

# Context-free Grammars and Languages

Context-free languages

$\Sigma = \{a, +, *, (\}, \}$

$E \rightarrow P + E$   
 $E \rightarrow P * E$   
 $E \rightarrow P$   
 $P \rightarrow (E)$   
 $P \rightarrow a$

Example strings in  $L(G)$  [  $a$ ,  $(a)$  ]

A CFG  $G = (V, \Sigma, R, S)$   
 $V$  a finite set ("variables")  
 $V \cap \Sigma = \emptyset$   
 ("start" or "sentence")  $S \in V$   
 ("rules")  $R$  is a finite subset of  $V \times (V \cup \Sigma)^*$

A derivation tree or parse tree in  $G$

Another string in  $L(G)$

1

2

**Notation**  $\rightarrow$  in rules (only; "produces" or "may be rewritten as")

$\Rightarrow$  "yields":  
 relation on strings in  $(V \cup \Sigma)^*$   
 $\alpha A \beta \Rightarrow \alpha \gamma \beta$   
 if  $A \rightarrow \gamma$  is a rule  
 for all  $\alpha, \beta \in (V \cup \Sigma)^*$   
 i.e., "context-free"

$\Rightarrow^*$  "derives" (reflexive, transitive closure of  $\Rightarrow$ ; "0 or more steps")

$\alpha \Rightarrow^* \beta$  means  $\exists \alpha_0, \alpha_1, \dots, \alpha_k$   $k \geq 0$   
 $\alpha = \alpha_0 \Rightarrow \alpha_1 \Rightarrow \alpha_2 \dots \Rightarrow \alpha_k = \beta$

$L(G) = \{ w \in \Sigma^+ \mid S \Rightarrow^* w \}$

3

**Example**

$G = (V, \Sigma, R, S)$   
 $V = \{S\}$   
 $\Sigma = \{a, b\}$   
 $R$ :  
 $S \rightarrow aSb \mid \epsilon$   
 $S \Rightarrow \epsilon$   
 $S \Rightarrow aSb \Rightarrow ab$   
 $S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aabb$   
 $\vdots$   
 $L(G) = \{ a^n b^n \mid n \geq 0 \}$

Note that  $L(G)$  is non-regular

4

Example

$$G_2 = (V, \Sigma, R, S)$$

$$V = \{S\}$$

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

R:

$$S \rightarrow aSa \mid bSb \mid \epsilon$$

$$S \Rightarrow \epsilon$$

$$S \Rightarrow aSa \Rightarrow aa$$

$$S \Rightarrow bSb \Rightarrow bb$$

$$S \Rightarrow aSa \Rightarrow abSba \Rightarrow abba$$

$$L(G_2) = \text{even length palindromes} \\ \{w w^R \mid w \in \Sigma^*\}$$

Example

We'll see later that  $L_{two} = \{ww \mid w \in \Sigma^*\}$  is *not* context free. At first glance, you might think that adding a new start symbol  $S'$  and a rule  $S' \rightarrow SS$  to  $G_2$  would generate  $L_{two}$ , but it doesn't; it generates all strings in  $L_{two}$  *plus* many others, since derivations from the two  $S$ 's are *not* coordinated. (Why not? It's context-free; what happens to one  $S$  can't influence the other.)

Example

$$G_3: \text{as above but add} \\ S \rightarrow a \mid b$$

$$L(G_3) = \text{all palindromes} \\ \{w \in \Sigma^+ \mid w = w^R\}$$

# Trees, Derivations and Ambiguity

A grammar

$$E \rightarrow P + E \\ P \rightarrow P * E \\ P \rightarrow (E) \\ P \rightarrow a$$

A tree



3 derivations correspond to same tree (same rules being used in the same places, just written in different orders in the linear derivation)

- 1)  $E \Rightarrow P+E \Rightarrow a+E \Rightarrow a+P \Rightarrow a+a$
- 2)  $E \Rightarrow P+E \Rightarrow P+P \Rightarrow a+P \Rightarrow a+a$
- 3)  $E \Rightarrow P+E \Rightarrow P+P \Rightarrow P+a \Rightarrow a+a$

But only one *leftmost* derivation corresponds to it (and vice versa). (more in HW?)

