Programming Language Concepts

CMPSC 461

PennState
College of Engineering

ELECTRICAL ENGINEERING
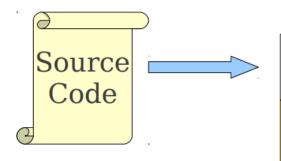AND COMPUTER SCIENCE

Grammar

Professor: Suman Saha

# Where we are?

Source Code

Lexical Analysis

**Syntax Analysis**

Semantic Analysis

IR Generation

IR Optimization

Code Generation

Optimization

Machine Code

# What is Syntax Analysis?

- After lexical analysis (scanning), we have a series of tokens.
- In Syntax analysis (or parsing), we want to interpret what those tokens mean.
- Goal: Recover the *structure* describe by that series of tokens.
- Goal: Report errors if those tokens do not properly encode a structure.

# Formal Languages

- An alphabet is a set Σ of symbols that act as letters.
- A language over Σ is a set of strings made from symbols in Σ.
- When scanning, our alphabet was ASCII or Unicode characters. We produced tokens.
- When parsing, our alphabet is the set of tokens produced by the scanner.

# The Limits of Regular Languages

- When scanning, we used regular expressions to define each token.
- Unfortunately, regular expressions are (usually) too weak to define programming languages.
  - Cannot define a regular expression matching all expressions with properly balanced parentheses.
  - Cannot define a regular expression matching all functions with properly nested block structure.
- We need a more powerful formalism.

# Grammars

- It is written in a metalanguage
- It defines all the legal strings of characters that can form a syntactically valid program

# Context-Free Grammars

- Context-Free Grammars
  - Developed by Noam Chomsky in the mid-1950s
  - Describe the syntax of natural languages
  - Define a class of languages called context-free languages
  - Was originally designed for natural languages

# Context-Free Grammars

- Using the notation of Backus-Naur Form (BNF) to describe CFG
- A grammar $G <N, T, P, S>$ consists of the following

  - A finite set $N$ of non-terminal symbols

  - A finite set $T$ of terminal symbols, that is disjoint from N

  - A finite set $P$ of production rules of the form

    $$A \rightarrow \omega$$

    where $\omega$ is a string of nonterminal and terminal

  - Start symbol

# Backus-Naur Form (BNF) Grammars

- A rule has a left-hand side (LHS), one or more right-hand side (RHS), and consists of terminal and nonterminal symbols
  - For instance
    - $<\text{binaryDigit}> \rightarrow 0$
    - $<\text{binaryDigit}> \rightarrow 1$
  - We can write $<\text{binaryDigit}> \rightarrow 0 \mid 1$

$$A \longrightarrow \omega$$

# Extended BNF Grammar

- Extended BNF simplifies writing a grammar by introducing metasymbols for iteration; option, and choice

- BNF

$$\rightarrow \quad \text{<expr>} := \quad \text{<expr>} + \text{<term>}$$
$$| \text{<expr>} - \text{<term>} \quad | \quad \text{<term>}$$

- EBNF

$$\text{<expr>} := \text{<expr>} \{(+ \mid -) \text{<term>}\}^* \mid \text{<term>}$$

$+ / -$

$e_1 = e_2 + 3$

$e_1 = e_2 - 3$

$e_1 = 3$

$e = e$
$e = e \, (+ t)$

# Extended BNF Grammar

- BNF

$$<ifStmt> ::= \quad if (<expr>) \ <stmt>$$
$$| \ if (<expr>) <stmt> \ else <stmt>$$

- EBNF

$$<ifStmt> ::= if (<expr>) \ <stmt> \ [else <stmt>]$$

0|1 (?)

# Extended BNF Grammar

- However, EBNF is any more powerful than BNF for formally describing language syntax

$$A \rightarrow x \{y\}^* z$$

- Equivalent to

$$A \rightarrow x A' z$$
$$A' \rightarrow \varepsilon \mid y A'$$

$$A \rightarrow x z$$
$$A \rightarrow x y z$$
$$A \rightarrow x y y z$$

$$A \rightarrow x A' z$$
$$\rightarrow x \varepsilon z$$
$$\rightarrow x z$$

$$A \rightarrow x A' z$$
$$\rightarrow x y A' z$$
$$\rightarrow x y z$$

$$A \rightarrow x A' z$$
$$\rightarrow x y A' z$$
$$\rightarrow x y y A' z$$
$$\rightarrow x y y z$$

# Derivation

- To determine that the given string of symbols belongs to grammar
- A sequence of steps where nonterminals are replaced by the right-hand side of a production is called a derivation.
  - Leftmost derivation
  - Rightmost derivation

- Sentential form vs Sentence
  - A *sentential form* is any string derivable from the start symbol.
  - A *sentence* is a sentential form consisting only of terminals

- If string α derives string ω, we write α ⇒* ω.

