# Context-free Grammars and Languages

# Context-free languages

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

Rules:
$$E \to P + E$$
$$E \to P * E$$
$$E \to P$$
$$P \to ( E )$$
$$P \to a$$

Example strings in L(G):
$$a$$
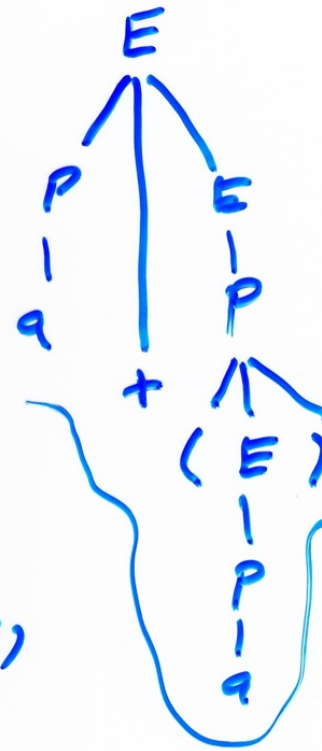$$( a )$$

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

$$a + ( a )$$

Another string in L(G)

2

$\longrightarrow$ in rules (only; "produces" or "may be rewritten as")

$\Longrightarrow$ "yeilds" :

relation on strings in $(\cup \cup \Sigma)^*$

$$\alpha A \beta \Longrightarrow \alpha \gamma \beta$$

if $A \to \gamma$ is a rule

for <u>all</u> $\alpha, \beta \in (\cup \cup \Sigma)^*$

$\longrightarrow$ i.e., "context-free"

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

$$\alpha \Rightarrow^* \beta \text{ means } \exists \; \alpha_0 \; \alpha_1, \cdots \alpha_k \quad k \geq 0$$

$$\alpha = \alpha_0 \Longrightarrow \alpha_1 \Longrightarrow \alpha_3 \cdots \alpha_k = \beta$$

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

3

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

$$V = \{S\}$$

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

$R$:

$$S \to aSb \mid \varepsilon$$

$$S \Rightarrow \varepsilon$$

$$S \Rightarrow aSb \Rightarrow ab$$

$$S \Rightarrow aSb \Rightarrow aaSbb \to aabb$$

$$\vdots$$

$$L(G) = \{a^n b^n \mid n \geq 0\}$$

Note that L(G) is non-regular

4

$G_2 = (V, \bar{\Sigma}, R, S)$

$V = \{S\}$

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

R :

$S \to aSa \mid bSb \mid \varepsilon$

$S \Rightarrow \varepsilon$

$S \Rightarrow a S a \Rightarrow a a$

$S \Rightarrow b S b \Rightarrow b b$

$S \Rightarrow a S a \Rightarrow a b S b a \Rightarrow a b b a$

$L(G_2) = $ even length palindromes

$\{w w^R \mid w \in \Sigma^* \}$

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
$$S' \to SS$$
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.)

5

$G_3$: as above but add
$$S \to a \mid b$$

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

# Trees, Derivations and Ambiguity

## A grammar

$E \rightarrow P + E$

$E \rightarrow P * E$

$E \rightarrow P$

$P \rightarrow ( E )$

$P \rightarrow a$

## A tree

```
      E
    / | \
   P  +  E
   |     |
   a     P
         |
         a
```

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 => P+E => a+E => a+P => a+a

2) E => P+E => P+P => a+P => a+a

3) E => P+E => P+P => P+a => a+a

But only one *leftmost* derivation corresponds to it

(and *vice versa*).  (more in HW?)

Another grammar for the same language:

E → E+E | E*E | (E) | a
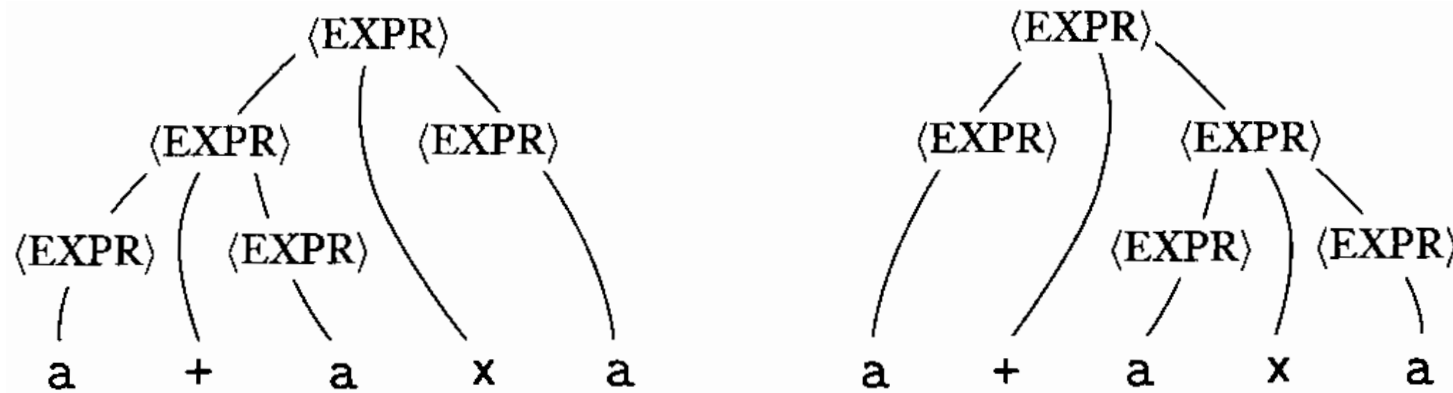


Fig 2.6: Two parse trees for a+a×a in grammar G₅

This grammar is *ambiguous*: there is a string in L(G₅) with two *different* parse trees, or, equivalently, with 2 different leftmost derivations. Note the pragmatic difference:
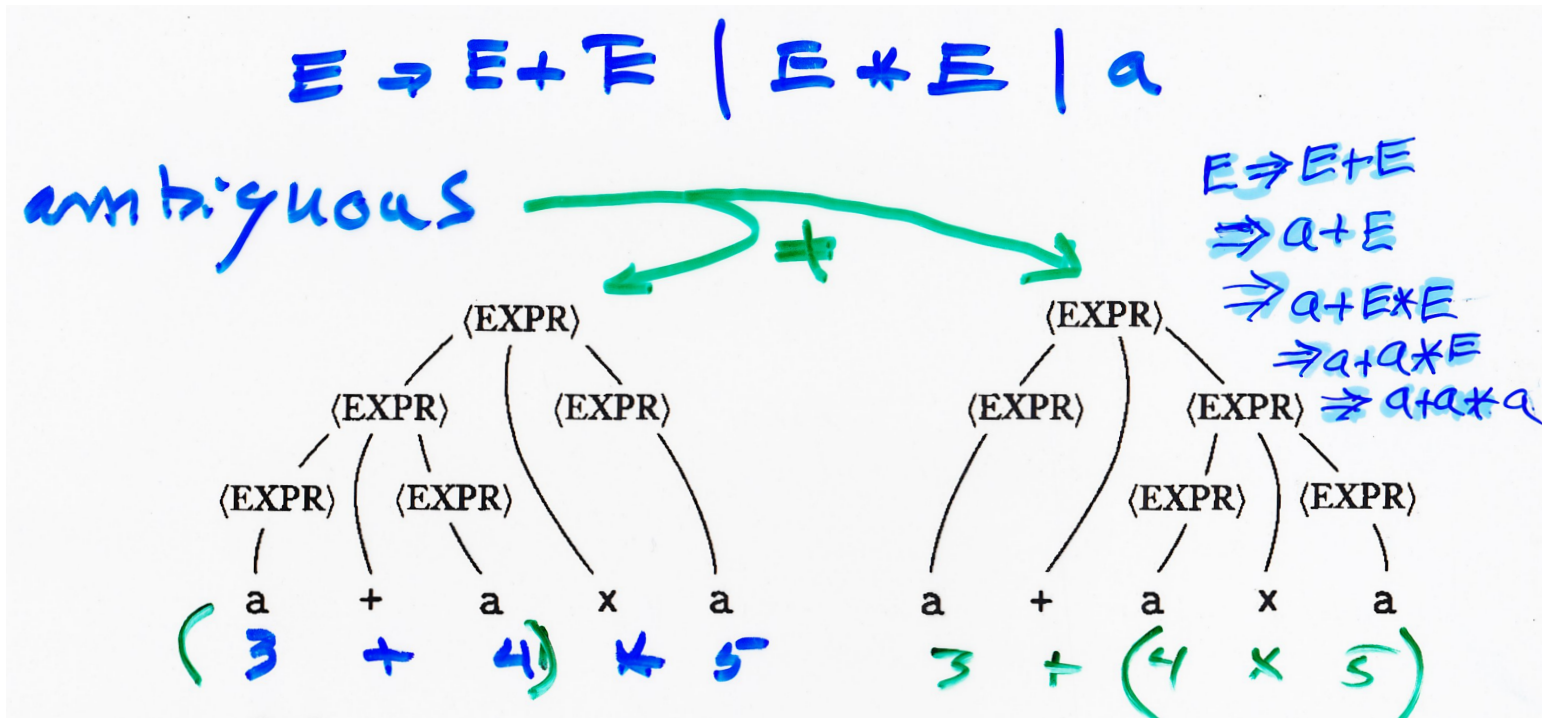in general, (a+a)*a != a+(a*a); which is "right"?