Another grammar for the same language:

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

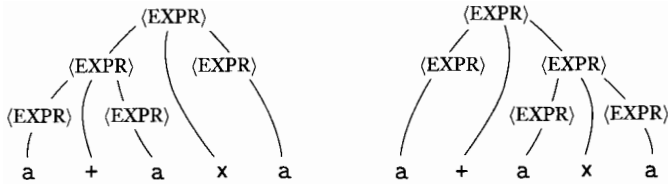


Fig 2.6: Two parse trees for a+a\*a in grammar G<sub>5</sub>

This grammar is *ambiguous*: there is a string in L(G<sub>5</sub>) with two different parse trees, or, equivalently, with 2 different leftmost derivations. Note the pragmatic difference: in general, (a+a)\*a != a+(a\*a); which is “right”?

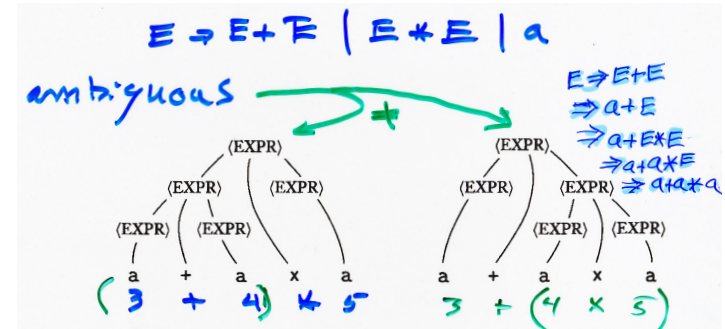
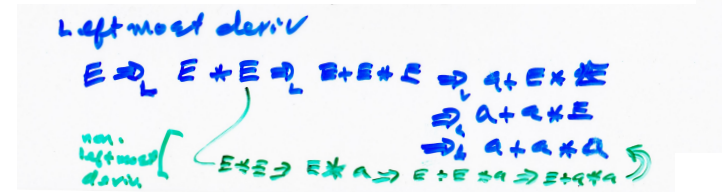


Fig 2.6: Two parse trees for a+a\*a in grammar G<sub>5</sub>

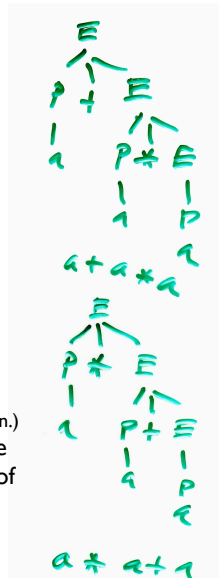


The “E, P” grammar again

$$\begin{aligned} E &\rightarrow P + E \\ P &\rightarrow P * E \\ P &\rightarrow (E) \\ P &\rightarrow a \end{aligned}$$

This grammar is *unambiguous*.

(Why? Very informally, the 3 E rules generate P((‘+’u\*)P)\* and only via a parse tree that “hangs to the right”, as shown.) But it has another undesirable feature: Parse tree structure does not reflect the usual precedence of \* over +. E.g., tree at lower right suggests “a \* a + a == a \* (a + a)”



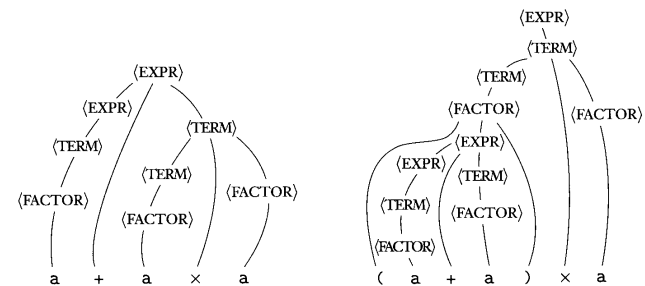
**EXAMPLE 2.4**

Consider grammar  $G_4 = (V, \Sigma, R, \langle \text{EXPR} \rangle)$ .

