

Context free grammars - II

Prof. A. Morzenti

SYNTAX TREES AND CANONICAL DERIVATIONS

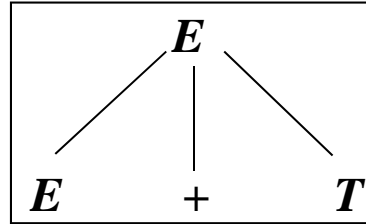
SYNTAX TREE: An oriented, sorted graph (children sorted from left to right) with no cycles, such that, for each pair of nodes, there is only one path connecting them

- 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

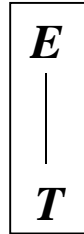
SUBTREE with root N : the tree having N as its root; it includes N *and* all its descendants

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

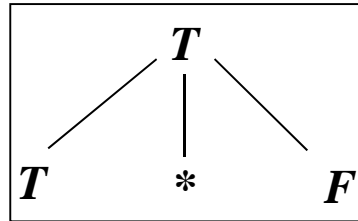
$$1. E \rightarrow E + T$$



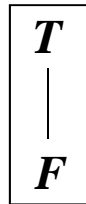
$$2. E \rightarrow T$$



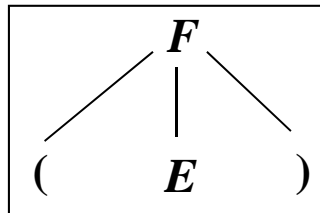
$$3. T \rightarrow T * F$$



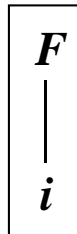
$$4. T \rightarrow F$$



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



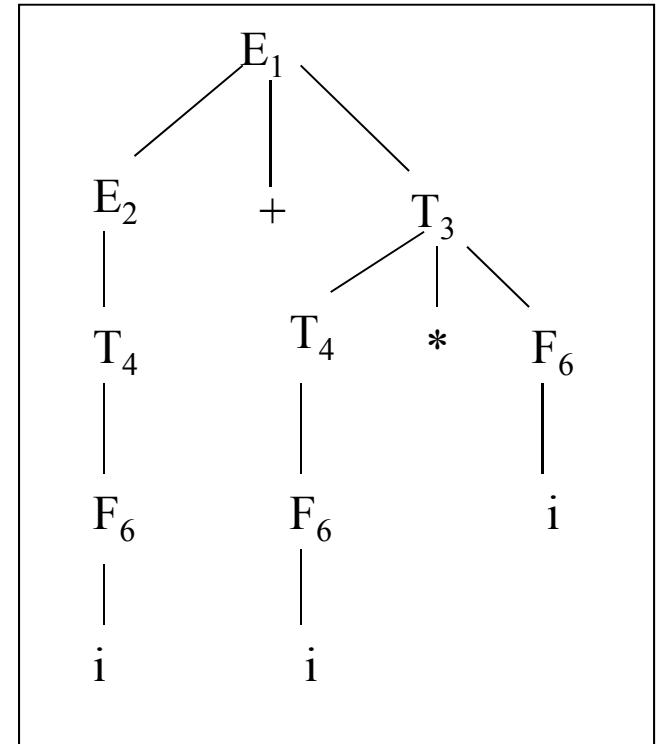
$$6. F \rightarrow i$$



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 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

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

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

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

(apparently strange structure but necessary to model operator precedence, and avoid ambiguity)

Exercise: specify the lang. of arithmetic expr. using regular expressions

obviously not possible ;-)

recursion makes expr. similar to lists with unbounded nesting level

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

«Rule of thumb» for modeling precedence:

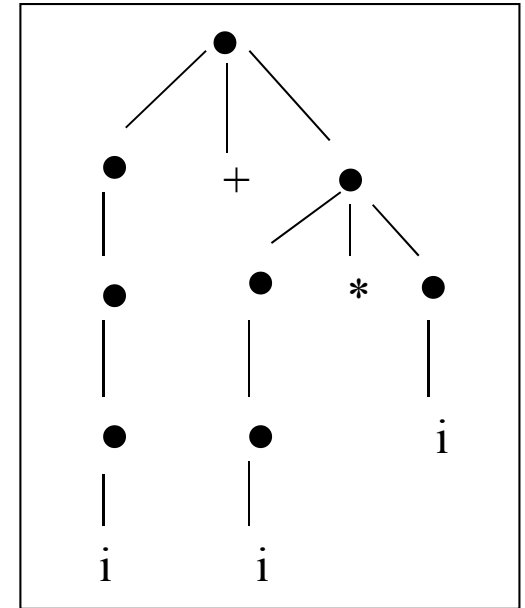
- nonterminals for **low-precedence** operators derived **first**

