# Free Grammars - A

*Translated and adapted by L. Breveglieri*

## CONTEXT-FREE GRAMMAR (FREE GRAMMAR)

CONTEXT-FREE LANGUAGE (FREE LANGUAGE)

LIMITATIONS OF REGULAR LANGUAGES

$$
\begin{array}{|l|}
\hline
begin \quad begin \quad begin \;\ldots end \quad end \quad end \\[4pt]
L_1 = \left\{ x \mid x = b^n e^n , n \geq 1 \right\} \\[4pt]
L_1 \quad \text{is not regular} \quad b^+ e^+ \quad (be)^+ \\
\hline
\end{array}
$$

SYNTAX – In order to define a language, one can use RULES that, after repeated application, allow to generate all and only the phrases of the language

The set of such rules constitutes a GENERATIVE GRAMMAR (or SYNTAX)

EXAMPLE − palindromes

$$L = \left\{ uu^R \mid u \in \{a,b\}^* \right\} = \{\varepsilon, aa, bb, abba, baab, ..., abbbba, ...\}$$

$frase \rightarrow \varepsilon$               − the empty string ε is a valid phrase

$frase \rightarrow a \; frase \; a$    − a phrase enclosed in *a*, *a* is a phrase

$frase \rightarrow b \; frase \; b$    − a phrase enclosed in *b*, *b* is a phrase

A derivation chain

$$frase \Rightarrow a \; frase \; a \Rightarrow ab \; frase \; ba \Rightarrow$$
$$\Rightarrow abb \; frase \; bba \Rightarrow abb\varepsilon bba = abbbba$$

The arrow "→" is a metasymbol, deserved to separate the left and right parts of a grammar rule

EXAMPLE: a list of palindromes

$$abba\ bbaabb\ aa$$

$$lista \to frase\ lista \qquad frase \to \varepsilon$$
$$lista \to frase \qquad\qquad frase \to a\ frase\ a$$
$$\qquad\qquad\qquad\qquad\qquad frase \to b\ frase\ b$$

Non-terminal symbols

lista (axiom), which defines the whole langauge phrase

frase (phrase), which defines the phrase components, i.e., the palindrome substrings that constitute the whole phrase

CONTEXT-FREE GRAMMAR (BNF - Backus Normal Form – of TYPE 2):
is defined by means of four entities

*V*         *non-terminal alphabet,* is a set of *non-terminal symbols* (or *syntactic classes*)
*Σ*         *terminal alphabet*, is the set of the *characters* that constitute the phrases
*P*         is a set of *syntactic rules* (also called *production rules* or simply *rules*)
*S* ∈ *V*     is a particular non-terminal, called *axiom*

> each rule of P is an ordered pair
> $$(X, \alpha), \;\; X \in V \;\; \text{e} \;\; \alpha \in (V \cup \Sigma)^{*}$$
> $$(X, \alpha) \in P \quad X \to \alpha$$
> $$X \to \alpha_1, X \to \alpha_2, \dots X \to \alpha_n$$
> $$X \to \alpha_1 \,|\, \alpha_2 \,|\, \dots \,|\, \alpha_n$$
> $$X \to \alpha_1 \cup \alpha_2 \cup \dots \cup \alpha_n$$

In order to avoid confusion, the metasymbols " → ", " | ", "∪ " and " ε " must not apper among the terminal and non-terminal symbols; moreover, the terminal and non-terminal alphabets must not share elements (i.e., must be disjoint sets)

1ˢᵗ CONVENTION

$$< if\_phrase > \rightarrow if \ < cond > then < phrase > else < phrase >$$

2ⁿᵈ CONVENTION

$$if\_phrase \rightarrow \textbf{if} \ cond \ \textbf{then} \ phrase \ \textbf{else} \ phrase$$

$$if\_phrase \rightarrow \text{'} if \text{'} \ cond \text{'} then \text{'} phrase \text{'} else \text{'} phrase$$

3ʳᵈ CONVENTION

$$F \rightarrow if \ C \ then \ D \ else \ D$$

CONVENTIONS USED IN THE FOLLOWING
- terminal symbols (or characters or letters) - { $a$, $b$, … }
- non-terminal symbols - { A, B, …}
- strings over Σ that contain only terminal characters - { $r$, $s$, …, $z$ }
- strings over (V U Σ) that contain terminal and non-terminal elements - { α, β,… }
- strings over V that contain only non-terminal symbols - σ

# DERIVATION AND GENERATED LANGUAGE

$$\beta, \gamma \in (V \cup \Sigma)^*$$

$\gamma$ derives from $\beta$ in the grammar $\quad G$

$\beta \underset{G}{\Rightarrow} \gamma \quad$ if $\quad A \to \alpha \quad$ is a rule of grammar $\quad G$

and $\quad \beta = hAd, \quad \gamma = h\alpha d$

$$\beta_0 \overset{n}{\Rightarrow} \beta_n \qquad \beta_0 \overset{*}{\Rightarrow} \beta_n \qquad \beta_0 \overset{+}{\Rightarrow} \beta_n$$

THE LANGUAGE GENERATED BY G WHEN STARTING FROM
A NON-TERMINAL "A" (one usually takes "S" for "A") IS $L_A(G)$

A *form* generated by G, starting
from the non-terminal A $\in$ V, is
a string $\alpha \in (V \cup \Sigma)^*$ such that

$$A \overset{*}{\Rightarrow} \alpha$$

$$L_A(G) = \left\{ x \in \Sigma^* \mid A \overset{+}{\Rightarrow} x \right\}$$

$$L(G) = L_S(G) = \left\{ x \in \Sigma^* \mid S \overset{+}{\Rightarrow} x \right\}$$

