# CONTEXT FREE GRAMMARS/LANGUAGES and PUSH DOWN AUTOMATA

We'll now study the class of "context free languages" (CFLs). This class
- includes the class of regular languages.
- is characterized by "context free grammars" (CFGs) in syntactic manner
- is characterized (equivalently) by "push down automata" (PDAs) in computational manner.

We start with grammars (CFGs).

Motivating example from natural language (English)

Sentence $\longrightarrow$ Subject Verb Object

Sentence $\longrightarrow$ Sentence and Sentence

Subject $\longrightarrow$ Jim | Lisa

Verb $\longrightarrow$ ate | threw | killed

Object $\longrightarrow$ apple | ball | rat

This grammar generates sentences such as
- Jim ate apple
- Jim ate apple and Lisa threw ball

- Lisa threw apple and Jim killed rat and Jim ate ball

**Fact** Programming Languages (such as C) are described by a grammar that generates all syntactically correct programs.

**Exercise** Look up the C-grammar.

## Example of a CFG

$$G: \quad S \rightarrow A$$
$$A \rightarrow 0A1$$
$$A \rightarrow \#$$

- 3 "rules".
- $S, A$ variables
- $S$ : start variable
- $0, 1, \#$ terminals

**Idea** Start with the string $S$.

At any step, replace a variable var by string $\alpha$ if var $\rightarrow \alpha$ is a rule. Stop if the string only has terminals and no variable.

**Eg:** $S \rightarrow A \rightarrow 0A1 \rightarrow 00A11 \rightarrow 00\#11$

We say that grammar $G$ <u>derives</u> string $00\#11$
            <u>generates</u>

Clearly this grammar generates precisely the set of strings (referred to as the language L(G) generated by the grammar) described

$$L(G) = \{ 0^n \# 1^n \mid n \geq 0 \}.$$

Note that L(G) is not regular.

Parse tree

The parse tree shows graphical (and structural) representation of the derivation. (of string 000#111 here)

Formal Definition of CFG

Def A CFG G is a 4-tuple $G = (\mathcal{V}, \Sigma, \mathcal{R}, S)$ where

- $\mathcal{V}$ is a finite set of variables
- $\Sigma$ " of terminals, disjoint from $\mathcal{V}$.
- $S \in \mathcal{V}$ is the start variable.

- $R$ is a finite set of rules of type

$$\text{Variable} \longrightarrow \text{String of variables and terminals}$$

i.e. $A \rightarrow (\mathcal{V} \cup \Sigma)^*$, $A \in \mathcal{V}$.

Usually the start variable appears as the left side of the topmost / first rule.

**Def** Suppose $u, v, w \in (\mathcal{V} \cup \Sigma)^*$ and $A \rightarrow w$ is a rule. Then we say that

$$u A v \implies u w v \qquad \text{"} u A v \text{ yields } u w v\text{"}.$$

I.e. one substitutes $A$ by the right side of the rule $w$.

**Def** $u \overset{*}{\implies} v$ if there is a (finite) sequence of strings $u_1, \cdots, u_k \in (\mathcal{V} \cup \Sigma)^*$ such that

$$u_1 = u \implies u_1 \implies u_2 \implies \cdots \implies u_k = v.$$

**Def** Language generated by grammar $G$ is

$$L(G) = \{ w \in \Sigma^* \mid s \overset{*}{\implies} w \}.$$

__Note__. L(G) is (only) over the alphabet $\Sigma$ of terminals.

__Example__

$$G: \quad \begin{array}{l} S \to (S) \\ S \to SS \\ S \to \varepsilon \end{array} \qquad \Sigma = \{\,(\,,\,)\,\}.$$

E.g. $\underline{S} \to (\underline{S}) \to (\underline{S}S) \to ((S)\underline{S})$

$\to ((\underline{S})(S)) \to (()(\underline{S})) \to (()())$

L(G) is the language of all "matched parentheses". The substituted variable is underlined in each step above.

__Example__

$$G_{Add}: \quad \begin{array}{l} E \to E+E \\ E \to 1 \mid 2 \mid 3 \end{array} \qquad \Sigma = \{\,+, 1, 2, 3\,\}$$

$L(G_{Add})$ is all expressions of type

$1 \qquad 1+2 \qquad 1+2+3 \qquad 2+1+1+3+2 \quad$ etc.