- nonterminals for **high-precedence** operators derived **later**

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

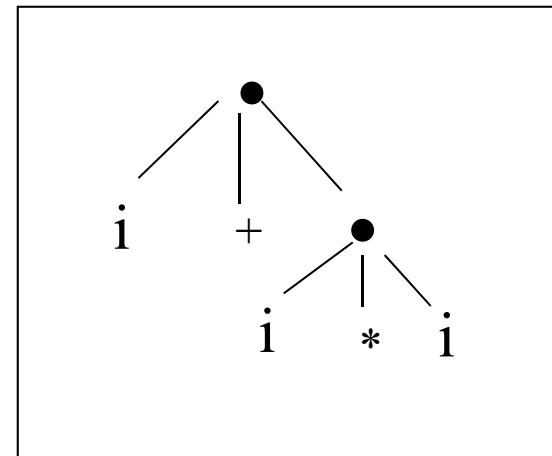
SKELETON TREE (only the frontier and the structure)

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



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

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



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

$$\begin{array}{l} \Sigma = \{a, c\} \\ S \rightarrow aScS \mid \varepsilon \\ \\ a \ a \ \underbrace{ac} \ c \ a \ a \ \underbrace{ac} \ c \ c \ c \\ \underbrace{\hspace{1.5cm}} \quad \underbrace{\hspace{1.5cm}} \\ \underbrace{\hspace{3.5cm}} \end{array}$$

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

Exercise: build the syntax tree for the sentence

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

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:
 is 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, that is ...

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 $S \rightarrow SS_I \mid \varepsilon$

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

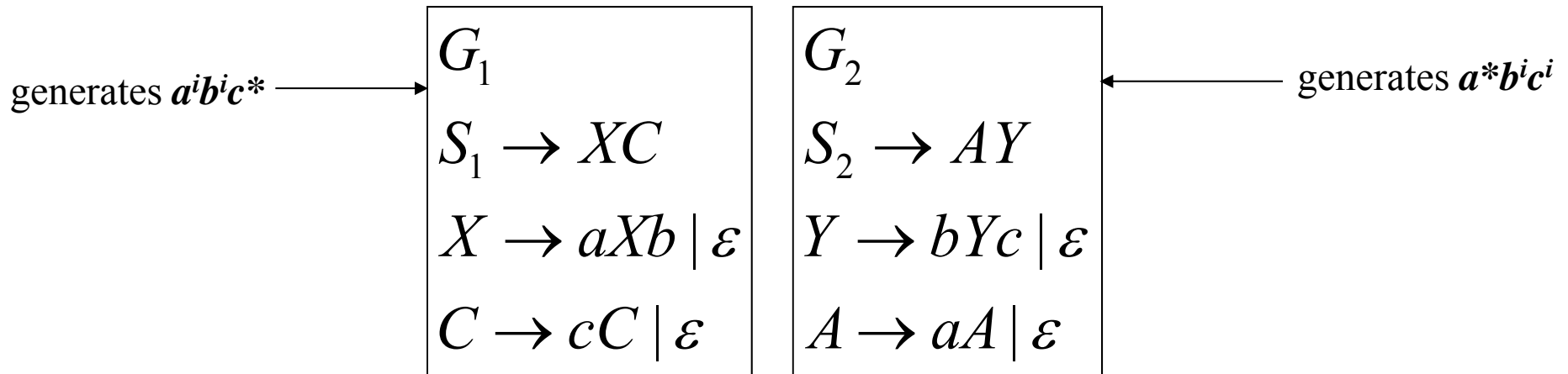
G for $(L_1)^+$ is obtained by adding to G_1 the rules $S \rightarrow SS_I \mid S_I$