EXAMPLE (the structure of a book). Grammar $G_l$ generates the structure of a book: it contains a preface ( $f$ ) and a series (denoted by the non-term. $A$) of one or more chapters, each one of which starts with a chapter title ( $t$ ) and contains a series ( $B$ ) of one or more lines ( $l$ )

$$S \rightarrow f\ A$$
$$A \rightarrow A\ t\ B \mid t\ B$$
$$B \rightarrow l\ B \mid l$$

Non-term. $A$ generates the forms $t\ B\ t\ B$ and eventually the string $t\ l\ l\ t\ l \in L_A(G_l)$

Non-term. $S$ generates the sentential forms $f\ A\ t\ l\ B$ and $f\ t\ B\ t\ B$

The language generated by the non-terminal $B$ is $L_B(G_l) = l^+$

The language $L(G_l)$, generated by $G_l$, is defined by the regular expression $f\ (\ t\ l^+\ )^+$

A language is CONTEXT-FREE (or simply FREE) if and only if there exists a free grammar that generates it

Two grammars G e G' are equivalent if and only if they generate the same language, i.e., iff the equality $L(G) = L(G')$ holds

$G_l$

$$S \rightarrow fX$$

$$X \rightarrow XtY \mid tY$$

$$Y \rightarrow lY \mid l$$

$G_{l2}$

$$S \rightarrow f\,A$$

$$A \rightarrow A\,t\,B \mid t\,B$$

$$B \rightarrow B\,l \mid l$$

$$Y \underset{G_l}{\overset{n}{\Rightarrow}} l^n \quad \text{and} \quad B \underset{G_{l2}}{\overset{n}{\Rightarrow}} l^n \qquad \text{generate the same language} \qquad L_B = l^+$$

ERRONEOUS GRAMMARS AND USELESS RULES: in writing a grammar, one needs to pay attention to that every non-terminal is defined and to that each of them actually contributes to generate the strings in the language

A GRAMMAR IS IN REDUCED FORM (or simply is REDUCED) if BOTH (not either one) the following conditions hold

Every non-terminal A is reachable from the axiom, i.e., there exists a derivation as follows

$$S \stackrel{*}{\Rightarrow} \alpha A \beta$$

Every non-terminal A generates a non-empty set of strings

$$L_A(G) \neq \emptyset$$

REDUCTION OF THE GRAMMAR: there exists an algorithm in two phases

The first phase identifies the non-terminal symbols that are undefined

The second phase identifies those that are unreachable

PHASE 1 – construct the complement set DEF of the defined non-term. symbols

$$DEF = V \setminus UNDEF$$

Initially DEF is set to contain the non-terminal symbols that are expanded by terminal rules

$$DEF := \left\{ A \mid (A \to u) \in P, \text{ with } u \in \Sigma^* \right\}$$

Then the following transformation is applied repeatedly, until it converges

$$DEF := DEF \cup \left\{ B \mid (B \to D_1 D_2 \ldots D_n) \in P \right\}$$

Each symbol $D_i$ $(1 \leq i \leq n)$ either is a terminal or belongs to DEF

A third property is often required for a grammar to be in the reduced form, as follows

The grammar G may not have circular derivations, which are inessential and moreover give raise to ambiguity (think of ambiguous regexps …)

$$A \overset{+}{\Rightarrow} A$$

$$\text{if } x \text{ is derivable as follows } A \Rightarrow A \Rightarrow x$$

$$\text{then } x \text{ is derivable also as follows } A \Rightarrow x$$

In the following, it will be shown why and how a circular derivation necessarily gives raise to an ambiguous behaviour

EXAMPLES - NOT REDUCED

$$1) \ \left\{ S \to aASb, \quad A \to b \right\}$$

$$2) \ \left\{ S \to a, \quad A \to b \right\} \quad \left\{ S \to a \right\}$$

$$3) \ \left\{ S \to aASb \,|\, A, \quad A \to S \,|\, c \right\} \quad \left\{ S \to aSSb \,|\, c \right\}$$

CAUTION: circularity may raise from a null rule, as follows

$$X \to XY \,|\, \dots \quad Y \to \varepsilon \,|\, \dots$$

CAUTION: even if properly reduced, a grammar may still exhibit redundancy, and therefore may behave ambiguously, as follows

the rule pairs (1, 4) and (2, 5) generate the same phrases

$$1) \ \ S \to aASb \qquad 4) \ \ A \to c$$

$$2) \ \ S \to aBSb \qquad 5) \ \ B \to c$$

$$3) \ \ S \to \varepsilon$$

RECURSION AND INFINITY OF THE LANGUAGE

Almost all the formal languages of practical interest are infinite ($=$ contain infinitely many strings). What is the feature that enables a free grammar to generate infinitely many strings, i.e., an infinite language?

