

CS419 Compilers Construction

A Simple One-Pass Compiler [Chapter 2]

Lecture 5

Edited by Dr. Ismail Hababeh

Adapted from slides by Dr. Mohammad Daoud

German Jordanian University

Originally from slides by Dr. Robert A. Engelen

Syntax Definition (Language Grammar) ²

- A grammar describes the **hierarchical structure** of **programming languages constructions**
- Example: in Java, the if-else structure has the form
if (expression) statement else statement

Let *expr* denotes an expression

stmt denotes a statement

→ means “can have the form”

Syntax Definition (Language Grammar) ³

The if-else **rule** structure can be expressed as:

$$stmt \rightarrow \mathbf{if} (expr) stmt \mathbf{else} stmt$$

This rule is called a **production**

Lexical elements like **if**, **else** and **()** are called **terminals**

Variables like **expr** and **stmt** are called **nonterminals**

Language Grammar

- A **context-free grammar** (CFG): a set of recursive rules used to generate patterns of strings.
- CFG consists of 4 components, or 4-tuple
 - A set of *non-terminals* (each nonterminal represents a set of strings of terminals)
 - A set of (*terminal* symbols)
 - A set of *productions* based on terminals and non-terminals.
 - A designated *start symbol*

Language Grammar

- **CFG** = *nonterminals* , *terminals*, *productions*, *start symbol*
- A context-free grammar G is defined as:
$$G = (V, M, P, S)$$
- V : nonterminals (variables) // *written in italic font*
- M : terminals // **written in bold**
- P : productions or rules
- S : start variable

Context-free Grammar - Example

Context-free grammar for simple expressions such as: $9 - 5 + 2$, $3 - 1$, or 7 (or list of digits separated by plus or minus signs) is defined as the 4-tuple:

$$G = \langle \{list, digit\}, \{+, -, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}, P, list \rangle$$

with productions $P =$

$list \rightarrow list + digit$ *// first production*

$list \rightarrow list - digit$ *// second production*

$list \rightarrow digit$ *// Third production*

$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

Alternative way to write the Context-free grammar - Example

$$G = \langle \{list, digit\}, \{+, -, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}, P, list \rangle$$

with productions $P =$

$$list \rightarrow list + digit \mid list - digit \mid digit$$

$$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

Note: \mid means OR

Grammar Derivation

- Given a CFG, we can determine the **set of all *strings*** (sequences of tokens) **generated by the grammar** using *derivation*
 - We begin with the start symbol
 - In each step, we replace one nonterminal in the current *sentential form (production)* with one of the right-hand sides of a production for that nonterminal
- The derivation can be used to prove that a string belongs to the grammar's language.

Grammar Derivation - Example

Derivation of the string $9 - 5 + 2$ from the *list* productions in the pervious example

\underline{list}
 $\Rightarrow \underline{list} + digit \quad // \text{ using the first production}$
sentential form $\Rightarrow \underline{list} - digit + digit \quad // \text{ using the second production}$
 $\Rightarrow \underline{digit} - digit + digit \quad // \text{ using the third production}$
 $\Rightarrow 9 - \underline{digit} + digit$
 $\Rightarrow 9 - 5 + \underline{digit}$
 $\Rightarrow 9 - 5 + 2$

This example is *leftmost derivation*, because we replaced the leftmost nonterminal (underlined) in each step

Parse Trees

- **Parsing** = *process of determining if a string of tokens can be generated by a certain grammar.*
- A **parse tree** shows how the **start symbol of a grammar derives a string in the language**
- Given a context-free grammar, the parse tree is a tree with the following properties:
 - The **root** of the tree is labeled by the **start symbol**
 - Each **leaf** of the tree is labeled by a **terminal** (token) or ϵ (ϵ denotes the *empty string*)
 - Each **interior node** is labeled by a **nonterminal**
 - If $A \rightarrow X_1 X_2 \dots X_n$ is a production, then node A has children X_1, X_2, \dots, X_n where X_i is a terminal, nonterminal, or ϵ

String Parsing - Example

Given the following Grammar:

$$G = \langle \{list, digit\}, \{+, -, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}, P, list \rangle$$

