

# Context free grammars - I

*Prof. A. Morzenti*

**NB:** many parts are well known from previous courses  
I will run quickly through them and discuss more deeply the really new ones

## LIMITS OF REGULAR LANGUAGES

Simple languages such as  $L = \{ a^n b^n \mid n > 0 \}$

representing basic syntactic structures like

begin begin ... begin ... end ... end end

are ***not*** regular

e.g.,  $b^+ e^+$  does not satisfy the constraint  $\#b = \#e$

$(be)^+$  does not ensure nesting

GRAMMARS – a more powerful means to define languages

through ***rewriting rules***

language phrases generated through repeated application of rules

The grammar is characterized by its set of rules

Example – Language of *palindromes*

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

a palindrome is...

$P \rightarrow \varepsilon$                       an empty palindrome

$P \rightarrow a P a$                       a palindrome surrounded by two  $a$ 's

$P \rightarrow b P b$                       a palindrome surrounded by two  $b$ 's

A chain of *derivation steps*:

$$P \Rightarrow a P a \Rightarrow ab P ba \Rightarrow abb P bba \Rightarrow abb \varepsilon bba = abbbba$$

**look out: distinguish the two *metasymbols***

$\rightarrow$  separates the left and right part of a rule

$\Rightarrow$  derivation *relation* (rewriting)

Example: a non-empty *list of* palindromes, ex: *abba bbaabb aa*

$$L \rightarrow P L$$

$$L \rightarrow P$$

$$P \rightarrow \varepsilon \quad P \rightarrow a P a \quad P \rightarrow b P b$$

***non terminal symbols:***

- $L$  (axiom, or start symbol)
- $P$  (defines the component palindrome substrings)

derivation of the string *abba bbaabb aa*

(spaces added for readability, nonterm to be expanded underlined)

$$\begin{aligned} L &\Rightarrow P \underline{L} \Rightarrow P P \underline{L} \Rightarrow \underline{P} P P \Rightarrow a \underline{P} a P P \\ &\Rightarrow ab \underline{P} ba P P \Rightarrow abba \underline{P} P \Rightarrow abba b \underline{P} b P \\ &\Rightarrow abba bb \underline{P} bb P \Rightarrow abba bba \underline{P} abb P \\ &\Rightarrow abba bbaabb \underline{P} \\ &\Rightarrow abba bbaabb a \underline{P} a \Rightarrow abba bbaabb aa \end{aligned}$$

## CONTEXT-FREE GRAMMAR

### (BNF - Backus Normal Form – TYPE 2 – FREE GRAMMAR)

defined by four entities

1.  $V$ , *non terminal alphabet*, is the set of nonterminal symbols
2.  $\Sigma$ , *terminal alphabet*, is the set of the symbols of which phrases/sentences are made
3.  $P$ , is the set of *rules* or *productions*
4.  $S \in V$ , is the specific nonterminal, called the *axiom (Start)*, from which derivations start

a rule is written as  $X \rightarrow \alpha$ , with  $X \in V$  and  $\alpha \in (V \cup \Sigma)^*$

rules with the same nonterminal  $X$ :  $X \rightarrow \alpha_1, X \rightarrow \alpha_2, \dots X \rightarrow \alpha_n$

can be written in brief as  $X \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$  (or  $X \rightarrow \alpha_1 \cup \alpha_2 \cup \dots \cup \alpha_n$ )

$\alpha_1, \alpha_2, \dots, \alpha_n$  are called the *alternatives* of  $X$

To avoid confusion,

the metasymbols ' $\rightarrow$ ', ' $\mid$ ', ' $\cup$ ', ' $\varepsilon$ ' cannot be terminal symbols,

terminal and nonterminal alphabets must be disjointed

NOTATIONS to distinguish terminal and nonterminal symbols

- angle brackets:

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

- **bold italic**:

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

- quotation marks ‘ ’

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

- upper- versus lower-case:

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

WE USUALLY ADOPT THESE CONVENTIONS:

- terminal characters -  $\{a, b, \dots\}$
- nonterminal characters -  $\{A, B, \dots\}$
- strings  $\in \Sigma^*$  (only terminals) -  $\{r, s, \dots, z\}$
- strings  $\in (V \cup \Sigma)^*$  (terminals and nonterminals) -  $\{\alpha, \beta, \dots\}$
- strings  $\in V^*$  (only **nonterminals**) -  $\sigma$

## TYPES OF RULES (RP = right part, LP = left part)

<b><u>Terminal</u></b> : RP contains only terminals, or the empty string	$\rightarrow u \mid \varepsilon$
<b><u>Empty (or null)</u></b> : RP is empty	$\rightarrow \varepsilon$
<b><u>Initial / Axiomatic</u></b> : LP is the axiom	$S \rightarrow$
<b><u>Recursive</u></b> : LP occurs in RP	$A \rightarrow \alpha A \beta$
<b><u>Left-recursive</u></b> : LP is prefix of RP	$A \rightarrow A \beta$
<b><u>Right-recursive</u></b> : LP is suffix of RP	$A \rightarrow \alpha A$
<b><u>Left- and right-recursive</u></b> : conjunction of two previous cases	$A \rightarrow A \beta A$
<b><u>Copy or categorization</u></b> : RP is a single nonterminal	$A \rightarrow B$
<b><u>Linear</u></b> : at most one nonterminal in RP	$\rightarrow u B v \mid w$
<b><u>Right-linear</u></b> (type 3): linear + nonterminal is suffix	$\rightarrow u B \mid w$
<b><u>Left-linear</u></b> (type 3): linear + nonterminal is prefix	$\rightarrow B v \mid w$
<b><u>Homogeneous normal</u></b> : $n$ nonterminals or just one terminal	$\rightarrow A_1 \dots A_n \mid a$
<b><u>Chomsky normal</u></b> (or homogeneous of degree 2): two nonterminals or just one terminal	$\rightarrow BC \mid a$
<b><u>Greibach normal</u></b> : one terminal possibly followed by nonterminals	$\rightarrow a \sigma \mid b$
<b><u>Operator normal</u></b> : two nonterminals separated by a terminal (operator); more generally, strings devoid of adjacent nonterminals	$\rightarrow A a B$