IN ORDER TO BE ABLE TO GENERATE INFINITELY MANY STRINGS, IT IS NECESSARY FOR THE RULES TO ALLOW DERIVATIONS OF UNBOUNDED LENGTH
THIS IS TO SAY THAT THE GRAMMAR IS RECURSIVE

$$A \overset{n}{\Rightarrow} xAy \quad n \geq 1 \qquad \text{is recursive}$$

$$\text{if} \quad n = 1 \qquad \text{is immediately recursive}$$

$$A \qquad \text{is a recursive nonterminal}$$

$$\text{if} \quad x = \varepsilon \qquad \text{is left recursive}$$

$$\text{if} \quad y = \varepsilon \qquad \text{is right recursive}$$

A NECESSARY AND SUFFICIENT CONDITION for a language L(G) to be infinite, where G is a grammar in the reduced form and does not have any circular derivations, is that G admits recursive derivations

PROOF OF NECESSITY: if there were not any recursive derivation, then every derivation would be of bounded length, hence L(G) would be finite (= would contain finitely many strings)

PROOF OF SUFFICENCE:
If there were such a derivation as

this derivation would be possible
but since G is in the reduced form,
it follows that A is reachable from S

hence from A at least one terminal
string $w$ can be derived

$$A \overset{n}{\Rightarrow} x \, A \, y$$

$$A \overset{+}{\Rightarrow} x^m \, A \, y^m \quad \text{for } m \geq 1 \text{ and } x, y \text{ not both empty}$$

$$S \overset{*}{\Rightarrow} u \, A \, v$$

$$A \overset{+}{\Rightarrow} w$$

whence there exist infinitely many derivations that yield different strings

$$S \overset{*}{\Rightarrow} u \, A \, v \overset{+}{\Rightarrow} u \, x^m \, A \, y^m \, v \overset{+}{\Rightarrow} u \, x^m \, w \, y^m \, v, \quad (m \geq 1)$$

A GRAMMAR DOES NOT CONTAIN RECURSION (or is RECURSION-FREE)
IF AND ONLY IF THE GRAPH OF THE BINARY RELATION *produce* IS ACYCLIC
(= does not contain any closed path)

G

EXAMPLE (for an infinite language)
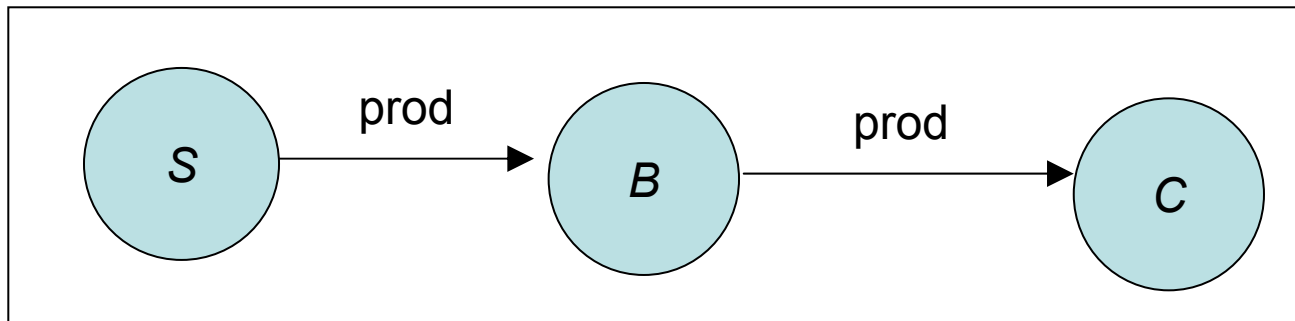
G generates the following finite language

$$S \rightarrow a\,B\,c$$
$$B \rightarrow a\,b \mid C\,a$$
$$C \rightarrow c$$

$$\{aabc, acac\}$$



$S$ — prod → $B$ — prod → $C$

EXAMPLE (arithmetic expression)

$$G = \left( \left\{ E, T, F \right\}, \left\{ i, +, *, ), ( \right\}, P, E \right)$$
$$E \rightarrow E + T \mid T \qquad T \rightarrow T * F \mid F \qquad F \rightarrow ( E ) \mid i$$
$$L(G) = \left\{ i, \, i + i + i, \, i * i, \, ( i + i ) * i, \, ... \right\}$$

F (factor)          is involved in an indirect recursion
E (expression)     is involved in an immediate left recursion
T (term)           is involved in an immediate left recursion

G is in the reduced form and does not contain any circular derivations,
therefore the generated language L(G) is infinite

# SYNTAX TREE AND CANONICAL DERIVATION

SYNTAX TREE: directed acyclic graph, such that for every pair of nodes there exists one and only one path (not necessarily directed) that connects them

THE SYNTAX TREE
> represents graphically the derivation process
> gives a *parent-child* relation or a *root-node-leaf* relation
> the sequence of leaves, scanned from left to right, is the so-called *frontier*
> the *degree* of a node (so-called node *arity*) is the length of the rule

SUBTREE with root N: the tree that has root N and includes all the descendants of N (the immediate siblings of N, the siblings of the siblings, and so on)