## Example

$$G_{Add \cdot Mult} : \quad E \rightarrow E + E \qquad \Sigma = \{+, \times, 1, 2, 3\}$$
$$E \rightarrow E \times E$$
$$E \rightarrow 1 \mid 2 \mid 3$$

$L(G_{Add \cdot Mult})$ has arithmetic expressions such

as $\quad 1 + 2 \times 3 \qquad 2 + 3 \times 1 + 2 \qquad$ etc.

$$\underline{E} \rightarrow \underline{E} \times E \rightarrow \underline{E} + E \times E \rightarrow 1 + E \times \underline{E}$$
$$\rightarrow 1 + \underline{E} \times 3 \rightarrow 1 + 2 \times 3$$

## Example

$$G_{Arith} : \quad E \rightarrow (E + E) \qquad \Sigma = \{+, \times, (, ), 1, 2, 3\}$$
$$E \rightarrow E \times E$$
$$E \rightarrow 1 \mid 2 \mid 3$$

$L(G_{Arith})$ has arithmetic expressions with

parentheses, e.g. $\qquad (1+2) \times 3$

$$\underline{E} \rightarrow \underline{E} \times E \rightarrow (\underline{E} + E) \times E \rightarrow (1 + \underline{E}) \times E$$
$$\rightarrow (1 + 2) \times \underline{E} \rightarrow (1 + 2) \times 3$$

# Example

$$G_{Arith\text{-}Int} : \quad E \to (E+E)$$
$$E \to E \times E$$
$$E \to I$$
$$I \to 0 \mid 1T \mid 2T \mid \cdots \mid 9T$$
$$T \to \varepsilon \mid T0 \mid T1 \mid \cdots \mid T9$$

This generates expressions with decimal integers

e.g. $(214 + 109) \times (51080 + 76541)$

---------------- $\times$ ----------------

It is easy to show that every regular language is context free.

**Theorem** Every regular language is context free.

**Proof** (by example). Suppose a language $L$ is recognized by a DFA such as

The following grammar G "simulates" the DFA.

G:
$$Q_1 \rightarrow 0\,Q_2 \mid 1\,Q_1$$
$$Q_2 \rightarrow 0\,Q_5 \mid 1\,Q_3$$
$$Q_3 \rightarrow 0\,Q_3 \mid 1\,Q_6$$
$$Q_4 \rightarrow 0\,Q_1 \mid 1\,Q_2$$
$$Q_5 \rightarrow 0\,Q_5 \mid 1\,Q_4$$
$$Q_6 \rightarrow 0\,Q_5 \mid 1\,Q_6$$

simulate transitions of the DFA

$$\boxed{\begin{array}{l} Q_5 \rightarrow \varepsilon \\ Q_6 \rightarrow \varepsilon \end{array}}$$

Add these only for accept states of the DFA.

Here $Q_1, Q_2, \cdots, Q_6$ are variables of the grammars and $Q_1$ is the start variable (state).

<u>Exercise</u> Convince yourself that $L = L(G)$ !

It is easily seen that the class of CFLs is closed under $\cup$, $\circ$, $*$.

<u>Theorem</u> If $L_1, L_2$ are context free languags, then so are $L_1 \cup L_2$, $L_1 \circ L_2$, $L_1^*$.