BTW: 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

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

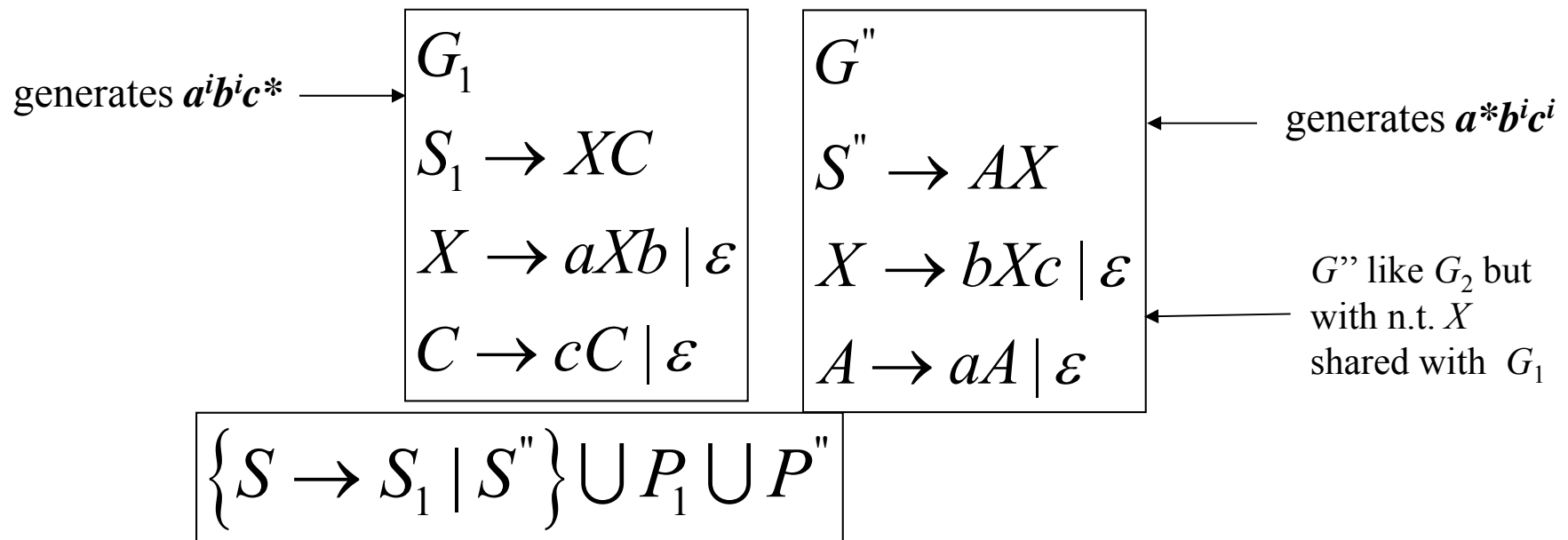
L_1 and L_2 are generated by the two grammars G_1 and G_2



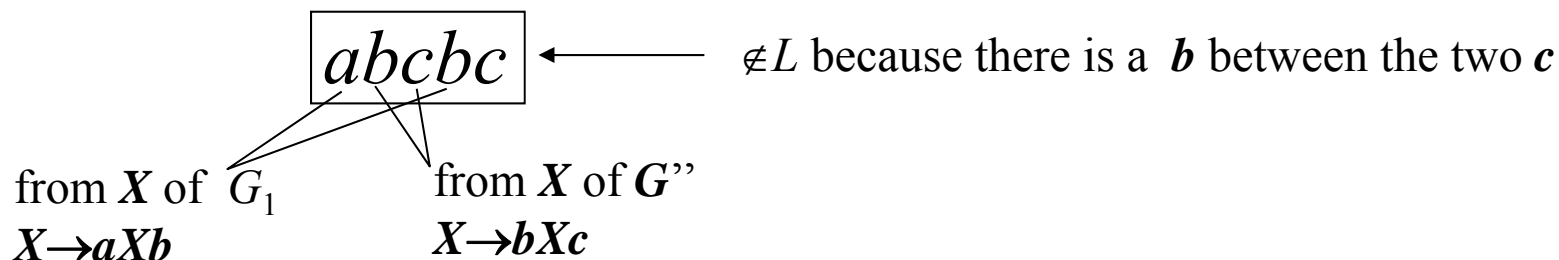
Notice that the nonterminal sets of grammars G_1 and G_2 are DISJOINTED

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



AMBIGUITY

Examples from natural language

“I saw the man with the binoculars.”

“La pesca è bella”

“half baked chicken”

“Questa è una rapina, dateci i soldi altrimenti *spariamo*.” “OK, allora *sparite*”...

In artificial, technical languages ambiguity must be ruled out (in the natural ones...).