$$E \Rightarrow E+E \mid E*E \mid a$$

ambiguous

$E \Rightarrow E+E$
$\Rightarrow a+E$
$\Rightarrow a+E*E$
$\Rightarrow a+a*E$
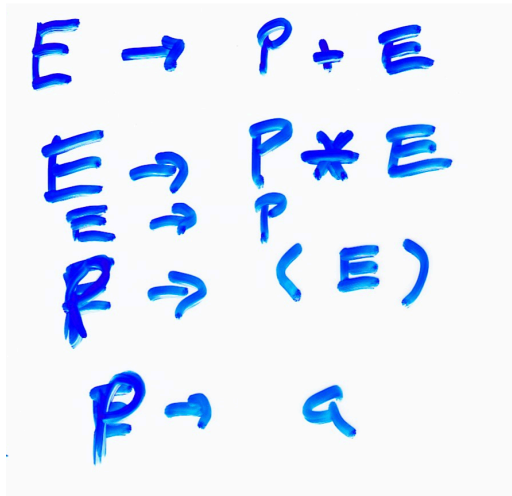$\Rightarrow a+a*a$



Fig 2.6: Two parse trees for a+a×a in grammar G₅

Leftmost deriv

$E \Rightarrow_L E*E \Rightarrow_L E+E*E \Rightarrow a+E*E$
$\Rightarrow_L a+a*E$
$\Rightarrow_L a+a*a$

non-leftmost deriv
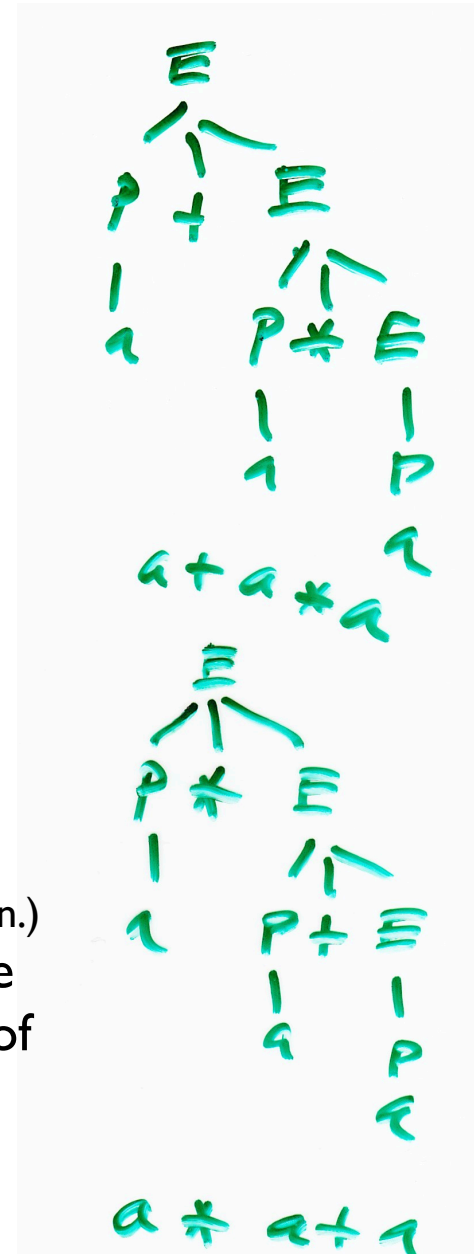
$E*E \Rightarrow E*a \Rightarrow E+E*a \Rightarrow E+a*a$