SYNTAX TREE: the root is the axiom and the frontier is the generated phrase

## grammar

## production subtrees

1. $E \rightarrow E + T$
2. $E \rightarrow T$
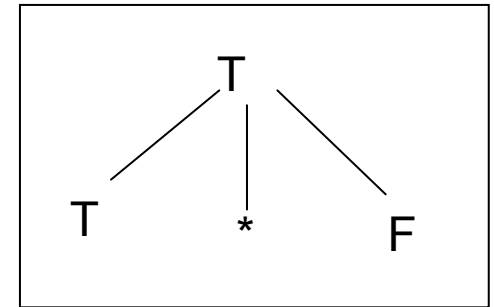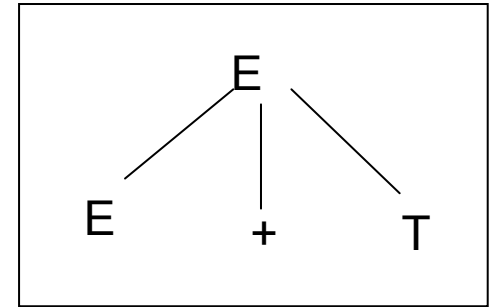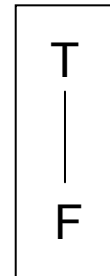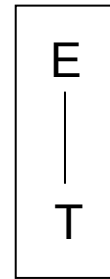3. $T \rightarrow T * F$
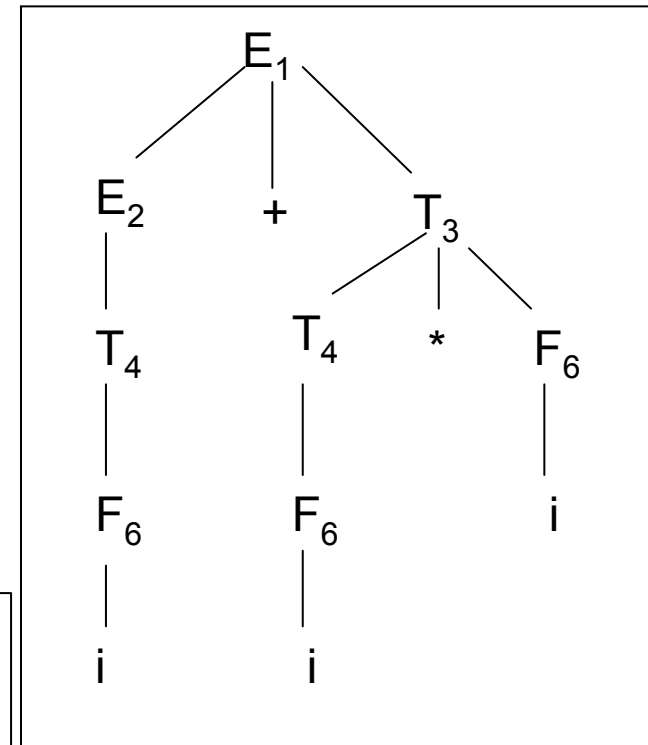4. $T \rightarrow F$
5. $F \rightarrow (E)$
6. $F \rightarrow i$

parenthesized representation of the syntax tree

$$[[[[i]_F]_T]_E + [[[i]_F]_T * [i]_F]_T]_E$$

LEFTMOST DERIVATION

$$E \underset{1}{\Rightarrow} E + T \underset{2}{\Rightarrow} T + T \underset{4}{\Rightarrow} F + T \underset{6}{\Rightarrow} i + T \underset{3}{\Rightarrow} i + T * F \underset{4}{\Rightarrow}$$
$$\underset{4}{\Rightarrow} i + F * F \underset{6}{\Rightarrow} i + i * F \underset{6}{\Rightarrow} i + i * i$$



syntax tree

RIGHTMOST DERIVATION

$$E \underset{1}{\Rightarrow} E + T \underset{3}{\Rightarrow} E + T * F \underset{6}{\Rightarrow} E + T * i \underset{4}{\Rightarrow} E + F * i \underset{6}{\Rightarrow} E + i * i \underset{2}{\Rightarrow}$$
$$\underset{2}{\Rightarrow} T + i * i \underset{4}{\Rightarrow} F + i * i \underset{6}{\Rightarrow} i + i * i$$

SKELETON TREE (only the frontier and the connections)

skeleton syntax tree

parenthesized representation of the syntax tree

$$[[[[i]]] + [[[i]] * [i]]]$$



CONDENSED SKELETON TREE (merge into one all the internal nodes aligned on a linear path, i.e., a path without branch points)

condensed
tree

$$[[i] + [[i] * [i]]]$$



condensed representation of the syntax tree

# LEFTMOST AND RIGHTMOST DERIVATION / FREE GRAMMAR

$$E \underset{s}{\overset{+}{\Rightarrow}} i + i * i$$

$$E \underset{d}{\overset{+}{\Rightarrow}} i + i * i$$

*legenda*

*s*: left
*d*: right

$$E \underset{s,d}{\overset{+}{\Rightarrow}} i + i * i$$

$$E \Rightarrow E + T \underset{d}{\Rightarrow} E + T * F \underset{s}{\Rightarrow} T + T * F \Rightarrow T + F * F \underset{d}{\Rightarrow} T + F * i \underset{s}{\Rightarrow}$$
$$\underset{s}{\Rightarrow} F + F * i \underset{d}{\Rightarrow} F + i * i \underset{d}{\Rightarrow} i + i * i$$

EVERY PHRASE OF A FREE GRAMMAR CAN BE GENERATED BY MEANS OF A LEFTMOST DERIVATION (and by a RIGHTMOST DERIVATION as well).

PARENTHESES LANGUAGE: artificial languages frequently contain nested structures, i.e., parentheses, where an element pair starts and ends some substructure (like for instance a substring), and where the pair may in turn contain a nested pair, and so on recursively down to arbitrary depth

| | |
|---|---|
| Pascal: | begin ... end |
| C: | { … } |
| XML: | < title > … < /title > |
| LaTeX: | \begin{equation} … \end{equation} |