## DERIVATIONS AND GENERATED LANGUAGE

Def. of *derivation relation* ' $\Rightarrow$ '

for  $\delta, \eta \in (V \cup \Sigma)^*$ ,  $A \rightarrow \alpha$  rule of  $G$

$\delta A \eta$  derives  $\delta \alpha \eta$  for grammar  $G$ ,  $\delta A \eta \xRightarrow{G} \delta \alpha \eta$ , or  $\delta A \eta \Rightarrow \delta \alpha \eta$

the rule  $A \rightarrow \alpha$  is applied in that *derivation step*, and  $\alpha$  *reduces to*  $A$

**power, reflexive and transitive closure of ' $\Rightarrow$ '**

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

If  $A \xRightarrow{*} \alpha$   $\alpha \in (V \cup \Sigma)^*$  called *string form generated by  $G$*

If  $S \xRightarrow{*} \alpha$   $\alpha$  called *sentential* or *phrase form*

If  $S \xRightarrow{*} s$   $s \in \Sigma^*$ ,  $s$  is called *phrase* or *sentence*

LANGUAGE GENERATED FROM  
NONTERMINAL  $A$  OR FROM AXIOM  $S$

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

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



Example: Grammar  $G_l$  generates the structure of a book: it contains

- a front page ( $f$ )
- a series (denoted by the nonterm.  $A$ ) of one or more chapters by means of a *left recursion*
- a chapter has a title  $t$  and a sequence  $B$  (right recursion) of one or more lines  $l$

$$\begin{array}{l} S \rightarrow fA \\ A \rightarrow AtB \mid tB \\ B \rightarrow lB \mid l \end{array}$$

from  $A$  one generates the string form  $tBtB$  and the phrase  $tlltl \in L_A(G_l)$

from  $S$  one generates the phrase forms  $fAtlB$ ,  $ftBtB$

The language generated from  $B$  is  $L_B(G_l) = l^+$

$L(G_l)$ , being generated by the context free grammar  $G_l$ , is **context free** or **free**

Notice: it is also regular, because it is defined also by the regular expression  $f(tl^+)^+$

Two grammars  $G$  and  $G'$  are equivalent if they generate the same language, that is,  $L(G) = L(G')$

$$\begin{array}{c}
 G_l \\
 \hline
 S \rightarrow fX \\
 X \rightarrow XtY \mid tY \\
 Y \rightarrow lY \mid l
 \end{array}$$

$$\begin{array}{c}
 G_{l_2} \\
 \hline
 S \rightarrow fA \\
 A \rightarrow AtB \mid tB \\
 B \rightarrow Bl \mid l
 \end{array}$$

The first rules have the same n.t., only renamed

$Y \xRightarrow{n} l^n$  in  $G_l$  and  $B \xRightarrow{n} l^n$  in  $G_{l_2}$  generate the same language  $L_B = l^+$  by means of *right* or *left* recursion

# ERRONEOUS GRAMMARS AND USELESS RULES

A grammar  $G$  is *clean* (or *reduced*) iff for every nonterminal  $A$  :

1.  $A$  is *reachable* from the axiom  $S$ , and hence contributes to the generation of the language; that is, there exists a derivation

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

2.  $A$  is *defined*, that is, it generates a non-empty language (we are not interested in the language  $\emptyset$ )

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

NB:  $L_A(G) = \emptyset$  includes also the case when no derivation from  $A$  terminates with a terminal string  $s$  (i.e.  $s \in \Sigma^*$ ), e.g.:  $P = \{ S \rightarrow aA, A \rightarrow bS \}$

GRAMMAR CLEANING: two steps algorithm:

The FIRST PHASE builds the set **UNDEF** of undefined nonterminals

The SECOND PHASE builds the set of unreachable nonterminals

PHASE 1- We first build the **complement** set  $DEF = V \setminus UNDEF$

$DEF$  is **initialized** from the *terminal rules* (the n.t. that immediately generate a terminal string)

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

The following **update** is repeated until a **fixed point** is reached :

$$DEF := DEF \cup \{ B \mid (B \rightarrow D_1 D_2 \dots D_n) \in P \wedge \forall i (D_i \in DEF \cup \Sigma) \}$$

Every  $D_i$  is already in  $DEF$  or it is a terminal

In algebra, the **fixed point** of a transformation

is an object that is transformed into itself

At each iteration, two cases can occur:

1. New nonterm. are found having the RP all with defined nonterm. or term., or
2. No new nonterm. is found, algorithm terminates (a *fixed point* has been reached)

nonterminals  $\in UNDEF = V \setminus DEF$  are eliminated

PHASE 2 – Consider the *produce* relation, defined as

$A$  produce  $B$  iff  $(A \rightarrow \alpha B \beta) \in P$ , with  $A \neq B$   $\alpha, \beta$  any string

$C$  is *reachable* from  $S$  iff there exists, in the graph of the produce relation, a path from  $S$  to  $C$

nonterminals that are not reachable can be eliminated

often another requirement is added for cleanness condition of a grammar  $G$  :

3.  $G$  must not allow for *circular derivations*: they are not essential and introduce *ambiguity*

if  $A \Rightarrow^+ A$  then

if the derivation  $A \Rightarrow^+ x$  is possible

then also  $A \Rightarrow^+ A \Rightarrow^+ x$  and many other similar ones exist

NB: *circular derivations* must not be confused with *recursive rules* and *derivations* !!

