
Chapter 2

A Simple One – Pass Compiler

The Entire Compilation Process

- Grammars for Syntax Definition
- Syntax-Directed Translation
- Parsing - Top Down & Predictive
- Pulling Together the Pieces
- The Lexical Analysis Process
- Symbol Table Considerations
- A Brief Look at Code Generation
- Concluding Remarks/Looking Ahead

Overview

Programming Language can be defined by describing

1. The syntax of the language
 - 1. What its program looks like*
 - 2. We use CFG or BNF (Backus Naur Form)*
2. The semantics of the language
 - 1. What its program mean*
 - 2. Difficult to describe*
 - 3. Use informal descriptions and suggestive examples*

Grammars for Syntax Definition

- A Context-free Grammar (CFG) Is Utilized to Describe the Syntactic Structure of a Language
- A CFG Is Characterized By:
 1. A Set of Tokens or Terminal Symbols
 2. A Set of Non-terminals
 3. A Set of Production Rules
Each Rule Has the Form
$$NT \rightarrow \{T, NT\}^*$$
 4. A Non-terminal Designated As the Start Symbol

Grammars for Syntax Definition

Example CFG

$$\textit{list} \rightarrow \textit{list} + \textit{digit}$$
$$\textit{list} \rightarrow \textit{list} - \textit{digit}$$
$$\textit{list} \rightarrow \textit{digit}$$
$$\textit{digit} \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

(the “ \mid ” means OR)

(So we could have written

$$\textit{list} \rightarrow \textit{list} + \textit{digit} \mid \textit{list} - \textit{digit} \mid \textit{digit})$$



Information

- ✓ A string of tokens is a sequence of zero or more tokens.
- ✓ The string containing with zero tokens, written as ϵ , is called empty string.
- ✓ A grammar derives strings by beginning with the start symbol and repeatedly replacing the non terminal by the right side of a production for that non terminal.
- ✓ The token strings that can be derived from the start symbol form the language defined by the grammar.

Grammars are Used to Derive Strings:

Using the CFG defined on the earlier slide, we can derive the string: $9 - 5 + 2$ as follows:

$list \rightarrow list + digit$

P1 : $list \rightarrow list + digit$

$\rightarrow list - digit + digit$

P2 : $list \rightarrow list - digit$

$\rightarrow digit - digit + digit$

P3 : $list \rightarrow digit$

$\rightarrow 9 - digit + digit$

P4 : $digit \rightarrow 9$

$\rightarrow 9 - 5 + digit$

P4 : $digit \rightarrow 5$

$\rightarrow 9 - 5 + 2$

P4 : $digit \rightarrow 2$

Grammars are Used to Derive Strings:

This derivation could also be represented via a Parse Tree
(parents on left, children on right)