Do not consider the specific way the marker symbols are encoded

The paradigm of parenthesis languages is known as the *Dyck Language*

Alphabet $\Sigma = \{ ')', '(', ']', '[' \}$  Phrases $()[[()[]]()]$

Parenthesis phrases can be equivalently defined by means of *cancellation rules*: repeatedly remove from the string any factor that consists of a pair of adjacent open and closed parentheses, as long as it is possible.
The original string is valid if and only if the final result is the *empty* string

$$[ \; ] \Rightarrow \varepsilon \quad ( \; ) \Rightarrow \varepsilon$$

DYCK LANGUAGE: open parentheses *a, b, ...,* closed parentheses *a', b', ...*

↓

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

$$S \rightarrow aSa'S \mid bSb'S \mid \varepsilon$$

NON-REGULAR LANGUAGE

↓

$a\,a\,\underbrace{a\,a'}\,a'\,a\,a\,\underbrace{a\,a'}\,a'\,a'\,a'$

$$L_1 = \{a^n c^n \mid n \geq 1\} = \{ac, aacc, ...\}$$

$$S \rightarrow aSc \mid ac$$

↓

$L_1$ is a subset of the Dyck language.
It does not admit more than one parenthesis nest

REGULAR COMPOSITION OF FREE LANGUAGES

The basic regular operations (union, concatenation and star), applied to free languages, still yield free languages

The family of free languaegs is closed with respect to the language operations *union*, *concatenation* and *star*

$$G_1 = \left( \Sigma_1, V_{N_1}, P_1, S_1 \right) \quad e \quad G_2 = \left( \Sigma_2, V_{N_2}, P_2, S_2 \right)$$

$$V_{N_1} \bigcap V_{N_2} = \varnothing \qquad S \notin (V_{N_1} \bigcup V_{N_2})$$

UNION

$$G = (\Sigma_1 \bigcup \Sigma_2, \{S\} \bigcup V_{N_1} \bigcup V_{N_2}, \{S \rightarrow S_1 \mid S_2\} \bigcup P_1 \bigcup P_2, S)$$

CONCATENATION

$$G = (\Sigma_1 \bigcup \Sigma_2, \{S\} \bigcup V_{N_1} \bigcup V_{N_2}, \{S \rightarrow S_1 S_2\} \bigcup P_1 \bigcup P_2, S)$$

STAR: G of $(L_1)^*$ is obtained by adding the rules $S \rightarrow S\ S_1 \mid \varepsilon$ to $G_1$

CROSS: G of $(L1)^+$ is obtained by adding the rules $S \rightarrow S\ S_1 \mid S_1$ to $G_1$

EXAMPLE: union of languages

$$L = \left\{ a^i b^i c^* \mid i \geq 0 \right\} \cup \left\{ a^* b^i c^i \mid i \geq 0 \right\} = L_1 \cup L_2$$

$G_1$

$S_1 \rightarrow XC$

$X \rightarrow aXb \mid \varepsilon$

$C \rightarrow cC \mid \varepsilon$

$G_2$

$S_2 \rightarrow AY$

$Y \rightarrow bYc \mid \varepsilon$

$A \rightarrow aA \mid \varepsilon$

$$\left\{ S \rightarrow S_1 \mid S^{''} \right\} \cup P_1 \cup P^{''}$$

CAUTION: if the hypothesis that the two non-terminal alphabets are disjoint does not hold, then the construction above yields a grammar that generates a superset of the real union language. For example, replacing $G_2$ with G'' would allow us to generate the invalid phrase shown aside

$G^{''}$

$S^{''} \rightarrow AX$

$X \rightarrow bXc \mid \varepsilon$

$A \rightarrow aA \mid \varepsilon$

$abcbc$

The family LIB of free languages is closed with respect to STRING MIRRORING

Given the grammar G of the language, the grammar $G_R$ that generates the mirror image of the language is obtained from G by mirroring the right member of every rule of G

REG and LIB are both closed with respect to union, concatenation and star

but later it will be proved that

*they do not behave in the same way as for complement and intersection*

AMBIGUITY: semantic versus syntactic
"Vedo un uomo in giardino con il cannocchiale"
(do I watch a man who has a spyglass or do I use a spyglass to watch a man ?)
"La pesca è bella" (the peach is fine or the sport of fishing is fine ?)
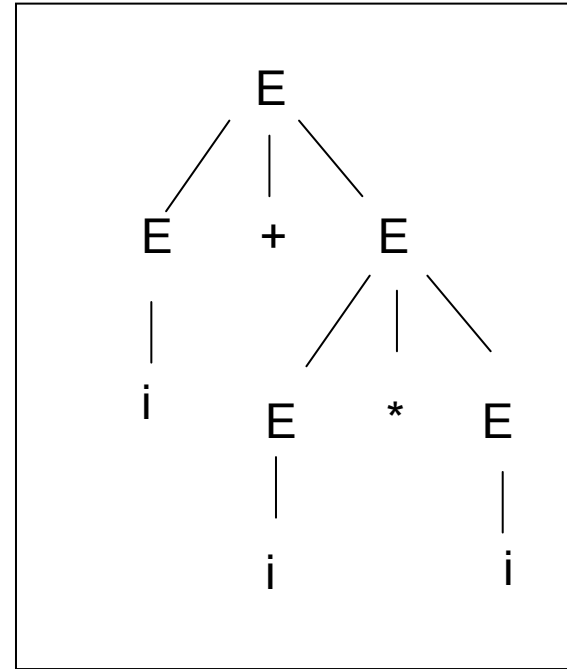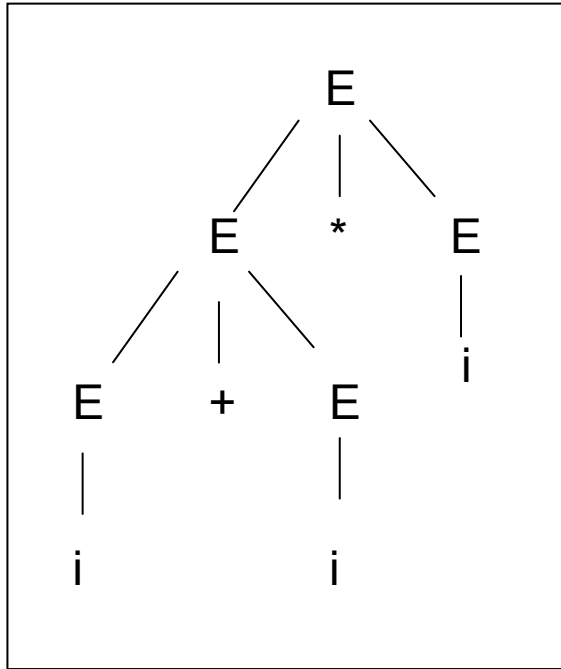"half baked chicken" (a chicken that is half baked or a half chicken that is baked ?)