## EXAMPLES OF GRAMMARS THAT ARE NOT CLEAN

- 1)  $\{ S \rightarrow aASb, A \rightarrow b \}$  ( $S$  does not generate any phrase, i.e.,  $L(S)=\emptyset$ )
- 2)  $\{ S \rightarrow a, A \rightarrow b \}$  ( $A$  not reachable)      ( $\{ S \rightarrow a \}$  equiv. clean version)
- 3)  $\{ S \rightarrow aASb \mid A, A \rightarrow S \mid b \}$  (circular on  $S$  and  $A$ )      ( $\{ S \rightarrow aSSb \mid b \}$  equiv. clean)

circularity can also derive from an empty rule

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

NB: even if clean, a grammar can have *redundant rules* (leading to ambiguity)

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

$$\begin{array}{ll} 1. S \rightarrow aASb & 4. A \rightarrow c \\ 2. S \rightarrow aBSb & 5. B \rightarrow c \\ 3. S \rightarrow \varepsilon & \end{array}$$

## RECURSION AND LANGUAGE *INFINITY*

most interesting languages are infinite

but what determines the ability of a grammar to generate an infinite language?

infinity of the language implies unbounded phrase length

therefore the grammar must be recursive

a *derivation*  $A \Rightarrow^n xAy$   $n \geq 1$  is *recursive*

if  $n = 1$  it is *immediately recursive*

$A$  is a *recursive nonterminal*

if  $x = \varepsilon$  then it is *left recursive* (l.r. derivation, l.r. nonterminal)

if  $y = \varepsilon$  then it is *right recursive* (r.r. derivation, r.r. nonterminal)

NB: circularity and recursiveness are (very) different notions

a grammar may be recursive (admit recursive derivations) but not circular

circular  $\Rightarrow$  recursive but it is **not** true that recursive  $\Rightarrow$  circular



**necessary and sufficient condition** for language  $L(G)$  to be infinite,  
 assuming  $G$  clean and devoid of circular derivations,  
 is that  $G$  allows for recursive derivations

**necessary condition:** if no recursive derivation was possible,  
 then every derivation would have limited length hence  $L(G)$  would be finite

**sufficient condition:**

$A \xRightarrow{n} xAy$  implies  $A \xRightarrow{+} x^m Ay^m$

for any  $m \geq 1$  with  $x, y \in \Sigma^*$  not both empty (because grammar is not circular)

Furthermore  $G$  clean implies

$S \xRightarrow{*} uAv$  ( $A$  reachable from  $S$ )

and  $A \xRightarrow{+} w$  (derivation from  $A$  terminates successfully)

therefore there exist nonterminals that generate an infinite language

$$S \xRightarrow{*} uAv \xRightarrow{+} ux^m Ay^m v \xRightarrow{+} ux^m wy^m v, (\forall m \geq 1)$$

a grammar does not have recursive derivations

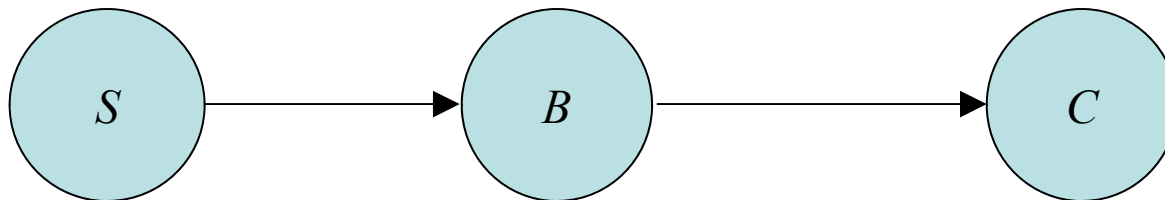
$\Leftrightarrow$  (if and only if)

the graph of the *produce* relation has no circuits

Example

$$\begin{array}{l} S \rightarrow aBc \\ B \rightarrow ab \mid Ca \\ C \rightarrow c \end{array}$$

finite language:  $\{ aabc, acac \}$



## Example (arithmetic expressions)

$$G = \left( \underbrace{\{E, T, F\}}_{\text{non term.}}, \underbrace{\{i, +, *, ), (\}}_{\text{term.}}, \underbrace{P}_{\text{productions}}, \underbrace{E}_{\text{axiom}} \right)$$

$$P = \{E \rightarrow E + T \mid T, \quad T \rightarrow T * F \mid F, \quad F \rightarrow (E) \mid i\}$$

$$L(G) = \{i, i + i + i, \quad i * i, \quad (i + i) * i, \quad \dots\}$$

$F$  (factor) has indirect recursion (non immediate)