We only consider SYNTACTIC AMBIGUITY:

A sentence x of a grammar G is ambiguous if it admits several distinct syntax trees

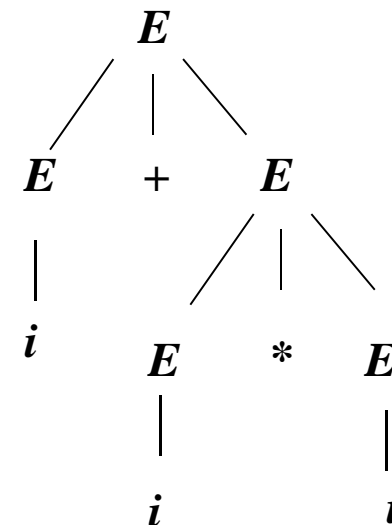
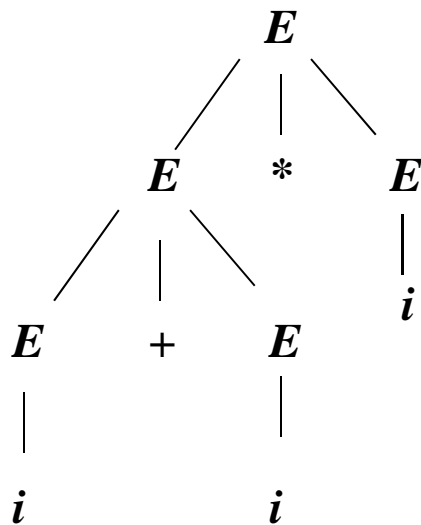
In such a case we say that the grammar G is ambiguous

EXAMPLE: a “simple” grammar G' for 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 sentence $i+i+i$ is also ambiguous

G' is ambiguous and **does not enforce** the usual **precedence of product over sum**

NB: G' is smaller than the previous G (p.3): it has only 1 n.t. and 4 rules
in grammar design there is oft a **trade-off** between size and ambiguity

The *degree of ambiguity of a sentence* x of a language $L(G)$
is the number of distinct trees of x compatible with G
for a *grammar* the degree of ambiguity
is the maximum among the degree of ambiguity of its sentences
Notice that the degree of ambiguity of sentences of a grammar can be unlimited

IMPORTANT PROBLEM: determine if a grammar is ambiguous

The *problem* is *undecidable*:
there does not exist a general algorithm that, given any free grammar,
terminates (after a finite number of steps) with the correct answer

The *absence of ambiguity* in a specific grammar can be shown
on a case by case basis, by hand, through *inductive reasonings*,
hence by analyzing a finite number of cases

Instead, to show ambiguity one can exhibit a *witness*: an ambiguous sentence

| |
|---|
| BEST APPROACH: AVOID AMBIGUITY IN THE GRAMMAR DESIGN PHASE |
|---|

Let us reconsider the arithmetic expression example

The degree of ambiguity is 2 for $i + i + i$ and 5 for $i + i * i + i$

| | | | | |
|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|
| $i + i * i + i$ | $i + i * i + i$ | $i + i * i + i$ | $i + i * i + i$ | $i + i * i + i$ |
| $\underbrace{\quad\quad\quad}$ | $\underbrace{\quad\quad\quad}$ | $\underbrace{\quad\quad\quad}$ | $\underbrace{\quad\quad\quad}$ | $\underbrace{\quad\quad\quad}$ |
| | $\underbrace{\quad\quad\quad}$ | $\underbrace{\quad\quad\quad}$ | $\underbrace{\quad\quad\quad}$ | $\underbrace{\quad\quad\quad}$ |

For longer sentences the degree of ambiguity increases with no limit

CATALOG OF AMBIGUOUS FORMS AND REMEDIES

1) AMBIGUITY FROM BILATERAL RECURSION: $A \rightarrow A \dots A$

Example 1: grammar G_1 generates $i + i + i$ in two different ways (check!).

$$G_1 : E \rightarrow E + E \mid i$$

But $L(G_1) = i (+i)^*$ is a regular language: other, simpler grammars are possible

Nonambiguous right-recursive grammar : $E \rightarrow i + E \mid i$

Nonambiguous left-recursive grammar : $E \rightarrow E + i \mid i$

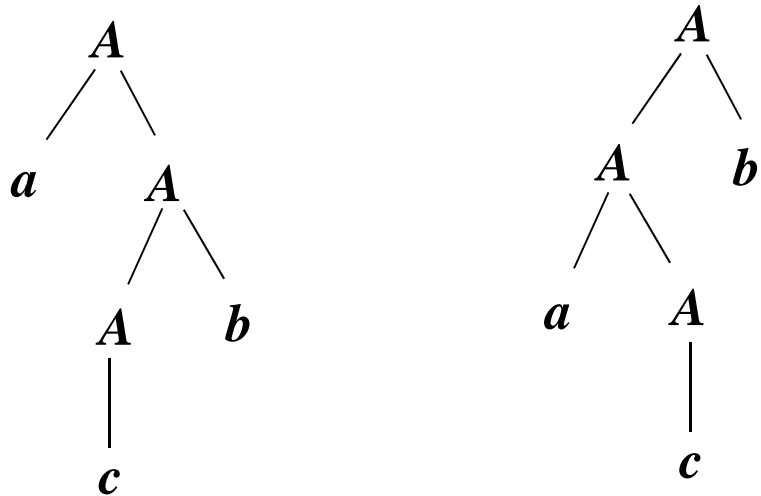
Example 2: Left and right recursion in different rules

