

# Free Grammars - I

*Translated and adapted by L. Breveglieri*

## CONTEXT-FREE GRAMMAR (FREE GRAMMAR)

## CONTEXT-FREE LANGUAGE (FREE LANGUAGE)

## LIMITATIONS OF REGULAR LANGUAGES

*begin begin begin ... end end end*

$$L_1 = \{x \mid x = b^n e^n, n \geq 1\}$$

$L_1$  is not regular  $b^+ e^+ (be)^+$

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

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

$frase \rightarrow \varepsilon$  – the empty string  $\varepsilon$  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

$$\begin{aligned} frase &\Rightarrow a \ frase \ a \Rightarrow ab \ frase \ ba \Rightarrow \\ &\Rightarrow abb \ frase \ bba \Rightarrow abb\varepsilon bba = abbbba \end{aligned}$$

The arrow “ $\rightarrow$ ” is a metasymbol, deserved to separate the left and right parts of a grammar rule

EXAMPLE: a list of palindromes

*abba bbaabb aa*

*lista*  $\rightarrow$  *frase lista*      *frase*  $\rightarrow$   $\varepsilon$

*lista*  $\rightarrow$  *frase*      *frase*  $\rightarrow$  *a frase a*

*frase*  $\rightarrow$  *b frase b*

Non-terminal symbols

*lista* (axiom), which defines the whole language phrase

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

EXAMPLE: the metalanguage of regular expressions

Regular expressions that define the regular languages over the terminal alphabet  $\Sigma = \{ a, b \}$ , are a language themselves, and in particular they are strings over the alphabet  $\Sigma_{r.e.} = \{ a, b, \cup, *, \emptyset, (, ) \}$

Syntax that generates  $G_{r.e.}$ .

- |                                 |  |
|---------------------------------|--|
| 1. $espr \rightarrow \emptyset$ | 4. $espr \rightarrow (espr \cup espr)$ |
| 2. $espr \rightarrow a$         | 5. $espr \rightarrow (espr\ espr)$     |
| 3. $espr \rightarrow b$         | 6. $espr \rightarrow (espr)^*$         |

Derivation example

$$\begin{aligned} espr &\Rightarrow (espr \cup espr) \Rightarrow ((espr\ espr) \cup espr) \Rightarrow ((a\ espr) \cup espr) \Rightarrow \\ &\Rightarrow ((a\ (espr)^*) \cup espr) \Rightarrow ((a\ ((espr \cup espr))^* \cup espr)) \Rightarrow \\ &\Rightarrow ((a\ ((a \cup espr))^* \cup espr)) \Rightarrow ((a\ ((a \cup b))^* \cup espr)) \Rightarrow \\ &\Rightarrow ((a\ ((a \cup b))^* \cup espr)) \Rightarrow ((a\ ((a \cup b))^* \cup b)) \\ L(a\{ab\}^* \cup b) &= \{b, a, aa, ab, aaa, aba, \dots\} \end{aligned}$$

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*)  
 $\Sigma$         *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*)  
 $S \in V$     is a particular non-terminal, called *axiom*

each rule of  $P$  is an ordered pair  
 $(X, \alpha)$ ,  $X \in V$  e  $\alpha \in (V \cup \Sigma)^*$   
 $(X, \alpha) \in P \quad X \rightarrow \alpha$   
 $X \rightarrow \alpha_1, X \rightarrow \alpha_2, \dots, X \rightarrow \alpha_n$   
 $X \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$   
 $X \rightarrow \alpha_1 \cup \alpha_2 \cup \dots \cup \alpha_n$

In order to avoid confusion, the metasymbols “ $\rightarrow$ ”, “ $\mid$ ”, “ $\cup$ ” and “ $\varepsilon$ ” must not appear among the terminal and non-terminal symbols; moreover, the terminal and non-terminal alphabets must not share elements (i.e., must be disjoint sets)