- Say, we have grammar

  Integer → Digit | Integer Digit

  Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

# Leftmost Derivation

- Say, we have grammar

  Integer → Digit | Integer Digit

  Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9


- 352 is an Integer?

- Say, we have grammar

  Integer → Digit | Integer Digit

  Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9


- 352 is an Integer?

  Integer          → Integer Digit

# Leftmost Derivation

- Say, we have grammar

    Integer → Digit | Integer Digit

    Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

- 352 is an Integer?

    Integer          → <span style="color:red">Integer</span> Digit

                          → <span style="color:#29ABE2">Integer Digit</span> Digit

# Leftmost Derivation

- Say, we have grammar

    Integer → Digit | Integer Digit

    Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

- 352 is an Integer?

    Integer        → Integer Digit

    → Integer Digit Digit

    → Digit Digit Digit

# Leftmost Derivation

- Say, we have grammar

  Integer → Digit | Integer Digit

  Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

- 352 is an Integer?

  Integer      → Integer Digit

            → Integer Digit Digit

            → <span style="color:red">Digit</span> Digit Digit

            → <span style="color:blue">3</span> Digit Digit

# Leftmost Derivation

- Say, we have grammar

Integer → Digit | Integer Digit

Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

- 352 is an Integer?

Integer     → Integer Digit

       → Integer Digit Digit

       → Digit Digit Digit

       → 3 Digit Digit

       → 3 5 Digit

# Leftmost Derivation

- Say, we have grammar

  Integer → Digit | Integer Digit

  Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

- 352 is an Integer?

  Integer      → Integer Digit

             → Integer Digit Digit

             → Digit Digit Digit

             → 3 Digit Digit

             → 3 5 Digit → 352

What if I choose

Integer → Digit

# Leftmost Derivation

- Say, we have grammar

    Integer → Digit | Integer Digit

    Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

- 352 is an Integer?

    Integer       → Integer Digit

                    → Integer Digit Digit

                    → Digit Digit Digit

                    → 3 Digit Digit

                    → 3 5 Digit → 352

$$\text{Integer} \Rightarrow^* 352$$

$$A \Rightarrow^* \text{---} \circlearrowright$$

# Rightmost Derivation

- Say, we have grammar

Integer → Digit | Integer Digit

Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

# Rightmost Derivation

- Say, we have grammar

  Integer → Digit | Integer Digit

  Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9


- 352 is an Integer?

# Rightmost Derivation

- Say, we have grammar

  Integer → Digit | Integer Digit

  Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

- 352 is an Integer?

  Integer        → Integer Digit

# Rightmost Derivation

- Say, we have grammar

    Integer → Digit | Integer Digit

    Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

- 352 is an Integer?

    Integer → Integer Digit

    → Integer 2

# Rightmost Derivation

- Say, we have grammar

    Integer → Digit | Integer Digit

    Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9


- 352 is an Integer?

    Integer      → Integer Digit

                     → Integer 2

                     → Integer Digit 2

# Rightmost Derivation

- Say, we have grammar

  Integer → Digit | Integer Digit

  Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

- 352 is an Integer?

  Integer        → Integer Digit

                  → Integer 2

                  → Integer Digit 2

                  → Integer 5 2

# Rightmost Derivation

- Say, we have grammar

  Integer → Digit | Integer Digit

  Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

- 352 is an Integer?

  Integer          → Integer Digit

                   → Integer 2

                   → Integer Digit 2

                   → Integer 5 2

                   → Digit 5 2

# Rightmost Derivation

- Say, we have grammar

    Integer → Digit | Integer Digit

    Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9


- 352 is an Integer?

    Integer → Integer Digit

    → Integer 2

    → Integer Digit 2

    → Integer 5 2

    → <span style="color:red">Digit</span> 5 2

    → <span style="color:#3399ff">3</span> 5 2

# Rightmost Derivation

- Say, we have grammar

  Integer → Digit | Integer Digit

  Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

- 352 is an Integer?

  Integer      → Integer Digit

            → Integer 2

            → Integer Digit 2

            → Integer 5 2

            → Digit 5 2

            → 3 5 2

$$\text{Integer} \Rightarrow^* 352$$

# The Language of a Grammar

- If *G* is a CFG with alphabet **Σ** and start symbol **S**, then the *language of G* is the set

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

- That is, L(G) is the set of strings derivable from the start symbol.

- Note: ω must be in $\Sigma^*$, the set of strings made from terminals. String involving nonterminals aren't in the language.

# Context-Free Languages

- A language L is called a context-free language (or CFL) if there is a CFG G such that L = L(G).