**Proof**  Let the grammars for $L_1, L_2$ be $G_1, G_2$ respectively.
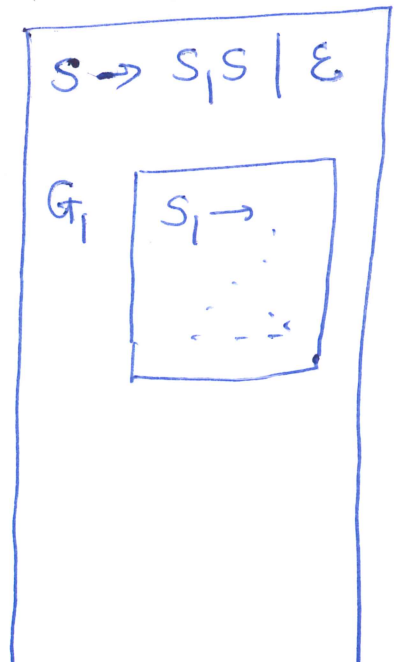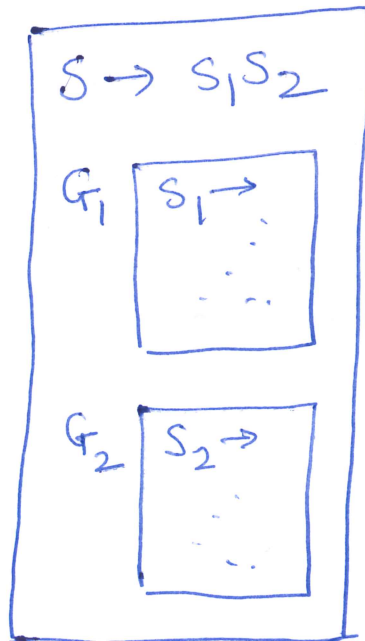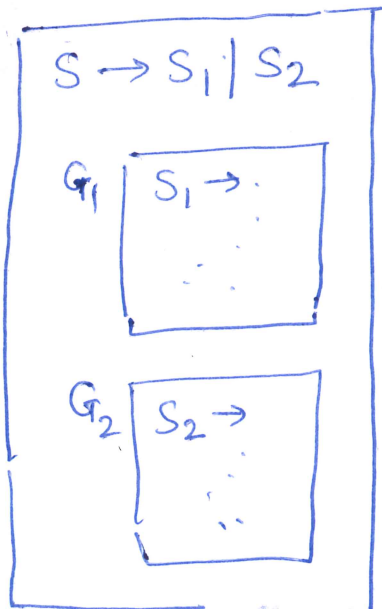
$G_1$

$S_1 \rightarrow$

$G_2$

$S_2 \rightarrow$

where $S_1, S_2$ are their start variables respectively. We may assume that their sets of variables are disjoint (but the set of terminals $\Sigma$ is the same and $L_1, L_2$ are languages over $\Sigma$). The grammar $G$, with a new start variable $S$, for languages $L_1 \cup L_2$, $L_1 \circ L_2$, $L_1^*$ (as the case may be) is:
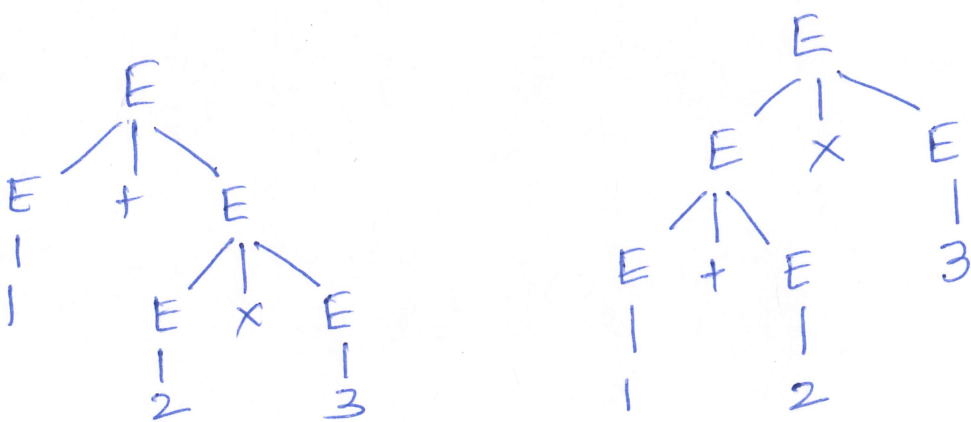
$L_1 \cup L_2$

$L_1 \circ L_2$

$L_1^*$

$L_1 \cup L_2$

$S \rightarrow S_1 \mid S_2$
$G_1 \quad S_1 \rightarrow$
$G_2 \quad S_2 \rightarrow$

$L_1 \circ L_2$

$S \rightarrow S_1 S_2$
$G_1 \quad S_1 \rightarrow$
$G_2 \quad S_2 \rightarrow$

$L_1^*$

$S \rightarrow S_1 S \mid \varepsilon$
$G_1 \quad S_1 \rightarrow$

# Ambiguity of Grammars

Consider the grammar:

$$G_{Add\cdot Mult} : \quad E \rightarrow E+E$$
$$E \rightarrow E \times E$$
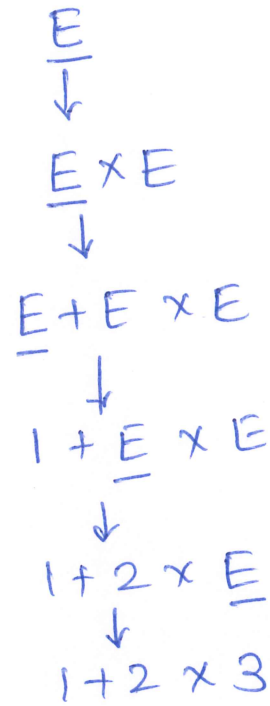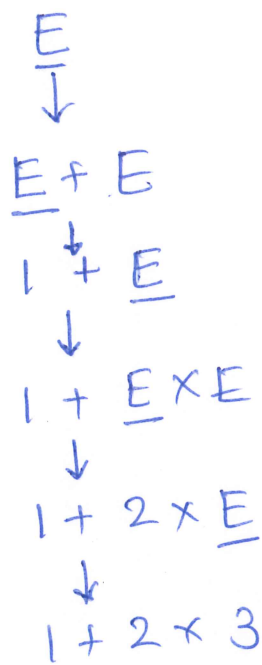$$E \rightarrow 1 \mid 2 \mid 3$$

and the string $1 + 2 \times 3$ generated by it.
There are two distinct parse trees that
yield this string.



It is not difficult to see that every
parse tree corresponds to a unique
"leftmost derivation" and vice versa. Here:

<u>Def</u> A derivation is said to be a leftmost
derivation if at every step, the
leftmost variable is substituted for
using a grammar rule.

The leftmost derivations corresponding to the parse trees above are, respectively:

$$E$$
$$\downarrow$$
$$\underline{E} + E$$
$$\downarrow$$
$$1 + \underline{E}$$
$$\downarrow$$
$$1 + \underline{E} \times E$$
$$\downarrow$$
$$1 + 2 \times \underline{E}$$
$$\downarrow$$
$$1 + 2 \times 3$$

$$E$$
$$\downarrow$$
$$E \times E$$
$$\downarrow$$
$$\underline{E} + E \times E$$
$$\downarrow$$
$$1 + \underline{E} \times E$$
$$\downarrow$$
$$1 + 2 \times \underline{E}$$
$$\downarrow$$
$$1 + 2 \times 3$$

The two derivations are "intrinsically different." If semantics are attached to the terminals $+, \times$ then these amount to whether in the expression $1 + 2 \times 3$, the addition is evaluated first or the multiplication. The grammar is said to be ambiguous –

**Def.** A grammar $G$ is ambiguous if there is some string $w \in L(G)$ that has two (or more) different leftmost derivations (or equivalently parse trees).
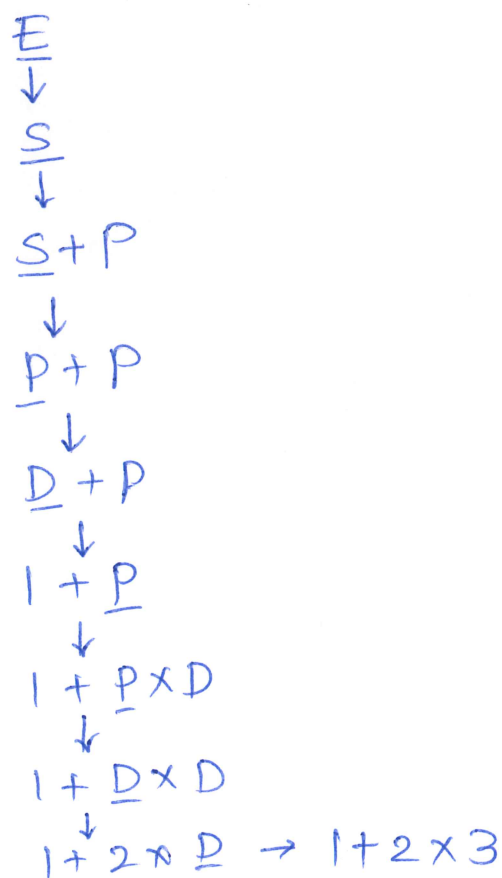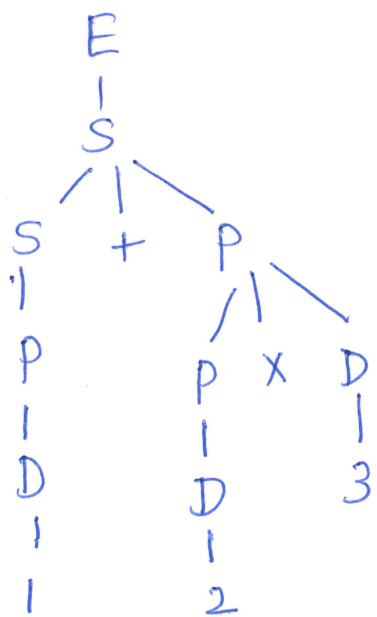
It is possible to rewrite this grammar as below so that the new grammar is unambiguous.