$$L(G_2) = a^*cb^*$$

G_2 admits derivations where the a and b in a given sentence are obtained in any order

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

Simplest example: two syntax trees for sentence acb



Example 2 (follows): Left and right recursion in different rules

First method to eliminate ambiguity:
Generate the two lists by distinct rules

$$L(G_2) = a^*cb^*$$

$$S \rightarrow AcB$$

$$A \rightarrow aA \mid \varepsilon$$

$$B \rightarrow bB \mid \varepsilon$$

Second method: force an *order in the derivations*:
ex., first generate the a 's then the b 's

$$L(G_2) = a^*cb^*$$

$$S \rightarrow aS \mid X$$

$$X \rightarrow Xb \mid c$$

2) AMBIGUITY FROM LANGUAGE UNION

If $L_1=L(G_1)$ and $L_2=L(G_2)$ share some sentences (nonempty intersection)

The grammar G for the union language (slide n 10), is ambiguous (slide n.11)

A sentence $x \in L_1 \cap L_2$ admits two distinct derivations,

One with the rules of G_1 and the other with the rules of G_2

It is ambiguous for a grammar G that includes all the rules

Instead the sentences that belong to $L_1 \setminus L_2$ and to $L_2 \setminus L_1$ are nonambiguous

Remedy: provide disjointed set of rules for $L_1 \cap L_2$, $L_1 \setminus L_2$, and $L_2 \setminus L_1$

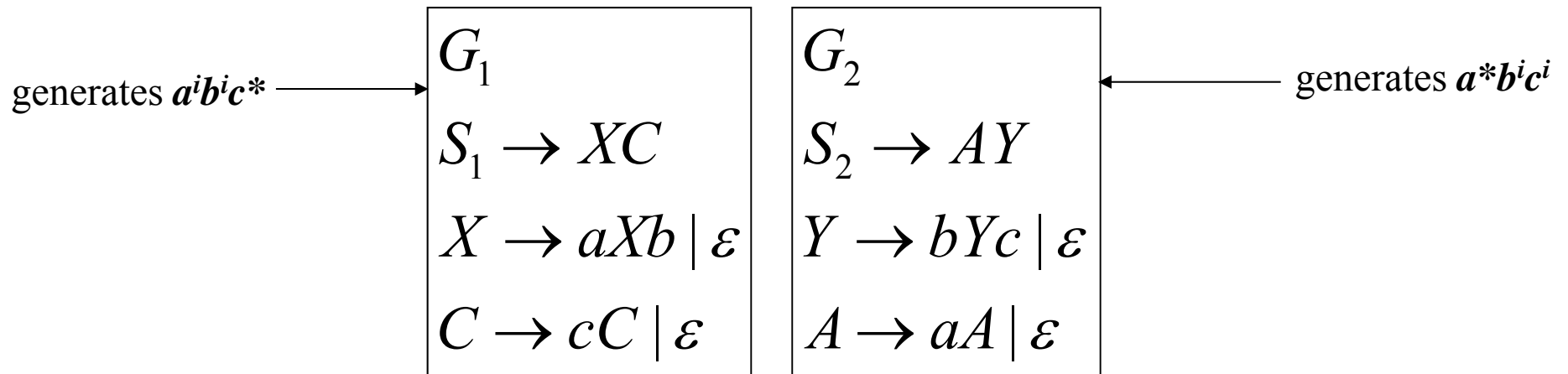
Examples in the textbook and in exam exercises

3) INHERENT AMBIGUITY

A language is INHERENTLY AMBIGUOUS if all its grammars are ambiguous

EXAMPLE:

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



The union grammar is ambiguous for the sentences $\varepsilon, abc, a^2 b^2 c^2, \dots$
(which constitute a non-context free language)