Consider SYNTACTIC AMBIGUITY. A phrase *x* of the language defined by the grammar G is said to be ambiguous if it admits two or more different syntactic trees (or, equivalently, two different derivations). A grammar G is said to be ambiguous if it generates (at least) an ambiguous string

EXAMPLE – grammar G' of the arithmetic expressions

$$E \rightarrow E + E \mid E * E \mid (E) \mid i$$

$$E \Rightarrow E * E \Rightarrow E + E * E \Rightarrow i + E * E \Rightarrow i + i * E \Rightarrow i + i * i$$

$$E \Rightarrow E + E \Rightarrow i + E \Rightarrow i + E * E \Rightarrow i + i * E \Rightarrow i + i * i$$

The phrase *i + i * i* is ambiguous. Therefore grammar G' is itself ambiguous

In fact, grammar G' is not compliant with the standard convention
that multiplication must precede addition

Previously, an unambiguous grammar G that generates the arithmetic expressions
has been shown. G' is equivalent to G (that is, L(G) = L(G')), but is smaller than G.
However, grammar G' is ambiguous. This is often the case: simplifying the grammar
frequently ends up with having an ambiguous behaviour, in general

The AMBIGUITY DEGREE of a phrase *x* in the language L(G) is the number of different syntax trees that *x* admits. The ambiguity degree of a string may be infinite (only if there are nullable non-terminal symbols)
The AMBIGUITY DEGREE of a grammar G is the maximum ambiguity degree of the strings generated by G. While every string of L(G) may have a finite ambiguity degree, the ambiguity degree of G may still turn out to be unbounded

RELEVANT PROBLEM: decide whether a grammar G is ambiguous or not.
UNFORTUNATELY, THIS PROBLEM IS UNDECIDABLE: there does not exist any algorithm to decide whether a given grammar G is ambiguous or not.
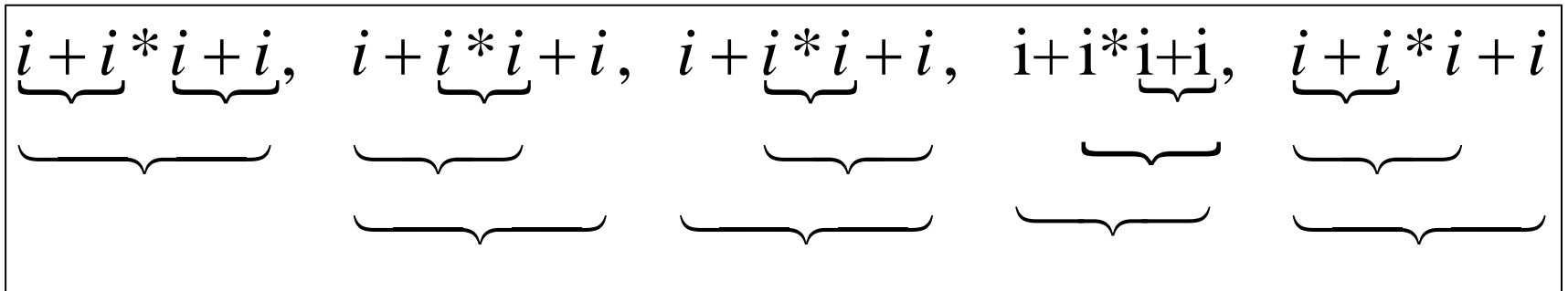This can be done for some grammars, but in general not for every grammar

AS AMBIGUITY MAY BE DIFFICULT TO IDENTIFY A POSTERIORI
IT IS ADVISABLE TO AVOID IT BY CONSTRUCTION (A FORTIORI)

EXAMPLE: resume the previous example

The phrase $i + i + i$ has ambiguity degree equal to 2

The phrase $i + i * i + i$ has ambiguity degree equal to 5 (see below)

$$\underbrace{i + \underbrace{i} * \underbrace{i + i}}, \quad i + \underbrace{i * i} + i, \quad i + \underbrace{i * i} + i, \quad i+i*\underbrace{i+i}, \quad \underbrace{i + i} * i + i$$

