

CMSC 141 - Automata and Language Theory

Handout: Context-Free Grammars

By: Gianina Renee V. Vergara

Context Free Grammars

- A context-free grammar, CFG is a 4-tuple (V, T, P, S) ...
- Each production is of the form $A \rightarrow \alpha$, where $A \in V$, $\alpha \in (V \cup T)^*$

Application

- Used in actual problems in programming languages and other areas of computer science.
- Used in the study of human language.
- Most compilers and interpreters contain a component called a parser that extracts the meaning of a program prior to generating a compiled code or performing the interpreted execution.

Example 1: CFG that tries to capture if-then-else statements. In this notation, the variables/non-terminals are those enclosed by "<>" symbols, and the terminals are in bold.

```
<start> → if (<condition>) then { <start> }  
<start> → if (<condition>) then {<start>} else {<start>}  
<start> → <stmt>  
<condition> → true | false  
<stmt> → printf("Hello World!");
```

Def'n: A derivation or parse tree is a graphic method of showing derivations of a string from the start symbol.

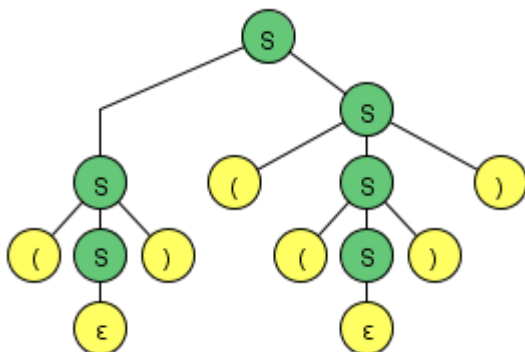
Def'n: a CFG is **ambiguous** iff it generates some sentence by two or more distinct leftmost (rightmost) derivations. Or, a CFG is ambiguous iff it generates some sentence by two or more distinct parse trees.

Example 2: Consider the grammar $G = (\{S\}, \{a,b\}, P, S)$. The set of production is:

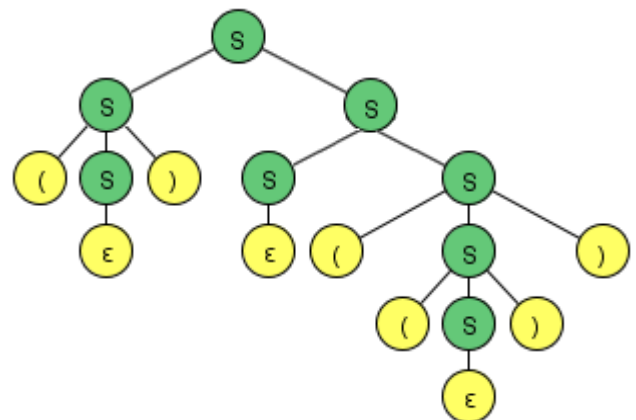
$S \rightarrow SS \mid (S) \mid \epsilon$

the string $w = ()()$ can be produced by G using two distinct parse trees:

a.



b.



Hence, this CFG is ambiguous. $L(G)$ is the language of all strings of properly nested parenthesis.

Def'n: Sometimes when we have an ambiguous grammar, we can find an unambiguous grammar that generates the same language. Some CFL, however, can only be generated by ambiguous grammars. Such grammars are called **inherently ambiguous**.

Brute Force Parser - Regular or Context Free Grammar with JFLAP

Multi-Input Run Grammar with JFLAP

Designing CFG's

Some basic techniques:

- CFG's for regular languages - RL are also CFL. It will be easier to generate the CFG from the equivalent DFA of the regular language.
- Matching - enforcing the fact that your grammar generates matching pairs of symbols.

Example 3: $L = \{0^n 1^{2n} \mid n \geq 0\}$. The context-free grammar for L is

$$S \rightarrow 0S11 \mid \epsilon$$

The grammar forces every 0 to match 11.

- Using recursive relationships - when you can represent strings in the language in terms of shorter strings in the language. The variable symbol should be placed in the location where it may recursively appear.

Example 4: L - the language of all strings of balanced brackets. Every string w in L can be viewed as either

- uv , for $u \in L$, $v \in L$, or
- (u) , for $u \in L$.

This gives the following grammar:

$$S \rightarrow SS \mid (S) \mid \epsilon$$

- Thinking of each non-terminal A as representing a language of all strings w derivable from A . Then, combining these languages into one, using the rules.

