

Lecture 9: Context Free Languages and Grammars

CSC 320: Foundations of Computer Science

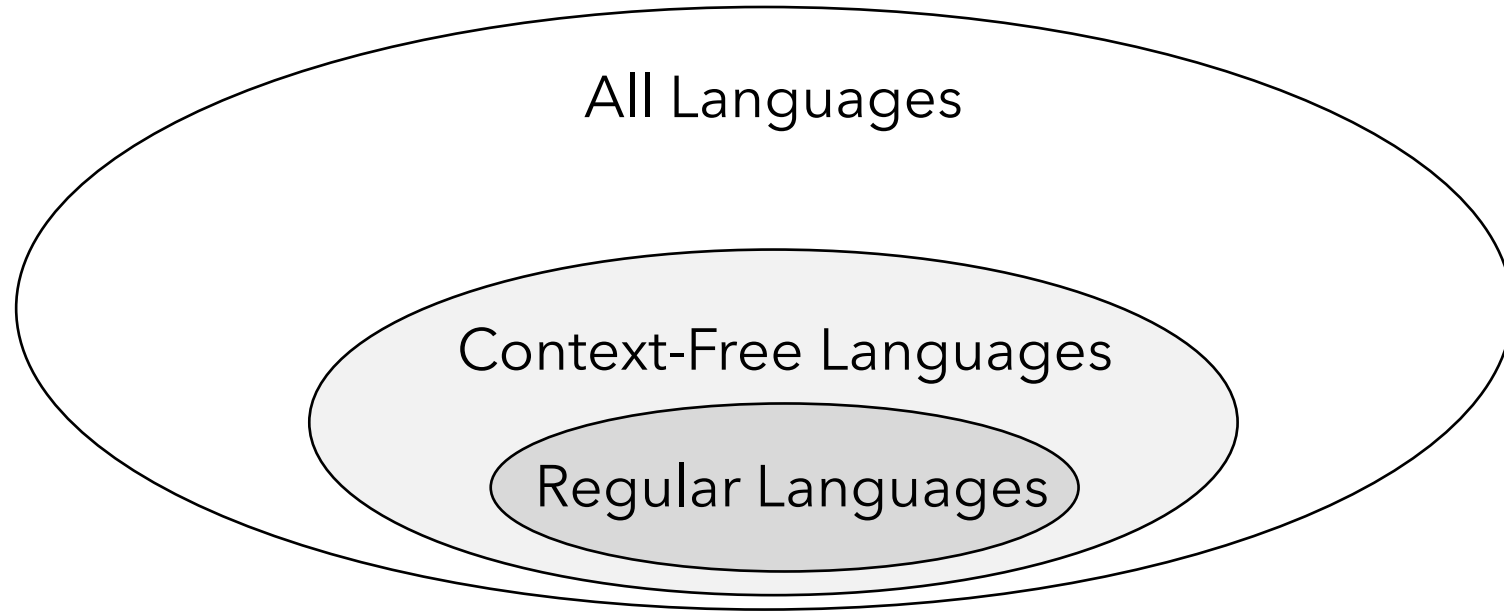
Quinton Yong

quintonyong@uvic.ca



**University
of Victoria**

Context-Free Languages



- The set of **regular languages** is the set of all languages recognized by DFAs, NFAs, and regular expressions
- The set of **context-free languages** is a **superset** of the regular languages

Grammars

Consider the following “more powerful” way to represent a language:

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$

This **grammar** G describes a language, which contains the strings which can be **derived** by the grammar.

We can derive a string from a grammar as follows:

$$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$$

Context-Free Grammars (CFG)

This is an example of a **context-free grammar** (CFG) G .

Start variable A $A \rightarrow 0A1$ substitution / production rule
 $A \rightarrow B$ variables $\{A, B\}$
 $B \rightarrow \#$ terminals $\{0, 1, \#\}$

A CFG consists of the following:

- A collection of **substitution rules** (also called **production rules**)
 - $\langle \text{variable} \rangle \rightarrow \langle \text{variable and terminals} \rangle$
- **Variable** symbols (use capital letters)
- **Terminal** symbols (alphabet symbols)
- One **start variable** (left-hand side of topmost substitution rule)

Context-Free Grammars (CFG)

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$

We **derive** strings using a CFG G as follows:

1. Write down the **start variable**

A

2. **Replace a variable** with the right side of a rule starting with that variable