Draw the parse tree of the string **9-5+2** according to the grammar G

list

$\Rightarrow \underline{list} + digit$

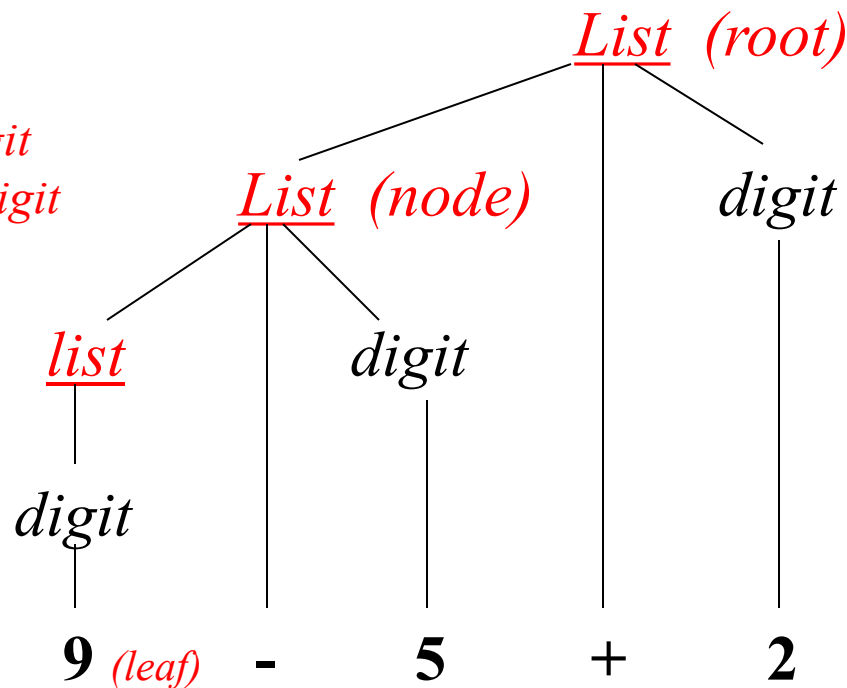
$\Rightarrow \underline{list} - digit + digit$

$\Rightarrow \underline{digit} - digit + digit$

$\Rightarrow \underline{9} - \underline{digit} + digit$

$\Rightarrow \underline{9} - \underline{5} + \underline{digit}$

$\Rightarrow \underline{9} - \underline{5} + \underline{2}$



The process of finding a parse tree for a given string of terminals is called *parsing that string*

The sequence of leaf's is the *yield* of the parse tree

String Parsing Ambiguity

Consider the following context-free grammar:

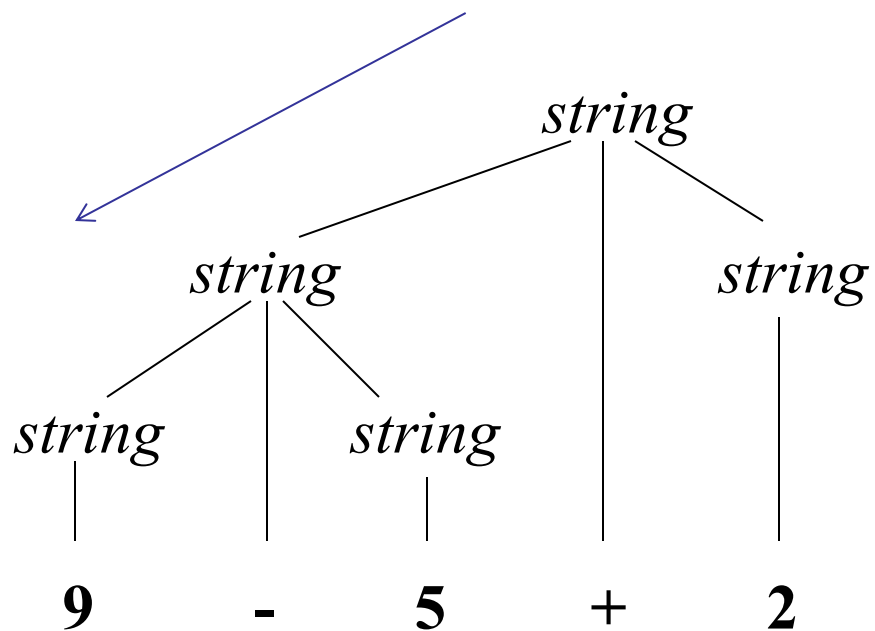
$$G = \langle \{string\}, \{+, -, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}, P, string \rangle$$