## 1<sup>st</sup> CONVENTION

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

## 2<sup>nd</sup> CONVENTION

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

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

## 3<sup>rd</sup> CONVENTION

$F \rightarrow \text{if } C \text{ then } D \text{ else } D$

## CONVENTIONS USED IN THE FOLLOWING

- terminal symbols (or characters or letters) -  $\{ a, b, \dots \}$
- non-terminal symbols -  $\{ A, B, \dots \}$
- strings over  $\Sigma$  that contain only terminal characters -  $\{ r, s, \dots, z \}$
- strings over  $(V \cup \Sigma)$  that contain terminal and non-terminal elements -  $\{ \alpha, \beta, \dots \}$
- strings over  $V$  that contain only non-terminal symbols -  $\sigma$

TERMINAL: the right part contains terminals or empty string	$\rightarrow u \mid \varepsilon$
EMPTY: the right part is the empty string	$\rightarrow \varepsilon$
INITIAL: the left part is the axiom	$S \rightarrow$
RECURSIVE: the left part appears in the right part	$A \rightarrow \alpha A \beta$
LEFT RECURSIVE: the left part is prefix of the right part	$A \rightarrow A \beta$
RIGHT RECURSIVE: the left part is suffix of the right part	$A \rightarrow \beta A$
TWO-SIDED RECURSIVE: superposition of previous cases	$A \rightarrow A \beta A$
COPY (or RENAMING): simply a change of name	$A \rightarrow B$
LINEAR: the right part contains at most one nonterminal	$\rightarrow u B v \mid w$
RIGHT LINEAR: like linear, but the nonterminal is suffix	$\rightarrow u B \mid w$
LEFT LINEAR: like linear, but the nonterminal is prefix	$\rightarrow B v \mid w$
CHOMSKY NORMAL: either one term. or two nonterm.s	$\rightarrow BC \mid a$
GREIBACH NORMAL: always a prefix term., max length 2	$\rightarrow a \sigma \mid b$
OPERATOR NORMAL: always an infix terminal	$\rightarrow A a B$



## DERIVATION AND GENERATED LANGUAGE

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

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

$\beta \xRightarrow[G]{\quad} \gamma$  if  $A \rightarrow \alpha$  is a rule of grammar  $G$

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

$$\beta_0 \xRightarrow{n} \beta_n \quad \beta_0 \xRightarrow{*} \beta_n \quad \beta_0 \xRightarrow{+} \beta_n$$

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

$$A \xRightarrow{*} \alpha$$

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

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

If the nonterminal  $A$  is the axiom, then the form is said to be *sentential*

A *phrase* is simply a sentential form that contains only terminal characters

EXAMPLE (the structure of a book). Grammar  $G_1$  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 (  $l$  )

$$\begin{aligned} S &\rightarrow f A \\ A &\rightarrow A t B \mid t B \\ B &\rightarrow l B \mid l \end{aligned}$$

Non-term.  $A$  generates the forms  $t B t B$  and eventually the string  $t l t l \in L_A(G_1)$

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  $L_B(G_1) = l^+$

The language  $L(G_1)$ , generated by  $G_1$ , 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

In the example before, the free language  $L(G_1)$  happens to be regular as well, but in the following it will be proved that in general the family LIB of context-free languages strictly contains that of regular languages

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

$G_1$

$$S \rightarrow fX$$

$$X \rightarrow XtY \mid tY$$

$$Y \rightarrow lY \mid l$$

$G_{l2}$

$$S \rightarrow f A$$

$$A \rightarrow At B \mid t B$$

$$B \rightarrow Bl \mid l$$

$$Y \xRightarrow[n]{G_l} l^n \quad \text{and} \quad B \xRightarrow[n]{G_{l2}} l^n \quad \text{generate the same language} \quad L_B = l^+$$

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

$$\begin{array}{c} * \\ S \Rightarrow \alpha A \beta \end{array}$$

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

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

Then the following transformation is applied repeatedly, until it converges

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

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

In every iteration of the first phase, two events are possible

- 1) new non-terminal symbols are identified, which have as their right hand side a string that contains terminal symbols or defined non-terminal symbols; such new symbols are added to the set DEF
- 2) the DEF set does not change, which means that convergence has been reached, and therefore the algorithm terminates

The non-terminal symbols that belong to the UNDEF set can be eliminated, along with the rules wherein they appear (in the left or right side)

PHASE 2 – The identification of the non-terminal symbols that are reachable starting from the axiom S can be reduced to the problem of the existence of a path in the graph associated with the following binary relation *produce*

$$\begin{array}{l} \textit{produce} \\ A \rightarrow B \text{ if and only if} \\ A \rightarrow \alpha B \beta \text{ with } A \neq B \quad \alpha, \beta \text{ are any strings} \end{array}$$

The non-terminal C is reachable from the axiom S if and only if the graph of the binary relation *produce* contains a path from S to C

The unreachable non-terminal symbols can be removed from the grammar, along with the rules wherein they appear (in the left or right side)

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 rise to ambiguity (think of ambiguous regexps ...)