- CFGs consist purely of production rules of the form $A \rightarrow \omega$. They do not have the regular expression operators * or ∪.

- However, we can convert regular expressions to CFGs as follows:

$$S \rightarrow a*b$$

aab

b

ab

- CFGs consist purely of production rules of the form $A \rightarrow \omega$. They do not have the regular expression operators * or ∪.

- However, we can convert regular expressions to CFGs as follows:

$$S \rightarrow Ab$$
$$A \rightarrow Aa \mid \varepsilon$$

$$S \rightarrow Ab$$
$$\varepsilon$$
$$\rightarrow b$$

$$S \rightarrow Ab$$
$$\rightarrow Aab$$
$$\rightarrow ab$$

$$S \rightarrow Ab$$
$$\rightarrow Aab$$
$$\rightarrow Aaab$$
$$\rightarrow aab$$

- CFGs consist purely of production rules of the form $A \to \omega$. They do not have the regular expression operators * or ∪.

- However, we can convert regular expressions to CFGs as follows:

$$S \to a(b \cup c^*) \qquad [a\,(b|c)]$$

$$\Rightarrow ab$$
$$\Rightarrow ac$$
$$\Rightarrow acc$$

- CFGs consist purely of production rules of the form $A \rightarrow \omega$. They do not have the regular expression operators * or ∪.

- However, we can convert regular expressions to CFGs as follows:

$$S \rightarrow aX$$
$$X \rightarrow b \mid C$$
$$C \rightarrow Cc \mid \varepsilon$$

$$S \rightarrow a\ X$$
$$\rightarrow a\ b$$

$$S \rightarrow a\ X$$
$$\rightarrow a\ c$$
$$\rightarrow a\ Cc$$
$$\rightarrow a\ c$$

# Regular Languages and CFLs

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

- **Proof Idea:** Use the construction from the previous slides to convert a regular expression for L into a CFG for L.

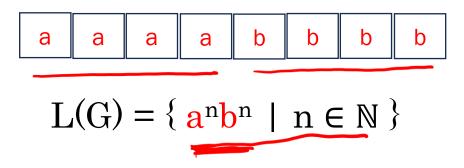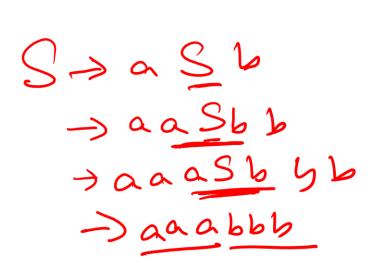- **Problem Set Exercise:** Instead, show how to convert a DFA/NFA into a CFG

- Consider the following CFG G:

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

- What strings can this generate?

$$S \rightarrow aSb$$
$$\rightarrow aaSbb$$
$$\rightarrow aaaSbbb$$
$$\rightarrow aaabbb$$

| a | a | a | a | b | b | b | b |

$$L(G) = \{\, a^n b^n \mid n \in \mathbb{N} \,\}$$

# Designing CFGs

- Like designing DFAs, NFAs, and regular expressions, designing CFGs is a craft.

- When thinking about CFGs:

  - Think recursively: Build up bigger structures from smaller ones.

  - Have a construction plan: Know in what order you will build up the string.

  - Store information in nonterminals: Have each nonterminal correspond to some useful piece of information.

- Let $\Sigma = \{a, b\}$ and let $L = \{w \in \Sigma^* \mid w$ is a palindrome $\}$
- We can design a CFG for L by thinking inductively:
- Base case: $\varepsilon$, $a$, and $b$ are palindromes.
- If $\omega$ is a palindrome, then $a\omega a$ and $b\omega b$ are palindromes.

$a\ b\ a\ a\ b\ a$

$$S \rightarrow \varepsilon \mid a \mid b \mid aSa \mid bSb$$

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

$a$
$b$
$\varepsilon$
$aSa$
$aa$

# Designing CFGs

- Let $\Sigma = \{\ (,\ )\ \}$ and let $L = \{w \in \Sigma^* \mid w$ is a string of balanced parentheses $\}$

- Some sample string in L

$$((()))$$

$$(())()$$

$$(()())(()())$$

$$(((((()))(())))$$

$$\varepsilon$$

$$()()$$

# Designing CFGs

- Let $\Sigma = \{$ (, ) $\}$ and let $L = \{w \in \Sigma^* \mid w$ is a string of balanced parentheses $\}$

- Let's think about this recursively.
  - **Base case:** the empty string is a string of balanced parentheses.
  - **Recursive step:** Look at the closing parenthesis that matches the first open parenthesis.