$0A1$

3. **Repeat step 2** until no variables remain

$0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$

Parse Tree

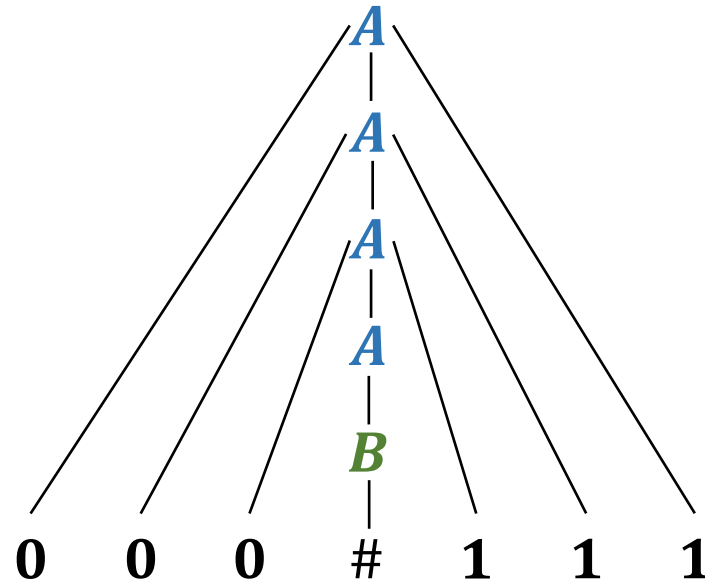
$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$

We can represent the **derivation** of a string pictorially using a **parse tree**

$A \rightarrow 0A1$

$A \rightarrow B$

$B \rightarrow \#$



Context-Free Grammars (CFG)

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$

As a shorthand, we can write **multiple substitution rules** with the same left-hand variable as a **single rule separated by |**

$$A \rightarrow 0A1 \mid B$$

$$B \rightarrow \#$$

Context-Free Grammars (CFG)

- **All strings** that can be generate by a **CFG G** is the **language of G**
- We write **$L(G)$** for the language of grammar **G**

Context-Free Grammars (CFG)

- **All strings** that can be generate by a **CFG** G is the **language of** G
- We write $L(G)$ for the language of grammar G
- What is the language of the following grammar G ?

$$A \rightarrow 0A1 \mid B$$

$$B \rightarrow \#$$

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

Formal Definition: Context-Free Grammar

A **context-free grammar** is a 4-tuple (V, Σ, R, S) where

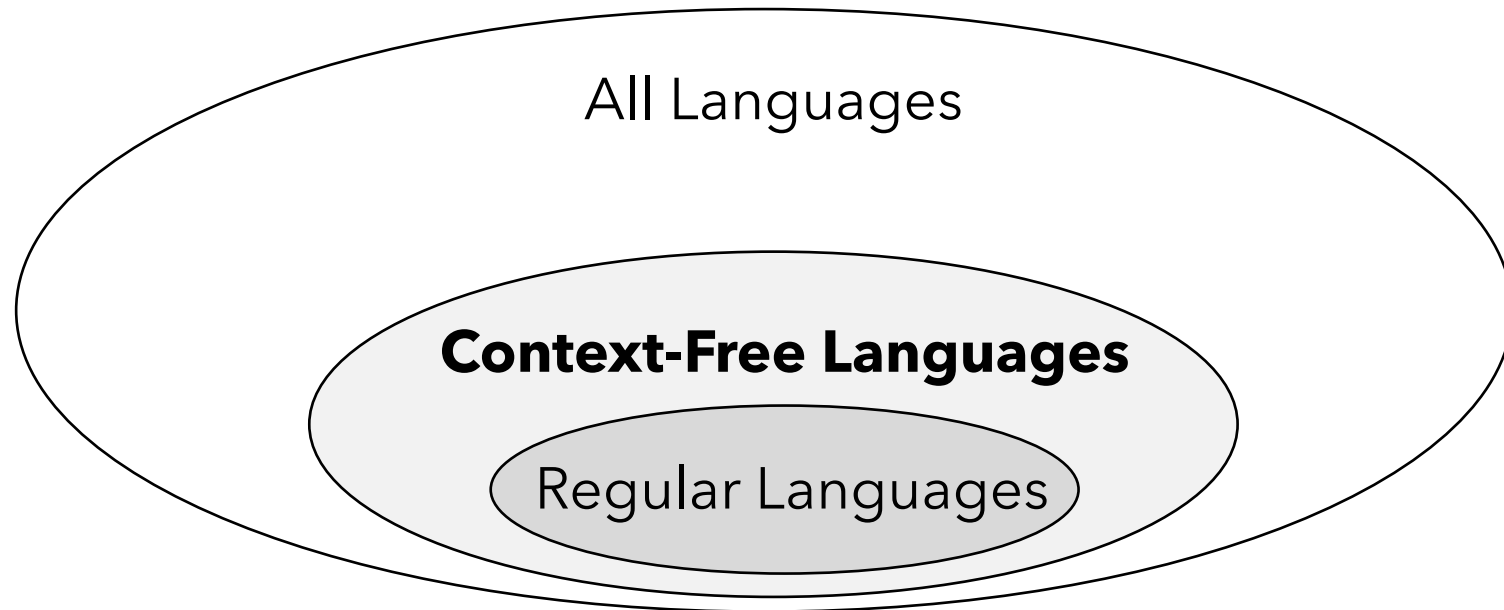
- V is a finite set called the **variables**
 - Σ is a finite set called the **terminals** (disjoint from V)
 - R is a finite set of (substitution / production) **rules**
 - Each rule is a variable to be substituted by a string of variables and terminals
 - $S \in V$ is the **start variable**
-
- **Note:** the right-hand side of a rule may be ε , but we don't add ε to Σ