$list \rightarrow list + digit$

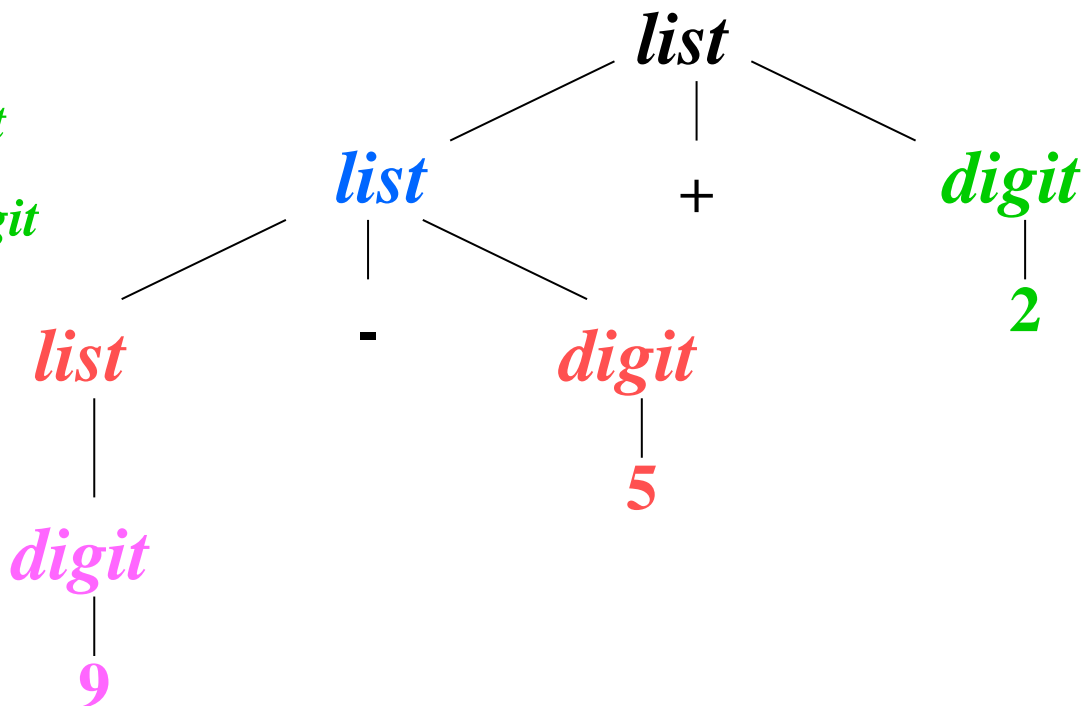
$\rightarrow list - digit + digit$

$\rightarrow digit - digit + digit$

$\rightarrow 9 - digit + digit$

$\rightarrow 9 - 5 + digit$

$\rightarrow 9 - 5 + 2$



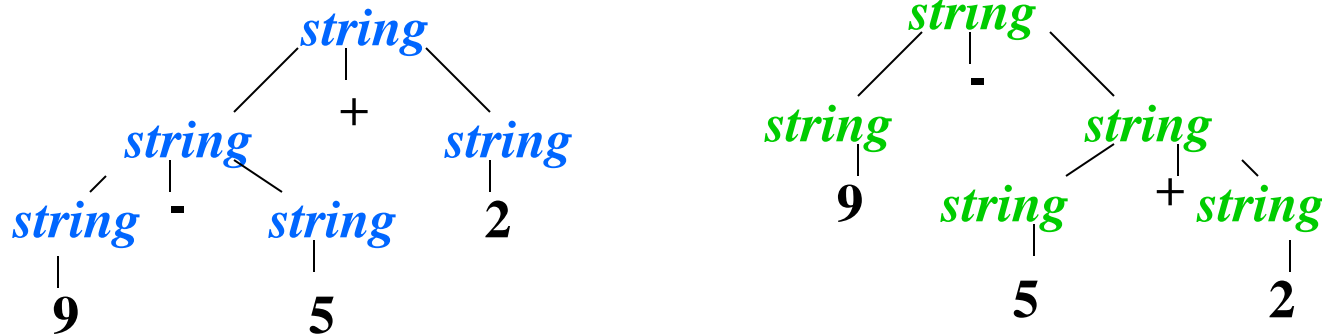
Defining a Parse Tree

- A parse tree pictorially shows **how** the start symbol of a grammar derives a string in the language.
- More Formally, a Parse Tree for a CFG Has the Following Properties:
 - Root Is Labeled With the **Start Symbol**
 - Leaf Node Is a Token or ϵ
 - Interior Node Is a **Non-Terminal**
 - If $A \rightarrow x_1x_2\dots x_n$, Then A Is an Interior; $x_1x_2\dots x_n$ Are Children of A and May Be **Non-Terminals** or **Tokens**

Other Important Concepts

Ambiguity

Two derivations (Parse Trees) for the same token string.



Grammar:

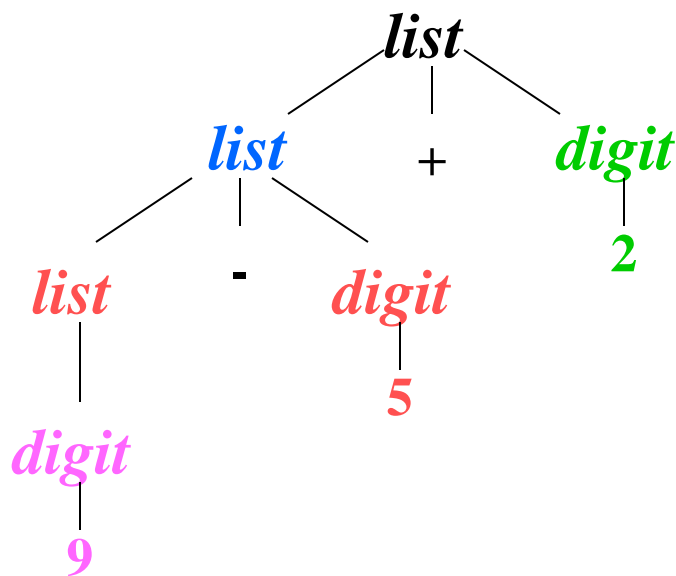
$string \rightarrow string + string / string - string \mid 0 \mid 1 \mid \dots \mid 9$

Why is this a Problem ?