$V$  is  $\{\langle \text{EXPR} \rangle, \langle \text{TERM} \rangle, \langle \text{FACTOR} \rangle\}$  and  $\Sigma$  is  $\{a, +, x, (, )\}$ . The rules are

$$\begin{aligned} \langle \text{EXPR} \rangle &\rightarrow \langle \text{EXPR} \rangle + \langle \text{TERM} \rangle \mid \langle \text{TERM} \rangle \\ \langle \text{TERM} \rangle &\rightarrow \langle \text{TERM} \rangle x \langle \text{FACTOR} \rangle \mid \langle \text{FACTOR} \rangle \\ \langle \text{FACTOR} \rangle &\rightarrow ( \langle \text{EXPR} \rangle ) \mid a \end{aligned}$$

The two strings a+a\*a and (a+a)\*a can be generated with grammar  $G_4$ . The parse trees are shown in the following figure.



A more complex grammar, again the same language. This one is unambiguous and its parse trees reflect usual precedence/associativity of plus and times.

$$L = \{ a^i b^j c^k \mid i=j \text{ or } j=k \}$$

$$\left. \begin{aligned} S &\rightarrow AC \mid DB \\ A &\rightarrow aAb \mid \epsilon \\ C &\rightarrow cC \mid \epsilon \\ D &\rightarrow aD \mid \epsilon \\ B &\rightarrow bBc \mid \epsilon \end{aligned} \right\} G$$

$$\begin{aligned} a^1 b^1 c^2 \\ a^2 b^2 c^2 \\ a^3 b^3 c^3 \end{aligned}$$

Can we always tweak the grammar to make it unambiguous?

No! Language L is a CFL; grammar at left. Easy to see this G is ambiguous—strings of the form  $a^n b^n c^n$  can come from the  $i=j$  (AC) or  $j=k$  (DB) path. Hard to prove, but true, that every G for this L is also ambiguous. Intuitively, a grammar can only match a's & b's or b's & c's, not both. As a related point,  $\{ a^n b^n c^n \mid n > 0 \}$  is not CFL.

G is ambiguous

L is inherently ambiguous, meaning every G for L is ambiguous

## Some closure results for CFLs

13

14

### Theorem:

The set of context-free languages is closed under  $\cup$ ,  $\cdot$ , and  $*$

### Corollary:

All regular languages are CFLs

Proof Sketch:

Directly give simple CFLs for  $\emptyset$ ,  $\{\epsilon\}$ , and  $\{a\}$  for each  $a \in \Sigma$ . Combine them using the above theorem.

(Aside:

We'll later prove that CFLs are *not* closed under intersection or complementation.)

15

### Proof: Closure under Concatenation

$$G_i = (V_i, \Sigma, R_i, S_i)$$

be 2 CFGs

with  $V_1 \cap V_2 = \emptyset$

lets  $\delta = V_1 \cup V_2$

Build new grammar

$$G = (V, \Sigma, R, S)$$

$$V = V_1 \cup V_2 \cup \{S\}$$

$$R = R_1 \cup R_2 \cup \{S \rightarrow S_1 S_2\}$$

$\forall x \in L_1 \forall y \in L_2$

$$S_1 \xRightarrow{G_1} x \ \& \ S_2 \xRightarrow{G_2} y$$

$$\therefore S \xRightarrow{G} S_1 S_2 \xRightarrow{G} x S_2 \xRightarrow{G} x y$$

$$\therefore L_1 \cdot L_2 \subseteq L(G)$$

16

Suppose  $S \Rightarrow_G^* w$

\*  $S \Rightarrow_G S_1 S_2 \Rightarrow_G^* w$  Then, for some  $x, y \in \Sigma^*$

$$S_1 S_2 \Rightarrow_L^* x S_2$$

using only rules from  $G_1$

$$x S_2 \Rightarrow_L^* xy = w$$

using only  $G_2$  rules

$$\left. \begin{array}{l} S_1 \Rightarrow^* x \text{ in } G_1 \\ S_2 \Rightarrow^* y \text{ in } G_2 \end{array} \right\} **$$

$L(G) \subseteq L_1 \cdot L_2$

A key issue in this direction of the proof is that, since  $V_1 \cap V_2 = \emptyset$ , there is no “crosstalk” between the two sub-grammars: any derivation in  $G$  from  $S_1$  is also a derivation in  $G_1$ , and likewise  $S_2/G_2$ , so derivation (\*) above in  $G$  can be split into (\*\*) in  $G_1$  &  $G_2$ .