$E$  (expression) has immediate left recursion and non immediate recursion

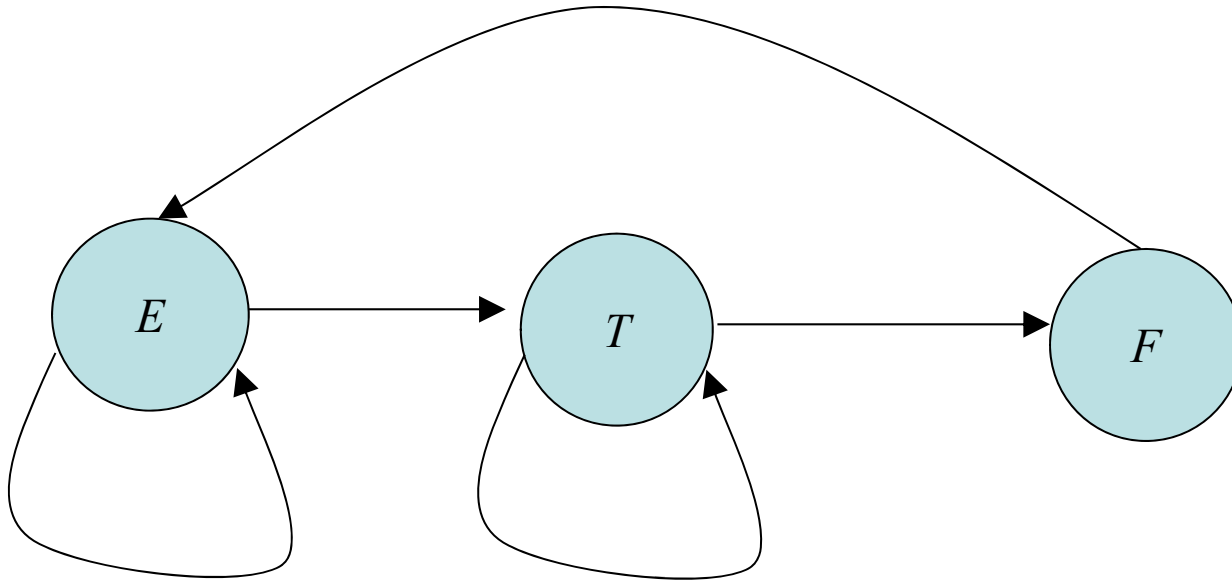
$T$  (term) has immediate left recursion and non immediate recursion

$G$  is clean, recursive, noncircular, hence the generated language is infinite

grammar has recursions

$\Leftrightarrow$

the graph of the produce relation has circuits



$$G = (\{E, T, F\}, \{i, +, *, ), ( \}, P, E)$$

$$P = \{E \rightarrow E + T \mid T, \quad T \rightarrow T * F \mid F, \quad F \rightarrow (E) \mid i\}$$

## SYNTAX TREES AND CANONICAL DERIVATIONS

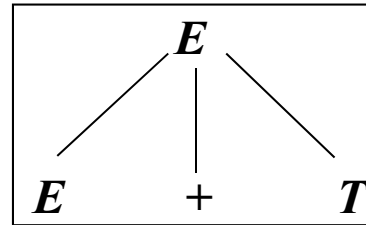
SYNTAX TREE: An oriented, sorted graph (children sorted from left to right) with no cycles, such that, for each pair of nodes, there is only one path connecting them

- it represents graphically the derivation process
- father-child relation / descendants / root node / leaf (or terminal) nodes
- degree of a node: number of its children
- root contains the axiom  $S$
- frontier of the tree (leaf sequence from left to right) contains the generated phrase

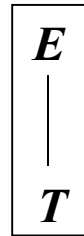
SUBTREE with root  $N$ : the tree having  $N$  as its root; it includes  $N$  *and* all its descendants

Expressions with sums and products:  $E$  expression,  $T$  term,  $F$  factor

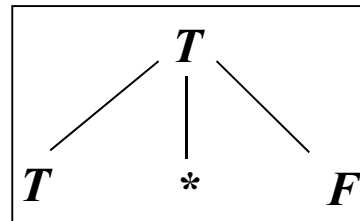
$$1. E \rightarrow E + T$$



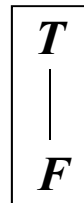
$$2. E \rightarrow T$$



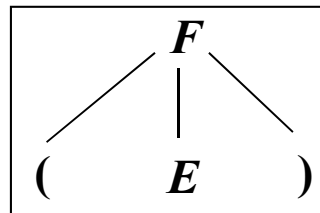
$$3. T \rightarrow T * F$$



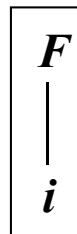
$$4. T \rightarrow F$$



$$5. F \rightarrow ( E )$$



$$6. F \rightarrow i$$



every rule application

generates a tree fragment

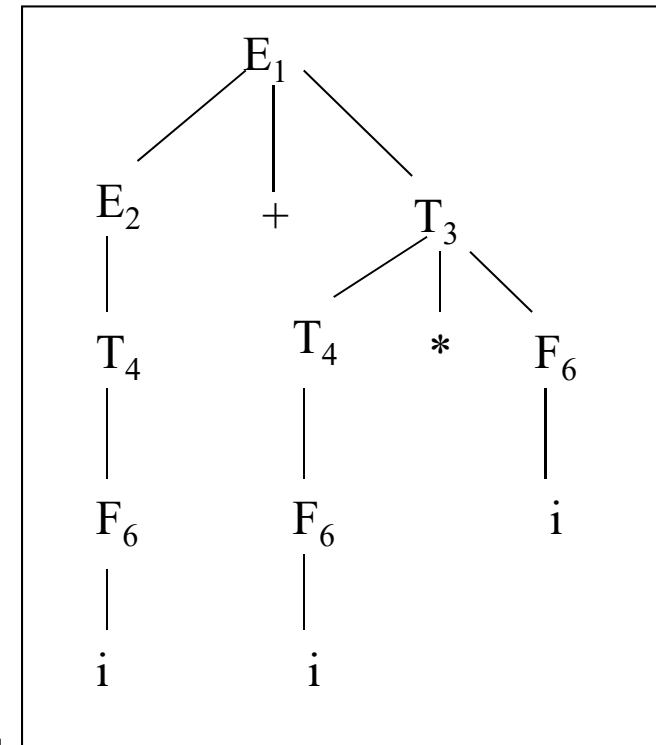
father + its children

Syntax tree for sentence  $i + i * i$

(n.t. labelled with the applied rule)