10

# The "E, P" grammar again

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

This grammar is *un*ambiguous.
  (Why? Very informally, the 3 E rules generate P(('+'∪'*')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
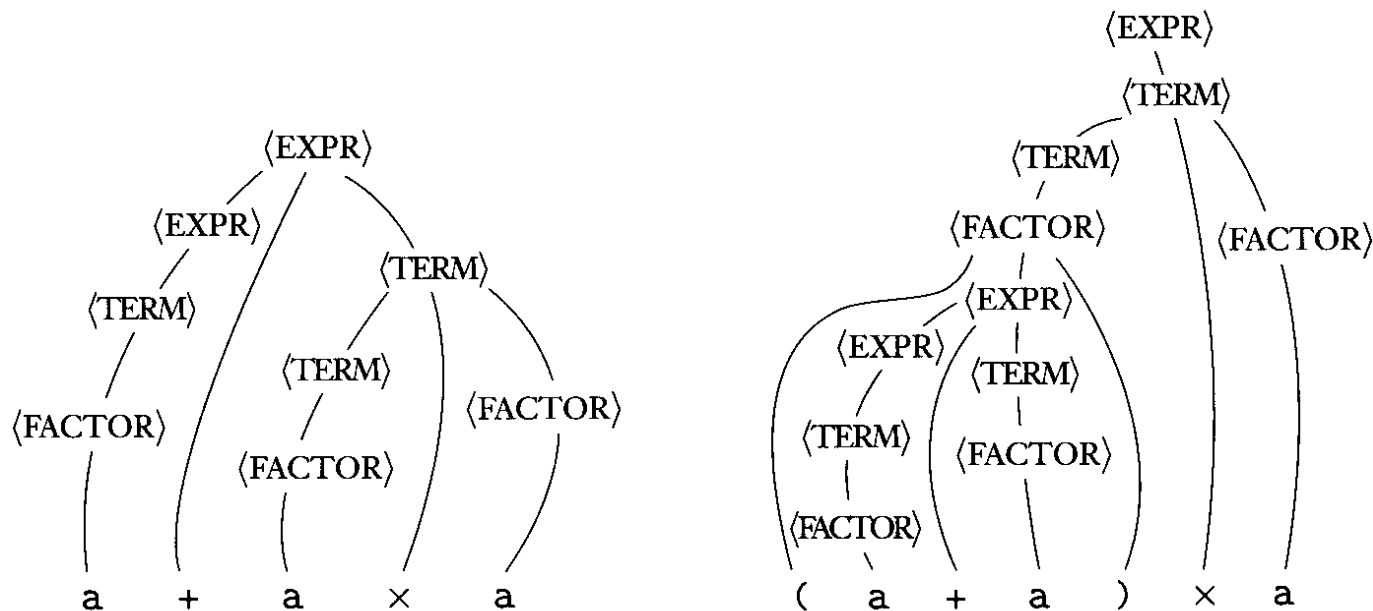* over +.  E.g., tree at lower right suggests
"a * a + a == a * (a + a)"



11

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 $\{\text{a}, +, \times, (,)\}$. The rules are

$$\langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle + \langle \text{TERM} \rangle \mid \langle \text{TERM} \rangle$$
$$\langle \text{TERM} \rangle \rightarrow \langle \text{TERM} \rangle \times \langle \text{FACTOR} \rangle \mid \langle \text{FACTOR} \rangle$$
$$\langle \text{FACTOR} \rangle \rightarrow (\langle \text{EXPR} \rangle) \mid \text{a}$$

The two strings a+axa and (a+a)xa 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.

12

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

$$S \rightarrow AC \mid DB$$

$$A \rightarrow aAb \mid \varepsilon$$

$$C \rightarrow cC \mid \varepsilon$$

$$D \rightarrow aD \mid \varepsilon$$

$$B \rightarrow bBc \mid \varepsilon$$

G

$$a^{10} b^{10} c^{22}$$

$$a^{10} b^{22} c^{22}$$

$$a^{10} b^{10} c^{10}$$

**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

13

# Some closure results for CFLs

## Theorem:
The set of context-free languages is closed under ∪, •, and *


## Corollary:
All regular languages are CFLs

Proof Sketch:
Directly give simple CFLs for $\varnothing$, $\{\varepsilon\}$, and $\{a\}$ for each a $\in$ Σ.  Combine them using the above theorem.

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

# Proof: Closure under Concatenation

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

be 2 CFG's

with $V_1 \cap V_2 = \phi$

let $S \notin V_1 \cup V_2$

Build new grammar

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

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

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

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

$$S_1 \underset{G_1}{\overset{*}{\Rightarrow}} x \;\; \& \;\; S_2 \underset{G_2}{\overset{*}{\Rightarrow}} y$$

$$\therefore \; S \underset{G}{\Rightarrow} S_1 S_2 \underset{G}{\overset{*}{\Rightarrow}} x S_2 \underset{G}{\overset{*}{\Rightarrow}} x y$$

$$\therefore \; L_1 \cdot L_2 \subseteq L(S)$$

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 $G_1$

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

using only $G_2$ rules

$$S_1 \Rightarrow^* x \quad \text{in } G_1$$
$$S_2 \Rightarrow y \quad \text{in } G_2$$  ] **

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

A key issue in this direction of the proof is that, since $V_1 \cap V_2 = \varnothing$, 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 $G_1$ & $G_2$.