Terminology

Given a grammar $G = (V, \Sigma, R, S)$

- Let u , v , and w be strings of variable and terminals
- Let $A \rightarrow w$ be a rule of G
- We say that uAv **yields** uwv , denoted $uAv \Rightarrow uwv$ (i.e. substitute one variable)
- We say that u **derives** v , denoted $u \xRightarrow{*} v$, if a sequence u_1, u_2, \dots, u_k exists for some $k \geq 0$ and $u \Rightarrow u_1 \Rightarrow u_2 \Rightarrow \dots \Rightarrow u_k \Rightarrow v$
- The **language of grammar** G is $L(G) = \{w \in \Sigma^* \mid S \xRightarrow{*} w\}$

Context-Free Languages



The class of **context-free languages** is the class of languages recognized by **context-free grammars**

Example 1

Produce a grammar for the language $\{0^n 1^n \mid n \geq 0\}$:

$$S \rightarrow 0S1 \mid \varepsilon$$

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

- $V = \{S\}$
- $\Sigma = \{0, 1\}$
- R is the set of rules

$$S \rightarrow 0S1 \mid \varepsilon$$

Example 2

Consider grammar $G = (\{S\}, \{a, b\}, R, S)$ where R is

$$S \rightarrow aSb \mid SS \mid \varepsilon$$

What is $L(G)$?

- Example strings in the language: $abab, aaabbb, aababb$
- May be easier to see if we set $a := ($ and $b :=)$
- Example strings in the language: $()(), ((())), (()())$

$L(G)$ is the language of all strings of properly nested parentheses

Example 3

Consider grammar $G = (V, \Sigma, R, E)$ with $V = \{E, F, T\}$, $\Sigma = \{a, +, \cdot, (,)\}$, and R is given by:

