phrases containing nested structures

The latter intersection filters the Dyck language by removing the phrases that contain two or more concatenated nested structures, leaving only the phrases with one nested structure

Both filtered languages are free, and the former one is regular as well

FREE GRAMMARS EXTENDED BY MEANS OF REGULAR EXPRESSIONS

REGULAR EXPRESSIONS ARE MORE READABLE AND PERSPICUOUS than free grammars, thanks to the iteration operators (star and cross) and to the choice operator (union)

To make free grammars more readable as well, it is often permitted to use regexps also in the right sides of the rules of the grammar. One thus obtains the so-called EXTENDED FREE GRAMMARS (Extended BNF or EBNF). The grammars of the technical languages (programming languages, description languages, etc) are usually presented in an extended form

The EBNF rules can be represented as syntax diagrams. A syntax diagram can be conceived as the flow graph of a syntax analysis algorithm

The family LIB of free languages is closed with respect to regular operators (union, star, concatenation); therefore the extended free grammars have the same generating power as the non-extended ones

EXAMPLE: list of declarations of variable names

character testo1, testo2; *real* temp, result;
integer alfa, beta2, gamma;
 $\Sigma = \{c, i, r, v, ', ' , ';'\}$ where v is a variable name
 $\left((c | i | r) v (, v)^* ; \right)^+$
 $S \rightarrow SE | E \quad E \rightarrow AF; \quad A \rightarrow c | i | r \quad F \rightarrow v, F | v$

The grammar is longer than the regexp, and makes less evident the existence of a hierarchy of two lists. Moreover, the names A, E and F are arbitrary

A FREE EXTENDED GRAMMAR (EBNF) $G = \{ V, \Sigma, P, S \}$ contains exactly $|V|$ rules, each of the form $A \rightarrow \eta$, where η is a regexp over the alphabet $V \cup \Sigma$. For better clarity it is allowed to use derived regular operators, like for instance cross, optionality, repetition, etc

EXAMPLE: a simplified Algol-like language

A more compact but less readable version. The non-terminal B cannot be eliminated, as it is indispensable to generate nested structures

$$B \rightarrow b[D]Ie$$

$$D \rightarrow ((c | i | r)v(,v)^*;)^+$$

$$I \rightarrow F(;F)^*$$

$$F \rightarrow a | B$$

$$B \rightarrow b \left[((c | i | r)v(,v)^*;)^+ \right] F(;F)^* e$$

$$B \rightarrow b((c | i | r)v(,v)^*;)^* F(;F)^* e$$

$$B \rightarrow b((c | i | r)v(,v)^*;)^* (a | B)(; (a | B))^* e$$

DERIVATION AND SYNTAX TREE IN THE EXTENDED GRAMMARS

The right side of an extended rule is a regexp that defines an infinite set of strings

Such strings could be imagined as the right sides of a grammar G' , equivalent to G , except that it contains infinitely many rules

$$A \rightarrow (aB)^+$$

$$A \rightarrow aB \mid aBaB \mid \dots$$

DERIVATION RELATION IN AN EXTENDED GRAMMAR

Similarly, one can define multiple step derivations and the generated language

given the strings $\eta_1 \quad \eta_2 \in (\Sigma \cup V)^*$

say that η_2 *derives* (immediately) from η_1

$\eta_1 \Rightarrow \eta_2$ if the two strings are factored as follows

$\eta_1 = \alpha A \gamma, \quad \eta_2 = \alpha \mathcal{G} \gamma$ and there is a rule

$$A \rightarrow e : \quad e \stackrel{*}{\Rightarrow} \mathcal{G}$$

note that $\eta_1 \quad \eta_2 (\quad \alpha A \gamma \quad \alpha \mathcal{G} \gamma)$ do not contain either regular operators or parentheses, and e is a regexp

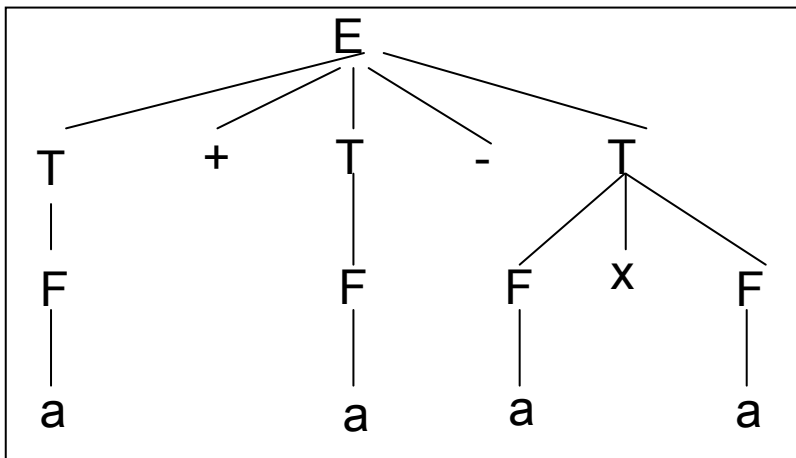
EXAMPLE: extended derivation for arithmetic expressions

The extended grammar G shown below generates the arithmetic expressions with the four standard operators (in infix notation), round brackets and the symbol a to represent generic variables or constants

$$E \Rightarrow [+|-]T((+|-)T)^* \quad T \rightarrow F((\times|/)F)^* \quad F \rightarrow (a|'('E')')$$

leftmost derivation

$$\begin{aligned} E &\Rightarrow T + T - T \Rightarrow F + T - T \Rightarrow a + T - T \Rightarrow a + F - T \Rightarrow \\ &\Rightarrow a + a - T \Rightarrow a + a - F \times F \Rightarrow a + a - a \times F \Rightarrow a + a - a \times a \end{aligned}$$



In EBNF, the arity of a tree node may be unlimited, in general

However, as the tree gets larger the depth gets lower

AMBIGUITY IN EXTENDED GRAMMARS

A non-extended but ambiguous grammar is ambiguous also when conceived as extended (as BNF is a special case of EBNF)

However, the presence of regexps in the rules may give rise to specific ambiguity forms

$a^*b \mid ab^*$ numbered as $a_1^*b_2 \mid a_3b_4^*$
is ambiguous because ab is
derivable as a_1b_2 or as a_3b_4
 $S \rightarrow a^*b \mid ab^*$ is ambiguous as well