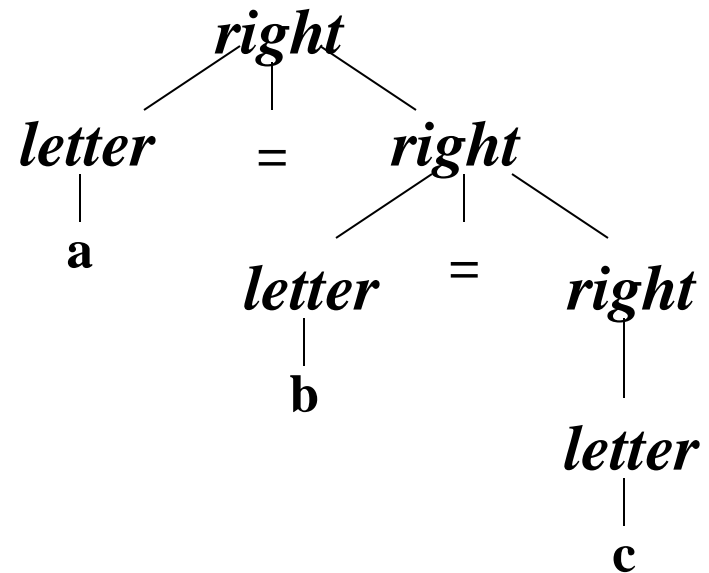
Other Important Concepts

Associativity of Operators

Left vs. Right



$list \rightarrow list + digit /$
 $\quad \quad \quad / list - digit / digit$
 $digit \rightarrow 0 | 1 | 2 | \dots | 9$



$right \rightarrow letter = right / letter$
 $letter \rightarrow a | b | c | \dots | z$

Embedding Associativity

- The language of arithmetic expressions with + -
 - (ambiguous) grammar that does not enforce associativity
$$\textit{string} \rightarrow \textit{string} + \textit{string} / \textit{string} - \textit{string} \mid 0 \mid 1 \mid \dots \mid 9$$
 - non-ambiguous grammar enforcing left associativity (parse tree will grow to the left)
$$\textit{string} \rightarrow \textit{string} + \textit{digit} / \textit{string} - \textit{digit} / \textit{digit}$$
$$\textit{digit} \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9$$
 - non-ambiguous grammar enforcing right associativity (parse tree will grow to the right)
$$\textit{string} \rightarrow \textit{digit} + \textit{string} / \textit{digit} - \textit{string} / \textit{digit}$$
$$\textit{digit} \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9$$

Other Important Concepts

Operator Precedence

What does

$9 + 5 * 2$

mean?

Typically

$\left\{ \begin{array}{l} () \\ * / \\ + - \end{array} \right.$ is precedence order

This can be
incorporated
into a grammar
via rules:

$expr \rightarrow expr + term \mid expr - term \mid term$

$term \rightarrow term * factor \mid term / factor \mid factor$

$factor \rightarrow \text{digit} \mid (expr)$

$\text{digit} \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid \dots \mid 9$

Precedence Achieved by:

expr & term for each precedence level

Rules for each are **left recursive** or **associate to the left**

Syntax for Statements

$stmt \rightarrow id := expr$
| if $expr$ then $stmt$
| if $expr$ then $stmt$ else $stmt$
| while $expr$ do $stmt$
| begin opt_stmts end

Ambiguous Grammar?

Syntax-Directed Translation

- Associate Attributes With Grammar Rules and Translate as Parsing occurs
- The translation will follow the parse tree structure (and as a result the structure and form of the parse tree will affect the translation).
- First example: Inductive Translation.
- **Infix** to **Postfix** Notation Translation for Expressions
 - Translation defined inductively as: **Postfix(E)** where E is an Expression.