- Let $\Sigma = \{ (, ) \}$ and let $L = \{w \in \Sigma^* \mid w$ is a string of balanced parentheses $\}$

- Let's think about this recursively.
  - **Base case:** the empty string is a string of balanced parentheses.
  - **Recursive step:** Look at the closing parenthesis that matches the first open parenthesis. Removing the first parenthesis and the matching parenthesis forms two new strings of balanced parentheses.

$$S \to (S)S \mid \varepsilon$$

$$S \to (S)S$$

$$((S)S)S$$

$$((S)S)(S)S \Rightarrow (())()$$

- Let $\Sigma = \{\, a, b\, \}$ and let $L = \{w \in \Sigma^* \mid w$ has the same number of $a$'s and $b$'s $\}$

$$S \Rightarrow aSa$$

- Is this a CFG for L?

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

$$S \to aSb$$
$$\to absba$$
$$= aabb$$

- Can you derive the string abba?

# Designing CFGs: A Caveat

- When designing a CFG for a language, make sure that it
  - generates all the strings in the language and
  - never generates a string outside the language.
- The first of these can be tricky – make sure to test your grammars!

- Is the following grammar a CFG for the language $\{ a^n b^n \mid n \in \mathbb{N} \}$?
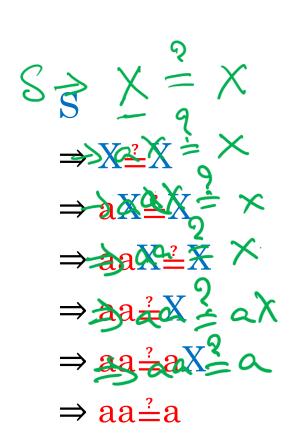
$$S \rightarrow aSb$$

- What strings can you derive?
  - Answer: None!

- What is the language of the grammar?
  - Answer: $\emptyset$

- When designing CFGs, make sure your recursion terminates!

**PennState**

- When designing CFGs, remember that each nonterminal can be expanded out independently of the others.

- Let $\Sigma = \{a, \stackrel{?}{=}\}$ and let $L = \{a^n \stackrel{?}{=} a^n \mid n \in \mathbb{N}\}$.

- Is the following a CFG for L?

$$S \to X \stackrel{?}{=} X$$
$$X \to aX \mid \varepsilon$$

$S$

$S \Rightarrow X \stackrel{?}{=} X$

$\Rightarrow X \stackrel{?}{=} X$

$\Rightarrow aX \stackrel{?}{=} X$

$\Rightarrow aaX \stackrel{?}{=} X$

$\Rightarrow aa \stackrel{?}{=} X \Rightarrow aX$

$\Rightarrow aa \stackrel{?}{=} aX \Rightarrow a$

$\Rightarrow aa \stackrel{?}{=} a$

- Let $\Sigma = \{a, \overset{?}{=}\}$ and let $L = \{a^n \overset{?}{=} a^n \mid n \in \mathbb{N}\}$. ✓

- To build a CFG for L, we need to be more clever with how we construct the string.

  - If we build the strings of a's independently of one another, then we can't enforce that they have the same length.

  - **Idea:** Build both strings of a's at the same time.

- Here's one possible grammar based on that idea:

$$S \to \overset{?}{=} \mid aSa$$

$S \to S \overset{\Rightarrow}{=} a\,Sa$

$\Rightarrow aSa$

$\Rightarrow aaSaa$

$\Rightarrow aaaSaaa$

$\Rightarrow aaa \overset{?}{=} aaa$

- Let $\Sigma = \{$void, int, double, name, (, ), ,, ;$\}$.

- Let's write a CFG for C-style function prototypes!

- Examples:
  - void name(int name, double name);
  - int name();
  - int name(double name);
  - int name(int, int name, int);
  - void name(void);

# Function Prototypes

- Here's one possible grammar:
  - $S \rightarrow$ Ret name (Args);
  - Ret $\rightarrow$ Type | void
  - Type $\rightarrow$ int | double
  - Args $\rightarrow$ ε | void | ArgList
  - ArgList $\rightarrow$ OneArg | ArgList, OneArg
  - OneArg $\rightarrow$ Type | Type name

- Fun question to think about: what changes would you need to make to support pointer types?

BLOCK → STMT
        | { STMTS }
STMTS → ε
        | STMT STMTS
STMT → EXPR;
        | if (EXPR) BLOCK
        | while (EXPR) BLOCK
        | do BLOCK while (EXPR);
        | BLOCK
        | …

EXPR → identifier
        | constant
        | EXPR + EXPR
        | EXPR – EXPR
        | EXPR * EXPR
        | …

# Reading and Exercises

## Reading

- Chapter: 2.2 (Michael Scott Book)

# References

Lecture Materials of CS 103, Stanford University