Linear representation of the tree (subscripts indicate the n.t. in the subtree root)

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



**LEFT DERIVATION**

(numbers denote the rule, expanded nonterm. is underlined)

$$\begin{aligned}
 E &\Rightarrow \underline{E} + T \Rightarrow \underline{T} + T \Rightarrow \underline{F} + T \Rightarrow i + \underline{T} \Rightarrow i + \underline{T} * F \Rightarrow \\
 &\Rightarrow i + \underline{F} * F \Rightarrow i + i * \underline{F} \Rightarrow i + i * i
 \end{aligned}$$

**RIGHT DERIVATION** (numbers denote the rule, expanded nonterm. is underlined)

$$\begin{aligned}
 E &\Rightarrow E + \underline{T} \Rightarrow E + T * \underline{F} \Rightarrow E + \underline{T} * i \Rightarrow E + \underline{F} * i \Rightarrow \underline{E} + i * i \Rightarrow \\
 &\Rightarrow \underline{T} + i * i \Rightarrow \underline{F} + i * i \Rightarrow i + i * i
 \end{aligned}$$

Derivations may be neither right nor left

$$\begin{array}{l}
 E \Rightarrow_{l,r} E + \underline{T} \Rightarrow_r \underline{E} + T * F \Rightarrow_l T + \underline{T} * F \Rightarrow T + F * \underline{F} \Rightarrow_r \underline{T} + F * i \Rightarrow_l \\
 \Rightarrow_l F + \underline{F} * i \Rightarrow_r \underline{F} + i * i \Rightarrow_{l,r} i + i * i
 \end{array}$$

However, for a **fixed** syntax tree of a sentence, there exist

a unique right derivation, and

a unique left derivation

matching that tree

Right and left derivation are useful to define *parsing* (i.e., syntax analysis) algorithms

A different question: does a given sentence have a unique syntax tree?

This determines the *ambiguity* of the grammar



$$G = \left( \underbrace{\{E, T, F\}}_{\text{non term.}}, \underbrace{\{i, +, *, ), (\}}_{\text{term.}}, \underbrace{P}_{\text{productions}}, \underbrace{E}_{\text{axiom}} \right)$$

$$P = \{E \rightarrow E + T \mid T, \quad T \rightarrow T * F \mid F, \quad F \rightarrow (E) \mid i\}$$

$$L(G) = \{i, i + i + i, \quad i * i, \quad (i + i) * i, \quad \dots\}$$

(this rule structure is necessary to model precedence, associativity, and avoid ambiguity)

IMPORTANT: as an exercise, generate many strings, to understand how the grammar works, and check that operator precedence is necessarily respected

### VERY IMPORTANT

«Rule of thumb» for modeling mutual precedence of operators:

- terminals for **low-precedence** operators derived **first** Ex: ‘+’
- terminals for **high-precedence** operators derived **later** Ex: ‘\*’

i.e., they are respectively closer to or farther from the axiom in the *produce* relation

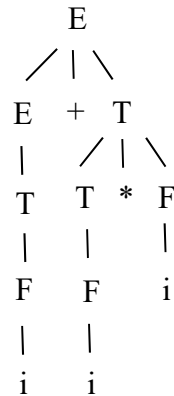
«Rule of thumb» for modeling associativity of operators:

- left-recursive rule  $\Rightarrow$  left-associative operator Ex: ‘+’ and ‘\*’ are left assoc.
- right-recursive rule  $\Rightarrow$  right-associative operator

With  $P = \{ E \rightarrow E + T \mid T, \quad T \rightarrow T * F \mid F, \quad F \rightarrow (E) \mid i \}$

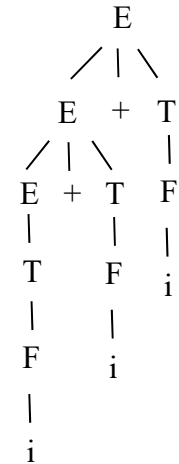
$i + i * i$  has syntax tree

‘\*’ higher precedence than ‘+’



$i + i + i$  has syntax tree

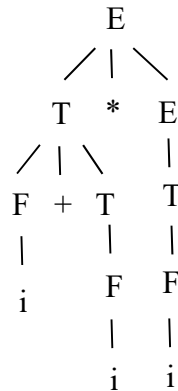
‘+’ left associative



With  $P = \{ E \rightarrow T * E \mid T, \quad T \rightarrow F + T \mid F, \quad F \rightarrow (E) \mid i \}$

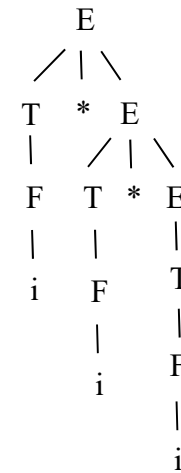
$i + i * i$  has syntax tree

‘+’ higher precedence than ‘\*’