with productions $P =$

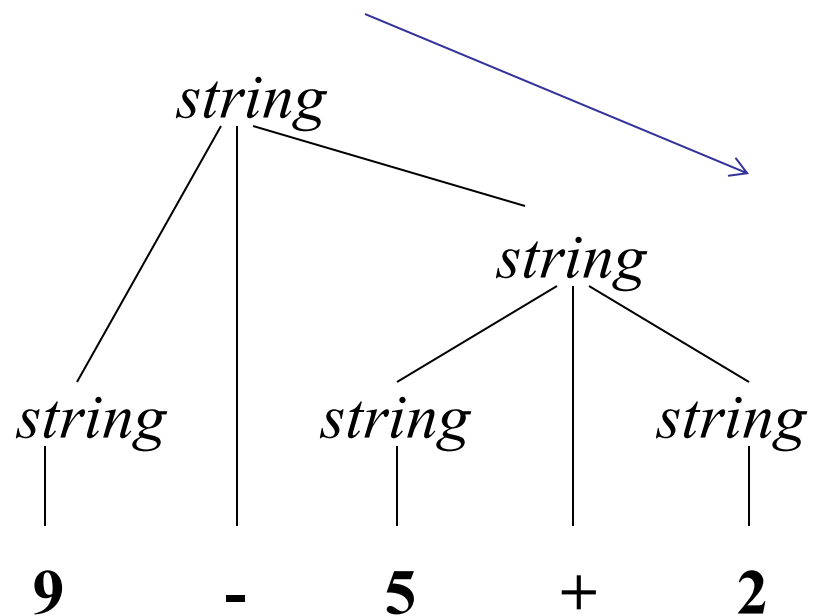
$$string \rightarrow string + string \mid string - string \mid 0 \mid 1 \mid \dots \mid 9$$

This grammar is *ambiguous*, because more than one parse tree generates the string 9-5+2

String Parsing Ambiguity - Example



$(9-5) + 2$



$9 - (5+2)$

Associativity of Operators

❖ *Left-associative operators* have *left-recursive* productions

$+$, $-$, $$, $/$ are left-associative*

Example: String $\mathbf{a+b+c}$ has the same meaning as $\mathbf{(a+b)+c}$

❖ *Right-associative operators* have *right-recursive* productions

$=$ is right-associative

Example: String $\mathbf{a=b=c}$ has the same meaning as $\mathbf{a=(b=c)}$

Precedence of Operators

Operators with higher precedence “bind more tightly”
Given the following productions of a grammar G:

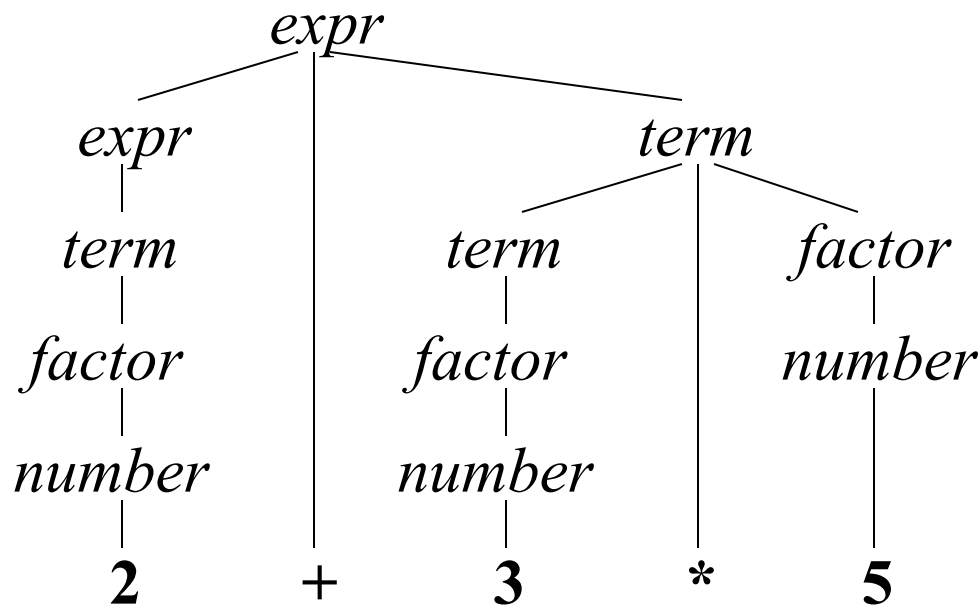
$expr \rightarrow expr + term \mid term$

$term \rightarrow term * factor \mid factor$

$factor \rightarrow number \mid (expr)$

$number \rightarrow 0 \mid \dots \mid 9$

String **2+3*5** has the same meaning as **2+(3*5)**



Grammar for a Subset of Java Statements

$$\begin{aligned} stmt &\rightarrow \mathbf{id} = expr; \\ &\quad | \mathbf{if} (expr) stmt \\ &\quad | \mathbf{if} (expr) stmt \mathbf{else} stmt \\ &\quad | \mathbf{while} (expr) stmt \\ &\quad | \mathbf{do} stmt \mathbf{while} (expr); \\ &\quad | \{\mathbf{stmts}\} \\ stmts &\rightarrow stmts stmt \\ &\quad | \varepsilon \end{aligned}$$