They are generated by G_1 which has rules ensuring that $|x|_a = |x|_b \dots$

$$\begin{array}{c} a \dots a \underbrace{ab} b \dots b \underbrace{bcc} \dots c \\ \underbrace{\hspace{10em}} \end{array}$$

\dots and also by G_2 , with rules ensuring that $|x|_b = |x|_c$

$$\begin{array}{c} a \dots a \underbrace{ab} b \dots b \underbrace{bc} c \dots c \\ \underbrace{\hspace{10em}} \end{array}$$

We *intuitively* argue that *any* grammar for L is ambiguous

Luckily inherent ambiguity is rare and can be avoided in technical languages

4) AMBIGUITY FROM CONCATENATION OF LANGUAGES

when the suffix of some sentence in the first language
is also a prefix of a sentence in the second one

G for the
concatenation of
 L_1 and L_2

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

Supposing G_1 and G_2 are not ambiguous, G is ambiguous if $\exists x' \in L_1$ and $x'' \in L_2$
such that there is a non empty string v such that

$$\begin{aligned} x' &= u'v \wedge u' \in L_1 & x'' &= vz'' \wedge z'' \in L_2 \\ u'vz'' &\in L_1.L_2 & \text{and is ambiguous:} \\ S &\Rightarrow S_1 S_2 \xRightarrow{+} u' S_2 \xRightarrow{+} u' vz'' \\ S &\Rightarrow S_1 S_2 \xRightarrow{+} u' v S_2 \xRightarrow{+} u' vz'' \end{aligned}$$

Example – Ambiguity in the concatenation of Dyck languages

$$\begin{aligned}\Sigma_1 &= \{a, a', b, b'\} & \Sigma_2 &= \{b, b', c, c'\} \\ aa'bb'cc' &\in L = L_1L_2 \\ G(L): \quad S &\rightarrow S_1S_2 \\ S_1 &\rightarrow aS_1a'S_1 \mid bS_1b'S_1 \mid \varepsilon \\ S_2 &\rightarrow bS_2b'S_2 \mid cS_2c'S_2 \mid \varepsilon \\ \underbrace{aa'}_{S_1} \underbrace{bb'cc'}_{S_2} & \quad \underbrace{aa'}_{S_1} \underbrace{bb'cc'}_{S_2}\end{aligned}$$

To prevent ambiguity one must avoid the shift of the substring from the suffix in the first language to the prefix in the second one

Easy solution: insert a **new** terminal symb. (e.g. ‘#’) acting as a separator

new: the symbol may not belong to any of the two alphabets

$L_1 \# L_2$ is generated by the «concatenation grammar» with the axiomatic rule $S \rightarrow S_1 \# S_2$.

The language is however modified

5) OTHER CASES OF AMBIGUITY

AMBIGUOUS REGULAR EXPRESSIONS

every sentence with two or more
 c is ambiguous

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

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

REMEDY: identify the leftmost c

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

OTHER CASE: LACK OF ORDER IN DERIVATIONS

Example: Every phrase with
>2 b 's is ambiguous

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

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

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

REMEDY: Impose an order in the derivation
First the rule that generates balanced b 's and c 's
Then the one that generates excess b 's

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

6) AMBIGUITY IN CONDITIONAL PHRASES: the (in)famous *dangling else* problem

$$S \rightarrow \underbrace{\text{if } b \text{ then } S \text{ else } S}_{\text{if } b \text{ then } \underbrace{\text{if } b \text{ then } a \text{ else } a}} \mid \text{if } b \text{ then } S \mid a$$

$\underbrace{\text{if } b \text{ then } \text{if } b \text{ then } a \text{ else } a}_{\text{if } b \text{ then } \text{if } b \text{ then } a \text{ else } a} \quad \leftarrow$

$$\begin{aligned} S &\rightarrow S_E \mid S_T \\ S_E &\rightarrow \text{if } b \text{ then } S_E \text{ else } S_E \mid a \\ S_T &\rightarrow \text{if } b \text{ then } S_E \text{ else } S_T \mid \text{if } b \text{ then } S \end{aligned}$$

REMEDY 1 Rule out this interpretation (i.e., use the *closest match interpretation*);

two n.t., S_E and S_T :

S_E always has the *else* (cannot have only the *then*), only for S_T it is possible not to have the *else* therefore only S_E can precede the *else*

$$\begin{aligned} S &\rightarrow \text{if } b \text{ then } S \text{ else } S \text{ endif} \\ &\quad \mid \text{if } b \text{ then } S \text{ endif} \mid a \end{aligned}$$

REMEDY 2 Modify the language, introduce a closing mark *endif*