$i * i * i$  has syntax tree

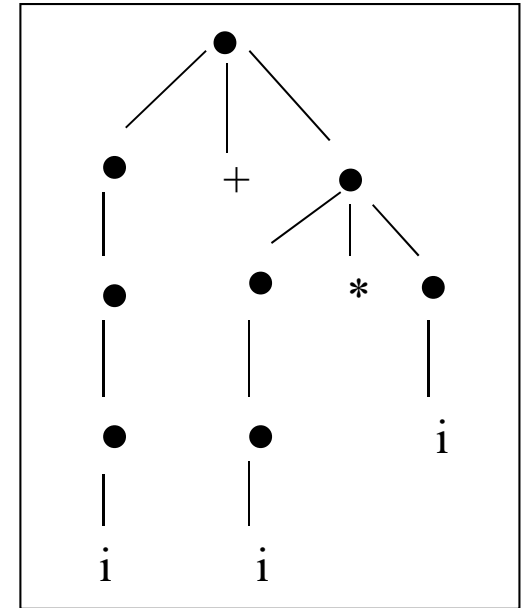
‘\*’ right associative



## SKELETON TREE (only the frontier and the structure)

linear representation

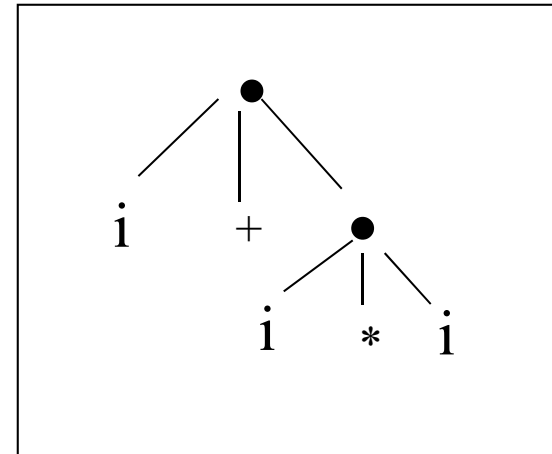
$$[[[[i]]] + [[[i]] * [i]]]$$



## CONDENSED SKELETON TREE (internal nodes on non-branching paths are merged)

linear representation

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



## STRONG (STRUCTURAL) AND WEAK EQUIVALENCE

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

$G$  and  $G'$  might assign different structures (syntax trees) to the same sentence

**the structure assigned to a sentence is important:**

**it is used by translators and interpreters**

STRONG or STRUCTURAL EQUIVALENCE of two grammars  $G$  and  $G'$

$L(G) = L(G')$  (weak eq.) **and**

$G$  and  $G'$  have the same *condensed skeleton trees*

strong eq.  $\rightarrow$  weak eq. **but** strong eq.  $\neq$  weak eq.

(hence  $\neg(\text{weak eq.} \rightarrow \text{strong eq.})$ )

strong eq. is **decidable**

weak eq. is **not decidable**

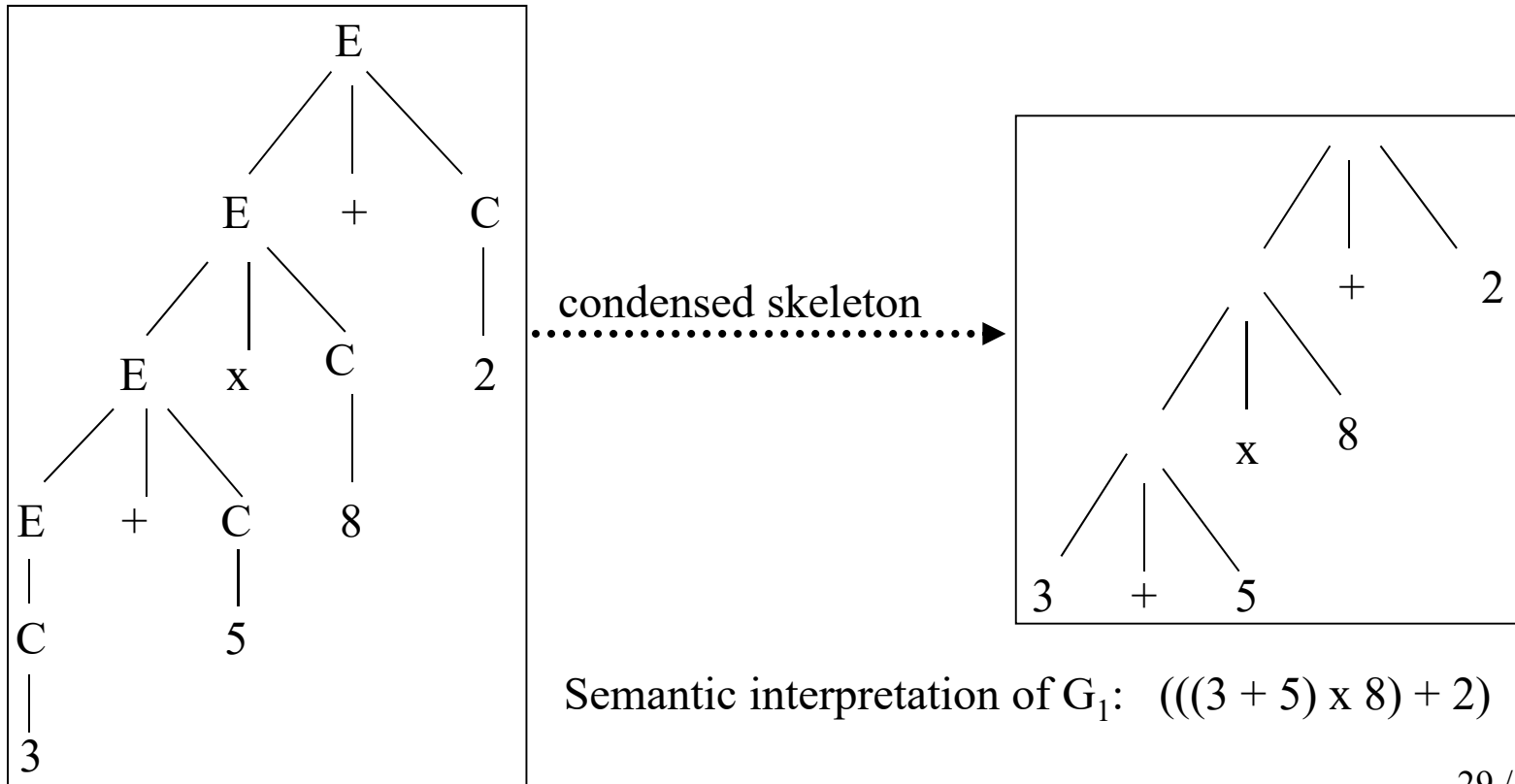
(it can happen that one can establish that  $G_1$  and  $G_2$  **are not** strongly equivalent