Typically, as the derived phrase gets longer, its ambiguity degree grows as well

AMBIGUITY OF TWO-SIDED RECURSION: a non-terminal symbol that is recursive both on the left and on the right (two-sided recursion) always allows two or more derivations, and thus gives rise to an ambiguous behaviour (see below)

$$A \overset{+}{\Rightarrow} A... \quad A \overset{+}{\Rightarrow} ...A$$

EXAMPLE 1: grammar $G_1$ generates the phrase $i + i + i$ in two different ways (by means of two different leftmost derivations)

ambiguous grammar

$$G_1 : \quad E \to E + E \mid i$$

notice that L is regular

non-ambiguous version

$$L(G_1) = i(+i)^*$$
$$E \to i + E \mid i$$
$$E \to E + i \mid i$$

EXAMPLE 2: ambiguity that originates form left and right recursion, separately

ambiguous grammar

ambiguity can be removed
by separately generating
the two lists

$$G_2 : \quad A \rightarrow aA \mid Ab \mid c$$

or one may choose to orderly generate
the two lists (first *a* then *b*, or viceversa)

$$L(G_2) = a^* c b^*$$
$$S \rightarrow AcB$$
$$A \rightarrow aA \mid \varepsilon$$
$$B \rightarrow bB \mid \varepsilon$$

$$L(G_2) = a^* c b^*$$
$$S \rightarrow aS \mid X$$
$$X \rightarrow Xb \mid c$$

EXAMPLE 3: the grammar that generates Polish postfix expressions that contain
addition and multiplication, has a left recursion (but not right) and is not ambiguous

$$S \rightarrow +SS \mid \times SS \mid i$$

AMBIGUITY OF UNION (maybe the simplest to understand)

If two languages $L_1(G_1)$ and $L_2(G_2)$ share some phrases (i.e., their intersection is not empty), the grammar G of the union language, constructed as shown before, is surely ambiguous

CAUTION: one needs to assume that the two non-terminal sets are disjoint, otherwise, as seen before, the union grammar would generate a superlanguage that strictly contains both languages

A phrase *x* that belongs to the union of $L_1$ and $L_2$, and admits two different derivations, the first only using rules of $G_1$ and the second only using rules of $G_2$, is ambiguous for the grammar G, as G contains both sets of rules. Only the phrases that belong to $L_1 \setminus L_2$ or to $L_2 \setminus L_1$, if any, would be generated in a non-ambiguous way

EXAMPLE 1: ambiguity of union may arise when in a programming language a special case, out of a general one, is treated by means of separate rules

$$E \to E + 1$$
$$E \to inc\ E$$

separate rule

$$E \to E + T \mid T \quad T \to V \mid C \quad V \to \quad 1$$
$$C \to 0 \mid 1B \mid ... \mid 9B \quad B \to 0B \mid 1B \mid ... \mid 9B \mid \varepsilon$$

basic grammar

EXAMPLE 2: ambiguity of union may arise when in a programming language the same operator has two different meanings. For example, in Pascal the operator "+" can indicate both integer addition and set-theoretic union

ambiguous grammar

$$E \to E + T \mid T \quad T \to V \quad V \to ...$$
$$E_{set} \to E_{set} + T_{set} \mid T_{set} \quad T_{set} \to V$$

EXAMPLE 3

$$G: \quad S \rightarrow bS \mid cS \mid D \quad D \rightarrow bD \mid cD \mid \varepsilon$$

$$L(G) = L_S(G) \bigcup L_D(G)$$

$$L_S(G) = \{b, c\}^* = L_D(G)$$

$$S \overset{+}{\Rightarrow} bbcD \Rightarrow bbc \quad S \Rightarrow D \overset{+}{\Rightarrow} bbcD \Rightarrow bbc$$

$$S \rightarrow bS \mid cS \mid \varepsilon$$

non-ambiguous version

EXAMPLE 4

$$S \rightarrow B \mid D \quad B \rightarrow bBc \mid \varepsilon$$

$$D \rightarrow dDe \mid \varepsilon$$

$B$ generates $b^n c^n, n \geq 0$

$D$ generates $d^n e^n, n \geq 0$

$\varepsilon$ is the only ambiguous phrase

$$S \rightarrow B \mid D \mid \varepsilon$$

$$B \rightarrow bBc \mid bc$$

$$D \rightarrow dDe \mid de$$

non-ambiguous version

# LANGUAGES THAT ARE INHERENTLY AMBIGUOUS

A language is said to be INHERENTLY AMBIGUOUS if every grammar that generates it is necessarily ambiguous, i.e., if all the equivalent grammars of the language happen to be ambiguous

EXAMPLE: here is a classical inherently ambiguous language

$$L = \left\{ a^i b^j c^k \mid (i, j, k \geq 0) \wedge ((i = j) \vee (j = k)) \right\}$$

$$L = \left\{ a^i b^i c^* \mid i \geq 0 \right\} \cup \left\{ a^* b^i c^i \mid i \geq 0 \right\} = L_1 \cup L_2$$

ambiguous union

Whatever grammar G is designed to generate L, the phrases $\varepsilon$, $abc$, $a^2b^2c^2$, ...
ALWAYS HAPPEN TO BE DERIVABLE IN TWO OR MORE WAYS.
This behaviour is caused by the very nature of the language L
and does not depend on the specific grammar G that generates L

# AMBIGUITY OF CONCATENATION

Concatenating two languages may give rise to ambiguity if there exists
a suffix of a phrase of the former language that is a prefix of a phrase
of the latter language

$$G = (\Sigma_1 \cup \Sigma_2, \{S\} \cup V_{N_1} \cup V_{N_2}, \{S \to S_1 S_2\} \cup P_1 \cup P_2, S)$$

Assume that $G_1$ and $G_2$ are not ambiguous, then G is ambiguous if there exist
two phrases x' $\in L_1$ and x" $\in L_2$, and a non-empty string $v$ such that

$$x' = u'v \wedge u' \in L_1 \quad x'' = vz'' \wedge z'' \in L_2$$

$$u'vz'' \in L_1.L_2 \quad \text{and is ambiguous}$$

$$S \Rightarrow S_1 S_2 \overset{+}{\Rightarrow} u'S_2 \overset{+}{\Rightarrow} u'vz''$$

$$S \Rightarrow S_1 S_2 \overset{+}{\Rightarrow} u'vS_2 \overset{+}{\Rightarrow} u'vz''$$