$$G'_{Add \cdot mult} : \quad E \to S$$
$$S \to S+P \mid P$$
$$P \to P \times D \mid D$$
$$D \to 1 \mid 2 \mid 3$$

'S' = 'Sum'
'P' = 'Product'
'D' = 'Digit'

The interpretation is that the overall arithmetic expression is viewed as "sum of products". Now the string $1+2\times3$ has only one parse tree and a leftmost derivation:

E
|
S
/ | \
S + P
|
P
|
D
|
1

P × D
|
D
|
2

3

E
↓
S
↓
S+P
↓
P+P
↓
D+P
↓
1+P
↓
1+P×D
↓
1+D×D
↓
1+2×D → 1+2×3

I.e. the parse tree on the left (before) is "preferred".

**Exercise** Convince yourself that
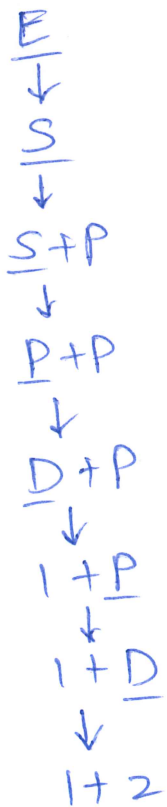
- $L(G'_{Add-Mult}) = L(G)$ and

- $G'_{Add-Mult}$ is unambiguous.

**Note** We emphasize that $G'_{Add-Mult}$ may yield two different derivations for the same string,

e.g. for the string $1+2$

$$E \to \underline{S} \to \underline{S}+P \to \underline{P}+P \to D+\underline{P} \to D+\underline{D} \to \underline{D}+2 \to 1+2$$

$$E \to \underline{S} \to S+\underline{P} \to S+\underline{D} \to \underline{S}+2 \to \underline{P}+2 \to \underline{D}+2 \to 1+2$$

However, there is <u>only one</u> <u>leftmost derivation</u> <u>and a parse tree</u>, i.e.

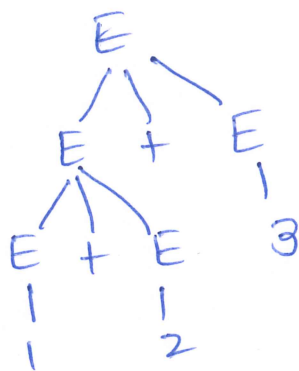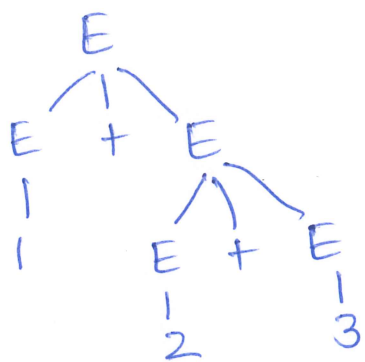

This is why we insist on leftmost derivations (or parse trees).

**Example** This grammar is ambiguous.

$$G_{Add}: \quad E \rightarrow E+E$$
$$E \rightarrow 1 \mid 2 \mid 3.$$

Since $1+2+3$ has two parse trees.



It is easy to rewrite an equivalent, unambiguous grammar.

$$G'_{Add}: \quad E \rightarrow E+D \mid D$$
$$D \rightarrow 1 \mid 2 \mid 3$$

**Example** If-then-else grammar as below is ambiguous.

$$S \rightarrow \text{if } E \text{ then } S$$
$$S \rightarrow \text{if } E \text{ then } S \text{ else } S$$
$$S \rightarrow \text{other}$$