but is unable to establish whether they are weakly equivalent or not)

Example: Structural equivalence of arithmetic expressions:  $3 + 5 \times 8 + 2$

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

NB: the grammar has *left recursion* hence the operators are *left-associative*



NB: copy (or categorization) rules

```

graph TD
    Root(( )) --- L1_1(( ))
    Root --- L1_2[+]
    Root --- L1_3(( ))
    L1_1 --- L2_1[3]
    L1_1 --- L2_2(( ))
    L2_2 --- L3_1[+]
    L2_2 --- L3_2(( ))
    L3_2 --- L4_1[5]
    L3_2 --- L4_2[x]
    L1_3 --- L2_3[+]
    L1_3 --- L2_4(( ))
    L2_4 --- L3_3[2]
    L2_4 --- L3_4[8]
  
```

```

graph TD
    E1[E] --- E2[E]
    E1 --- P1[+]
    E1 --- T1[T]
    E2 --- E3[E]
    E2 --- P2[+]
    E2 --- T2[T]
    E3 --- T3[T]
    T3 --- C1[C]
    C1 --- 3[3]
    T2 --- C2[C]
    C2 --- 5[5]
    T2 --- X[x]
    T2 --- C3[C]
    C3 --- 8[8]
    T1 --- C4[C]
    C4 --- 2[2]
  
```

Semantic interpretations:  $G_1: (((3 + 5) \times 8) + 2)$   
 $G_2: ((3 + (5 \times 8)) + 2)$

**NB: this is how operator priority can be enforced through grammars**

# PARENTHESIS LANGUAGES

structures with pairs of opening / closing marks

***nested***: inside a pair there can be other parenthesized structures (recursion)

(nested) structures can also be placed in sequences at the same level of nesting

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

Abstracting away from the type of parentheses, the paradigmatic language is

## *The Dyck language*

Ex. alphabet:

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

Sentence example:

$$O[[O[]]O]$$

**DYCK LANGUAGE** with opening parenthesis  $a, \dots$ , closing parenthesis  $c, \dots$   
 Grammar is surprisingly simple

$$\begin{array}{l} \Sigma = \{a, c\} \\ S \rightarrow aScS \mid \varepsilon \\ \\ a \ a \ \underbrace{ac} \ c \ a \ a \ \underbrace{ac} \ c \ c \ c \\ \underbrace{\hspace{1.5cm}} \quad \underbrace{\hspace{1.5cm}} \\ \underbrace{\hspace{4cm}} \end{array}$$

The language is not linear ( $>1$  nonterm. in the right part)

Exercise: build the syntax tree for the sentence

$a \ a \ a \ c \ c \ a \ a \ c \ c \ c \ c$

## LINEAR NON REGULAR LANGUAGE:

$$\begin{array}{l} L_1 = \{a^n c^n \mid n \geq 1\} = \{ac, aacc, \dots\} \\ S \rightarrow aSc \mid ac \end{array}$$

$L_1$  is a proper subset of the Dyck Language:  
 it does not admit many nested structures at the same level  
 (that is, strings of type  $acacac$  are ruled out)



# REGULAR COMPOSITION OF FREE LANGUAGES

Applying the union, concatenation, Kleen star operations to free languages ...  
one obtains free languages, therefore ...

*The family of free languages is closed under union, concatenation, and star*

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

$$V_{N_1} \cap V_{N_2} = \emptyset \quad S \notin (V_{N_1} \cup V_{N_2})$$

NB: one needs disjoint nonterm. sets and a new axiom

UNION:

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

CONCATENATION:

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

STAR:  $G$  for  $(L_1)^*$  is obtained by adding to  $G_1$  the rules  $S \rightarrow SS_1 \mid \varepsilon$

CROSS '+' (not necessary because '+' is derived; this is added for simplicity):