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

if  $x$  is derivable as follows  $A \Rightarrow A \Rightarrow x$

then  $x$  is derivable also as follows  $A \Rightarrow x$

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

## EXAMPLES - NOT REDUCED

- 1)  $\{ S \rightarrow aASb, A \rightarrow b \}$
- 2)  $\{ S \rightarrow a, A \rightarrow b \} \quad \{ S \rightarrow a \}$
- 3)  $\{ S \rightarrow aASb \mid A, A \rightarrow S \mid c \} \quad \{ S \rightarrow aSSb \mid c \}$

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

$$X \rightarrow XY \mid \dots \quad Y \rightarrow \varepsilon \mid \dots$$

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)  $S \rightarrow aASb$
- 2)  $S \rightarrow aBSb$
- 3)  $S \rightarrow \varepsilon$
- 4)  $A \rightarrow c$
- 5)  $B \rightarrow c$



## 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

$A \overset{n}{\Rightarrow} xAy$	$n \geq 1$	is recursive
if $n = 1$		is immediately recursive
$A$		is a recursive nonterminal
if $x = \varepsilon$		is left recursive
if $y = \varepsilon$		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

$$\begin{array}{l} \overset{n}{A \Rightarrow x A y} \\ + \\ A \Rightarrow x^m A y^m \text{ for } m \geq 1 \text{ and } x, y \text{ not both empty} \\ * \\ S \Rightarrow u A v \\ + \\ A \Rightarrow w \end{array}$$

whence there exist infinitely many derivations that yield different strings

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

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)

EXAMPLE (for an infinite language)

G generates the following finite language

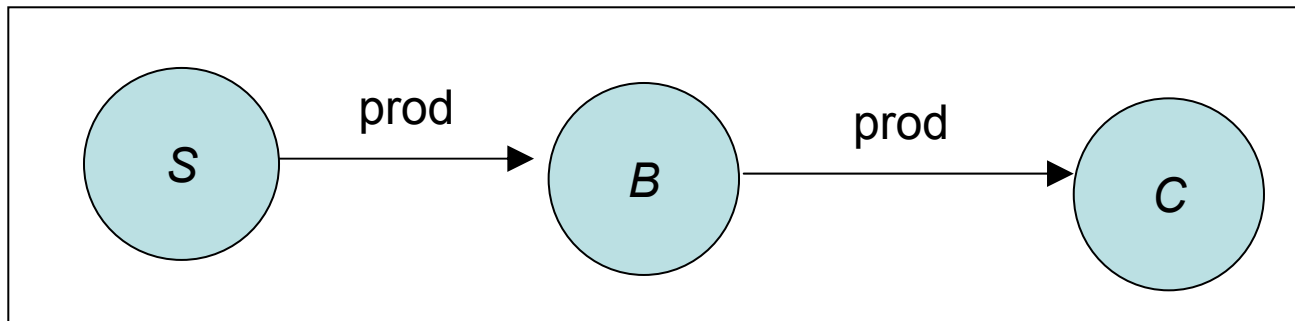
$\{aabc, acac\}$

G

$S \rightarrow a B c$

$B \rightarrow a b \mid C a$

$C \rightarrow c$



### EXAMPLE (arithmetic expression)

$$\begin{aligned} G &= ( \{ E, T, F \}, \{ i, +, *, ), ( \}, P, E ) \\ E &\rightarrow E + T \mid T \quad T \rightarrow T * F \mid F \quad F \rightarrow ( E ) \mid i \\ L(G) &= \{ i, i + i + i, i * i, (i + i) * i, \dots \} \end{aligned}$$

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