Example 5: $L = \{0^n 1^m \mid n \neq m\}$.

Let $L_1 = \{0^n 1^m \mid n > m\}$, and $L_2 = \{0^n 1^m \mid n < m\}$. Then, if S_1 generates L_1 , and S_2 generates L_2 , our grammar will be

$$S \rightarrow S_1 \mid S_2$$

L_1 is just the language of strings $0^n 1^m$ with one or more extra 0's in front. So,

$$S_1 \rightarrow 0S_1 \mid 0A$$

$$A \rightarrow 0A1 \mid \epsilon$$

L_2 is just the language of strings $0^n 1^m$ with one or more extra 1's in the end. So,

$$S_2 \rightarrow S_21 \mid A1$$

$$A \rightarrow 0A1 \mid \epsilon$$

$$L = L_1 \cup L_2$$

$$S \rightarrow S_1 \mid S_2$$

$$S_1 \rightarrow 0S_1 \mid 0A$$

$$S_2 \rightarrow S_21 \mid A1$$

$$A \rightarrow 0A1 \mid \epsilon$$

Chomsky Normal Form (CNF)

- A CFG is in CNF if every rule is of the form: $A \rightarrow BC$ or $A \rightarrow a$, where a is any terminal and A , B , C are any non-terminal - (*in some books B and C may not be the start variable*). In addition we permit the rule $S \rightarrow \epsilon$.
 1. Ensure S does not appear on RHS, add new S_0 if necessary and copy all rules of S .
 2. Eliminate all ϵ rules, except $S \rightarrow \epsilon$.
 3. Eliminate unit rules
 4. Convert rules to proper form.

- **Eliminating ϵ rules:**

- Pick a ϵ -rule, $Y \rightarrow \epsilon$, and remove it.
- Given a rule $X \rightarrow \alpha$ where α contains n occurrences of Y , replace it with 2^n rules in which 0, ..., n occurrences are replaced by ϵ . (Do not add $X \rightarrow \epsilon$ if previously removed.)

e.g.

1. Remove $Y \rightarrow \epsilon$..
2. Replace X with ϵ occurrences of Y ..
3. New rules..

$X \rightarrow aYcbY,$	$X \rightarrow aYcbY \mid a\epsilon cbY \mid aYcb\epsilon \mid$	$X \rightarrow aYcbY \mid acbY \mid aYcb \mid$
$Y \rightarrow a \mid \epsilon$	$a\epsilon cb\epsilon,$	$acb,$
	$Y \rightarrow a$	$Y \rightarrow a$

- **Eliminating unit rules:**

- Pick a unit rule $A \rightarrow B$ and remove it.
- For every rule $B \rightarrow u$, add rule $A \rightarrow u$ unless this is a unit rule that was previously removed

e.g.

1. Remove $S \rightarrow A$..
2. add occurrences of A to S

$S \rightarrow A,$	$S \rightarrow aB \mid b,$
$A \rightarrow aB \mid b,$	$A \rightarrow aB \mid b,$
$B \rightarrow aS$	$B \rightarrow aS$

CNF Example

	remove $X \rightarrow \epsilon$	remove $S \rightarrow X$	follow format ✓	CNF
$S \rightarrow XX$	$S \rightarrow XX$	$S \rightarrow XX$	$S \rightarrow XX$ ✓	$S \rightarrow XX$
$X \rightarrow aXb$	$S \rightarrow X$	$S \rightarrow aXb$	$S \rightarrow AXB$	$S \rightarrow AB$
$X \rightarrow \epsilon$	$S \rightarrow \epsilon$	$S \rightarrow ab$	$S \rightarrow AB$ ✓	$A \rightarrow a$
	$X \rightarrow aXb$	$S \rightarrow \epsilon$	$A \rightarrow a$ ✓	$B \rightarrow b$
	$X \rightarrow ab$	$X \rightarrow aXb$	$B \rightarrow b$ ✓	$S \rightarrow \epsilon$
		$X \rightarrow ab$	$S \rightarrow \epsilon$ ✓	$X \rightarrow AB$
			$X \rightarrow AXB$	$S \rightarrow AZ$
			$X \rightarrow AB$ ✓	$Z \rightarrow XB$
				$X \rightarrow AZ$

[Transform Grammar with JFLAP](#)