EXAMPLE 1 – concatenation of Dyck languages

$$\Sigma_1 = \{a, a', b, b'\} \quad \Sigma_2 = \{b, b', c, c'\}$$

$$aa'bb'cc' \in L = L_1 L_2$$

$$G(L): \quad S \rightarrow S_1 S_2$$

$$S_1 \rightarrow aS_1 a' S_1 \mid bS_1 b' S_1 \mid \varepsilon$$

$$S_2 \rightarrow bS_2 b' S_2 \mid cS_2 c' S_2 \mid \varepsilon$$

$$\overbrace{aa'bb'}^{S_1}\overbrace{cc'}^{S_2} \qquad \overbrace{aa'}^{S_1}\overbrace{bb'cc'}^{S_2}$$

In order to eliminate ambiguity, it is necessary to exclude the case when a suffix moves from the former language to the latter one. A solution consists of inserting a separator between the two languages. Such a separator may not belong to the alphabets of the two languages. The concatenation language $L_1 \# L_2$ is generated by the additional axiomatic rule $S \rightarrow S_1 \# S_2$

EXAMPLE 2 – encoding – the relationship between ambiguity and the uniqueness of encoding in information theory

A message is a sequence of symbols out of the set Γ = { A, B,…, Z }, with possible repetitions. Such symbols are then encoded into strings of letters out of the set of terminal symbols Σ. Most frequently Σ is the binary alphabet, Σ = { 0, 1 }

$$\Gamma = \left\{ \overbrace{A}^{01}, \overbrace{C}^{10}, \overbrace{E}^{11}, \overbrace{R}^{001}, \right\}$$

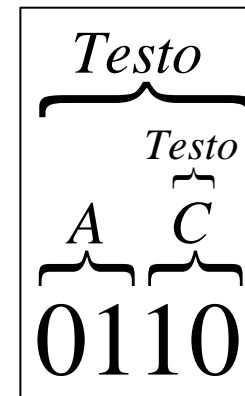$$ARRECA : 01\ 001\ 001\ 11\ 10\ 01$$

$$Testo \rightarrow ATesto \mid CTesto \mid ETesto \mid RTesto \mid A \mid C \mid E \mid R$$

$$A \rightarrow 01 \quad C \rightarrow 10 \quad E \rightarrow 11 \quad R \rightarrow 001$$

unambiguous grammar

This grammar is not ambiguous: every message encoded onto Σ admits only one syntax tree, and therefore it can be decoded in a unique way

A bad choice of the set of encoding strings causes the grammar to be ambiguous

ambiguous grammar

$$\Gamma = \left\{ \overbrace{A}^{00}, \overbrace{C}^{01}, \overbrace{E}^{10}, \overbrace{R,}^{010} \right\}$$

$$Testo \rightarrow ATesto \mid CTesto \mid ETesto \mid RTesto \mid A \mid C \mid E \mid R$$

$$A \rightarrow 00 \quad C \rightarrow 01 \quad E \rightarrow 10 \quad R \rightarrow 010$$

$$ARRECA = 00\ 010\ 010\ 10\ 01\ 00$$

$$ACAEECA = 00\ 01\ 00\ 10\ 10\ 01\ 00$$

The *theory of codes* studies these and also other aspects, to identify conditions ensuring that a code (= a set of encoding strings) is decodable in a unique way

OTHER AMBIGUOUS SITUATIONS – regexps may be themselves ambiguous

ambiguous grammar

EXAMPLE 1

Every phase that contains
two or more *c* is ambiguous.
To remove ambiguity,
one can impose that the mandatory *c* is the leftmost one

$$S \rightarrow DcD \quad D \rightarrow bD \mid cD \mid \varepsilon$$

$$\{b, c\}^* \, c \, \{b, c\}^*$$

unambiguous version

$$S \rightarrow BcD \quad D \rightarrow bD \mid cD \mid \varepsilon \quad B \rightarrow bB \mid \varepsilon$$

EXAMPLE 2 – Fixing the order of
application of the rules – The rule to
generate two *b* can be applied before
or after the rule that generates one *b*

$$S \rightarrow bSc \mid bbSc \mid \varepsilon$$

$$S \Rightarrow bbSc \Rightarrow bbbScc \Rightarrow bbbcc$$

$$S \Rightarrow bSc \Rightarrow bbbScc \Rightarrow bbbcc$$

ambiguous grammar

$$S \rightarrow bSc \mid D \quad D \rightarrow bbDc \mid \varepsilon$$

unambiguous version

# AMBIGUITY IN CONDITIONAL PHRASES

$$S \rightarrow if\ b\ then\ S\ else\ S\ |\ if\ b\ then\ S\ |\ a$$

$$if\ b\ then\ if\ b\ then\ a\ else\ a$$

$$if\ b\ then\ if\ b\ then\ a\ else\ a$$

so-called problem
of the
"dangling else"

$$S \rightarrow S_E\ |\ S_T \quad S_E \rightarrow if\ b\ then\ S_E\ else\ S_E\ |\ a$$
$$S_T \rightarrow if\ b\ then\ S_E\ else\ S_T\ |\ if\ b\ then\ S$$

Only $S_E$ can precede *else*

unambiguous version 1

$$S \rightarrow if\ b\ then\ S\ else\ S\ end\_if\ |$$
$$|\ if\ b\ then\ S\ end\_if\ |\ a$$

unambiguous version 2

Use the additional
keyword *end_if* to
mark the end of the
conditional phrase