Here if, then, else, other are terminals. 'S' stands for statement, 'other' stands for other statement, 'E' stands for expression (or

logical condition) that we don't really care about.
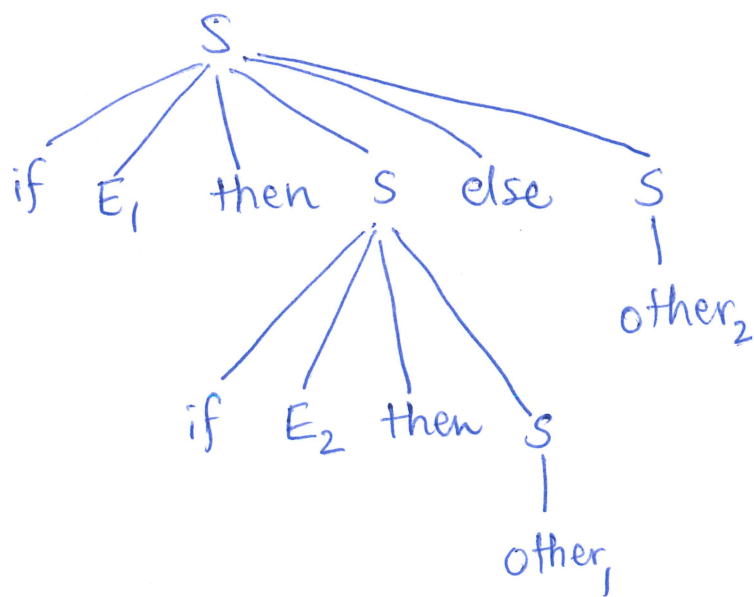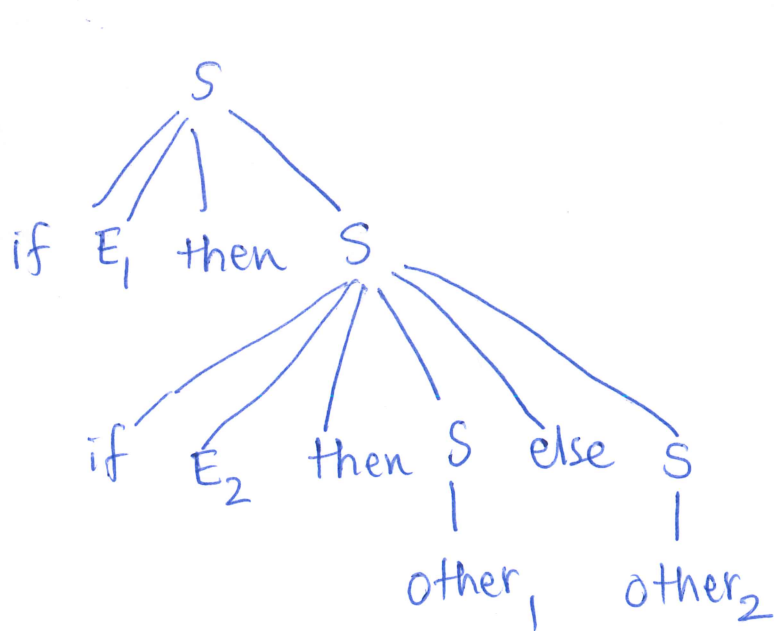
Consider the following statement.

if $E_1$ then if $E_2$ then other$_1$ else other$_2$

There are two ways to (semantically) interpret the statement:

if $E_1$ then { if $E_2$ then other$_1$ else other$_2$ }

if $E_1$ then { if $E_2$ then other$_1$ } else other$_2$

These lead to two different parse trees:
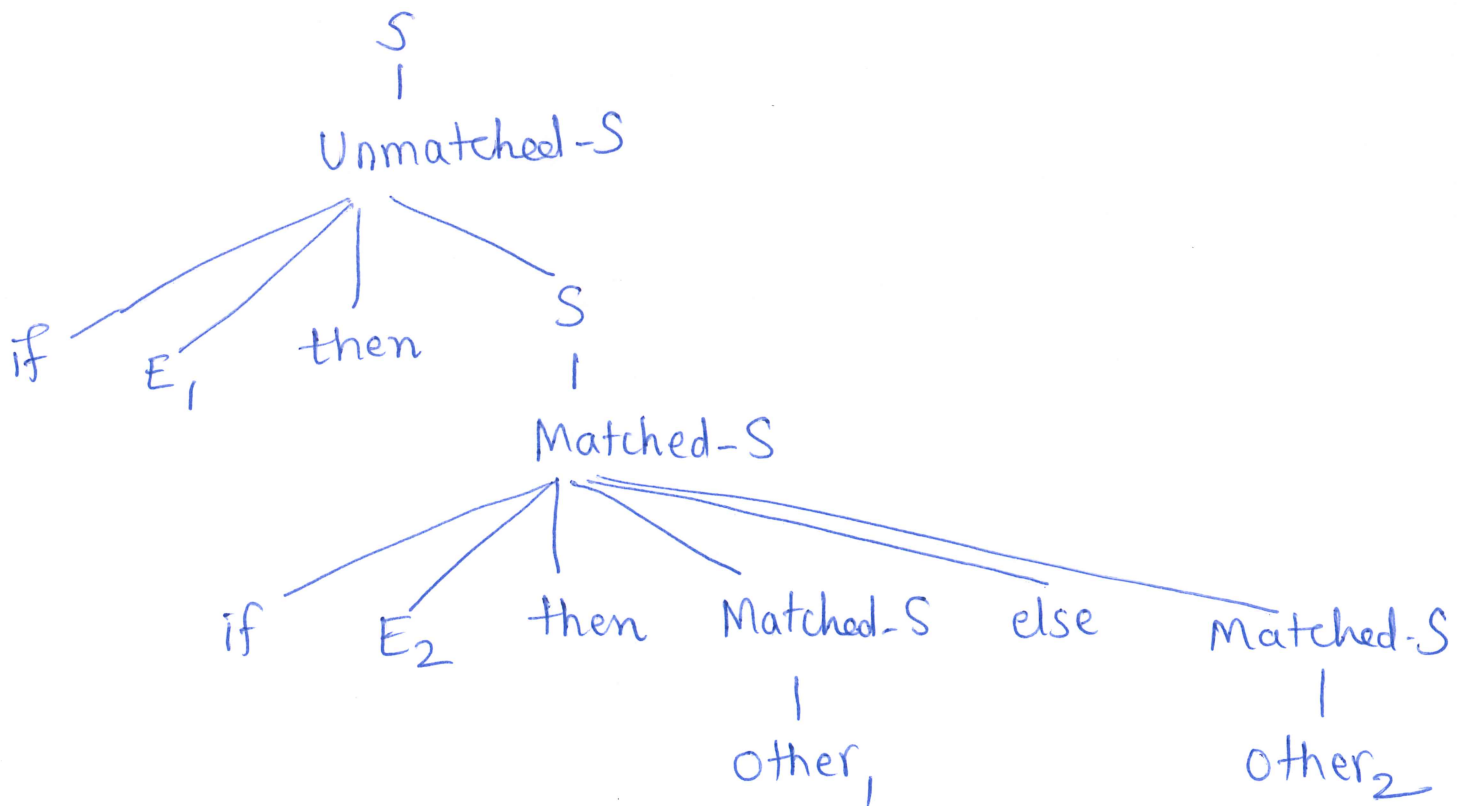


Hence the grammar is ambiguous.

There is a text-book manner to rewrite an equivalent unambiguous grammar.

$S \rightarrow$ Matched-S | Unmatched-S

Matched-S $\rightarrow$ if E then Matched-S else Matched-S
| other

Unmatched-S $\rightarrow$ if E then S |
if E then Matched-S else Unmatched-S

The statement before has the parse tree:

S
|
Unmatched-S
/ | \
if  E$_1$  then   S
|
Matched-S
/ | \ \ \
if  E$_2$  then  Matched-S  else  Matched-S
|                    |
other$_1$            other$_2$

I.e. the left parse tree (before) is "preferred".

Sometimes it is impossible to write a grammar for a CFL that is unambiguous.

Def A CFL is called inherently ambiguous if _every_ grammar for it is ambiguous.

Example $L = \{a^i b^j c^k \mid i=j \text{ or } j=k, \; i,j,k \geq 0\}$ is inherently ambiguous. We'll not prove this in this class.

Note $L_1 = \{a^i b^j c^k \mid i=j\}$
$L_2 = \{a^i b^j c^k \mid j=k\}$ are both context free.

However, as we prove later,
$$L_1 \cap L_2 = \{a^i b^j c^k \mid i=j=k\} \text{ is not}$$
context free. Thus the class of CFLs is _not_ closed under intersection and hence also _not_ closed under complements (why? $\overline{A} \cup \overline{B} = \overline{A \cap B}$). In this regard, CFLs differ from regular languages.