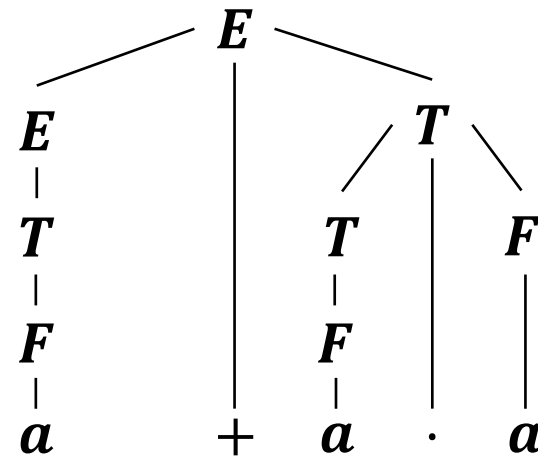
$$E \rightarrow E + T \mid T$$

$$T \rightarrow T \cdot F \mid F$$

$$F \rightarrow (E) \mid a$$

Show how we can derive the string $a + a \cdot a$ using a **parse tree**

Parse trees for this grammar
reflect order of operations



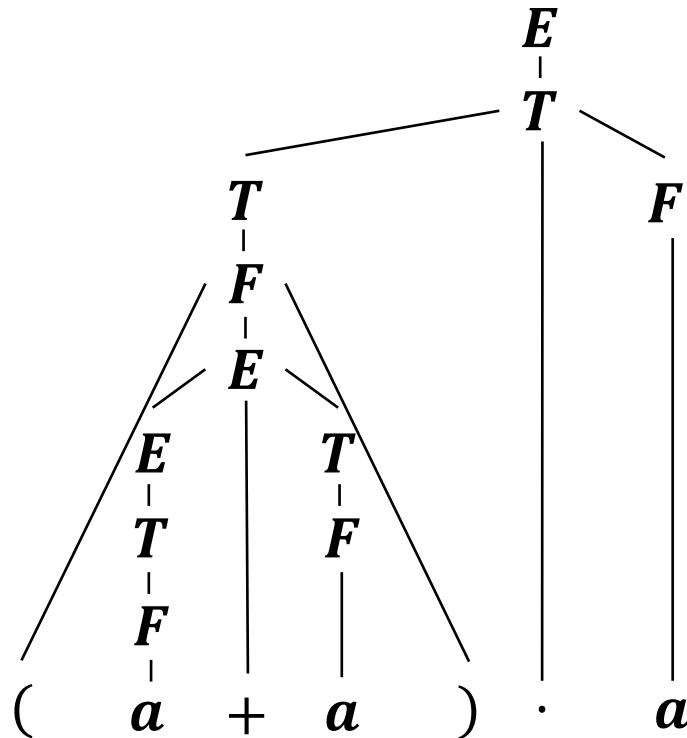
Example 3b

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T \cdot F \mid F$$

$$F \rightarrow (E) \mid a$$

Show how we can derive the string $(a + a) \cdot a$ using a **parse tree**



Example CFG for Programming Languages

Func \rightarrow *Id* (*PrmtrList*) *Stmt*

PrmtrList $\rightarrow \epsilon \mid \dots$

Stmt $\rightarrow \{ \textit{StmtList} \} \mid \textit{Id} = \textit{Expr}; \mid \textbf{if} (\textit{Expr}) \textit{Stmt} \mid \dots$

StmntList $\rightarrow \textit{Stmnt} \mid \textit{StmntList} \textit{Stmt}$

Expr $\rightarrow \textit{Id} \mid \textit{Num} \mid \textit{Expr} \textit{Optr} \textit{Expr}$

Id $\rightarrow \mathbf{a} \mid \mathbf{b} \mid \mathbf{c} \mid \dots$

Num $\rightarrow \mathbf{0} \mid \mathbf{1} \mid \mathbf{2} \mid \dots$

Optr $\rightarrow + \mid - \mid > \mid \dots$

Exercise

Produce a context-free grammar for the language

$\{ w \in \{0, 1\}^* \mid w \text{ has an equal number of 0s and 1s} \}$

Leftmost Derivation

- When deriving a string by substituting variables, there's **normally no rule** on which variable to substitute first
- However, it is useful to **add some structure** to substituting variables

Leftmost derivation: We call a derivation in grammar G a **leftmost derivation** if at every step, the **leftmost remaining variable** is replaced

$$E \rightarrow E + T \mid T$$

leftmost derivation of $a + a$

$$T \rightarrow T \cdot F \mid F$$

$$E \Rightarrow E + T \Rightarrow T + T \Rightarrow F + T \Rightarrow a + T \Rightarrow a + F \Rightarrow a + a$$

$$F \rightarrow (E) \mid a$$

Ambiguous Grammars

- A string is derived **ambiguously** in a CFG if it has at **least two different leftmost derivations**
 - i.e. can derive the string in multiple ways even if always substituting the leftmost variable
- A context-free grammar is **ambiguous** if there **exists a string** that can be derived ambiguously

Ambiguous Grammars Example

Consider grammar $G = (V, \Sigma, R, E)$ with $V = \{E, F, T\}$, $\Sigma = \{a, +, \cdot, (,)\}$, and R is given by:

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

Is this grammar ambiguous? **Yes**

Consider the string $a \cdot a + a$

- Leftmost derivation 1:

$$E \Rightarrow E \cdot E \Rightarrow a \cdot E \Rightarrow a \cdot E + E \Rightarrow a \cdot a + E \Rightarrow a \cdot a + a$$

- Leftmost derivation 2:

$$E \Rightarrow E + E \Rightarrow E \cdot E + E \Rightarrow a \cdot E + E \Rightarrow a \cdot a + E \Rightarrow a \cdot a + a$$