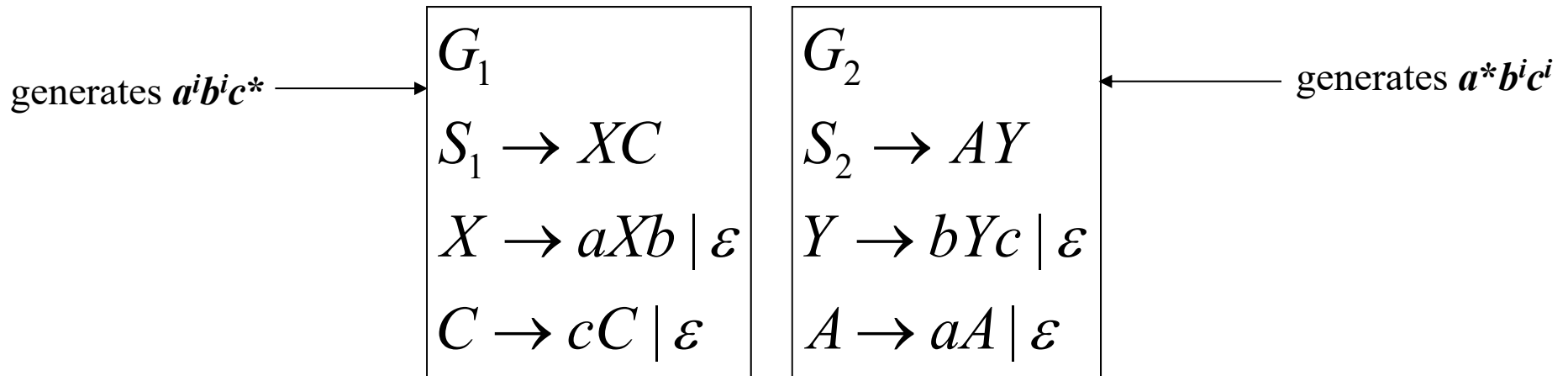
$G$  for  $(L_1)^+$  is obtained by adding to  $G_1$  the rules  $S \rightarrow SS_1 \mid S_1$

The **MIRROR LANGUAGE** of  $L(G)$ ,  $(L(G))^R$  is generated by the **mirror grammar**,  
obtained by reversing the right part of the rules

EXAMPLE: Union of free languages

$$L = \{a^i b^j c^k \mid i = j \vee j = k\} = \{a^i b^i c^* \mid i \geq 0\} \cup \{a^* b^i c^i \mid i \geq 0\} = L_1 \cup L_2$$

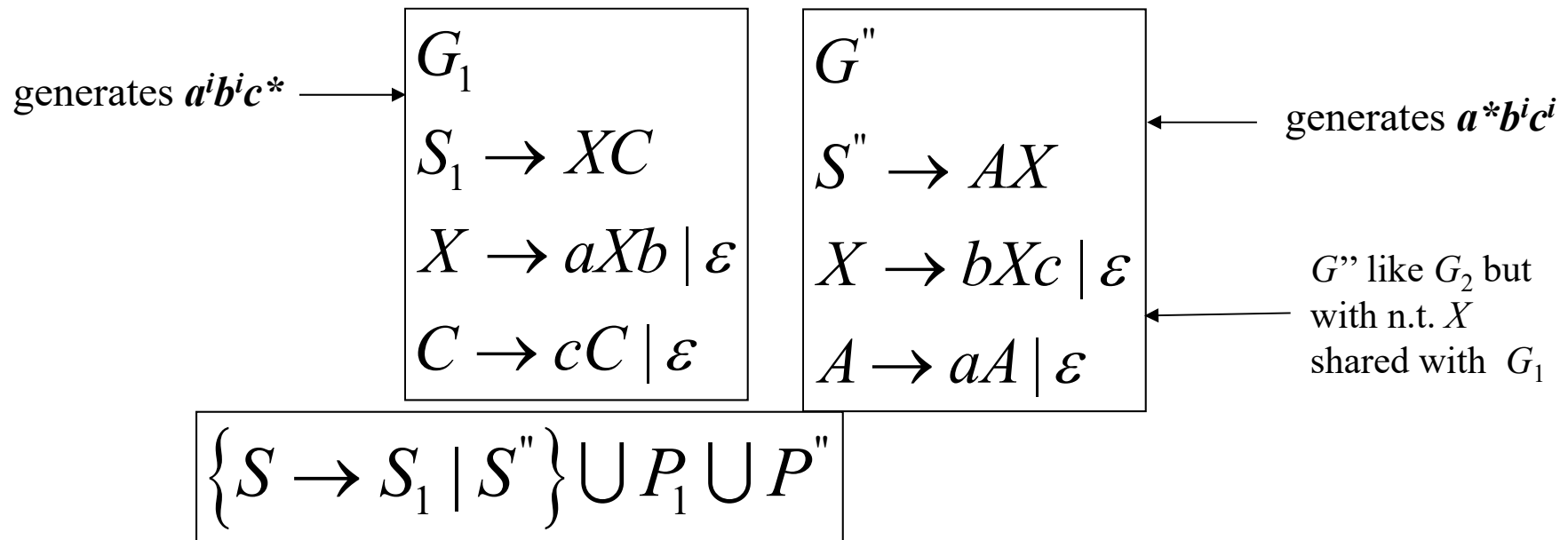
$L_1$  and  $L_2$  are generated by the two grammars  $G_1$  and  $G_2$



Notice that the nonterminal sets of grammars  $G_1$  and  $G_2$  are DISJOINTED

EXAMPLE: Union of free languages (follows)

What happens if the nonterm. sets are not disjoint? Let us use  $G''$  instead of  $G_2$



If nonterm. are not disjoint, the grammar generates a **superset** of the union language:  
spurious additional sentences are generated