Rules

1. *If E is a variable or constant then* **Postfix(E) = E**
2. *If E is E1 op E2 then* **Postfix(E)**
= Postfix(E1 op E2) = Postfix(E1) Postfix(E2) op
3. *If E is (E1) then* **Postfix(E) = Postfix(E1)**

Examples

Postfix((9 - 5) + 2)
 = **Postfix((9 - 5)) Postfix(2) +**
 = **Postfix(9 - 5) Postfix(2) +**
 = **Postfix(9) Postfix(5) - Postfix(2) +**
 = **9 5 - 2 +**

Postfix(9 - (5 + 2))
 = **Postfix(9) Postfix((5 + 2)) -**
 = **Postfix(9) Postfix(5 + 2) -**
 = **Postfix(9) Postfix(5) Postfix(2) + -**
 = **9 5 2 + -**

Syntax-Directed Definition

- Each Production Has a Set of **Semantic Rules**
- Each Grammar Symbol Has a Set of **Attributes**
- For the Following Example, String Attribute “*t*” is Associated With Each Grammar Symbol

$$\textit{expr} \rightarrow \textit{expr} - \textit{term} \mid \textit{expr} + \textit{term} \mid \textit{term}$$

$$\textit{term} \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid \dots \mid 9$$

- recall: What is a Derivation for $9 + 5 - 2$?

$$\begin{aligned} \textit{list} &\rightarrow \textit{list} - \textit{digit} \rightarrow \textit{list} + \textit{digit} - \textit{digit} \rightarrow \textit{digit} + \textit{digit} - \textit{digit} \\ &\rightarrow 9 + \textit{digit} - \textit{digit} \rightarrow 9 + 5 - \textit{digit} \rightarrow 9 + 5 - 2 \end{aligned}$$

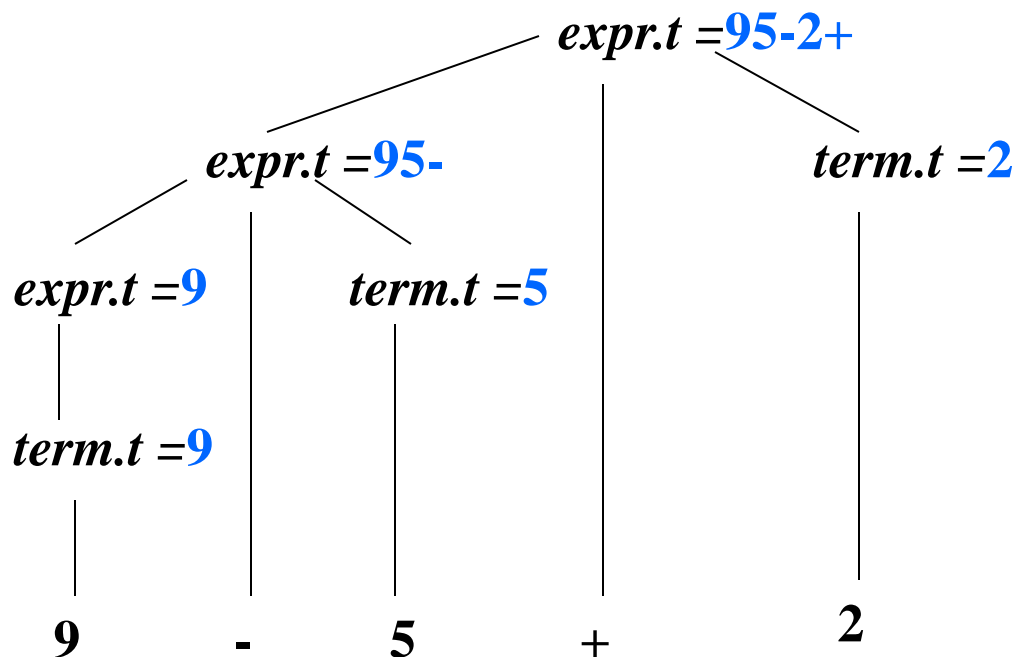
Syntax-Directed Definition (2)

- Each Production Rule of the CFG Has a Semantic Rule

Production	Semantic Rule
$expr \rightarrow expr + term$	$expr.t := expr.t \parallel term.t \parallel '+'$
$expr \rightarrow expr - term$	$expr.t := expr.t \parallel term.t \parallel '-'$
$expr \rightarrow term$	$expr.t := term.t$
$term \rightarrow 0$	$term.t := '0'$
$term \rightarrow 1$	$term.t := '1'$
....
$term \rightarrow 9$	$term.t := '9'$

- Note:** Semantic Rules for $expr$ define t as a “synthesized attribute” i.e., the various copies of t obtain their values from “children t ’s”

Semantic Rules are Embedded in Parse Tree



- It starts at the root and recursively visits the children of each node in left-to-right order
- The semantic rules at a given node are evaluated once all descendants of that node have been visited.
- A parse tree showing all the attribute values at each node is called annotated parse tree.

Translation Schemes

Embedded Semantic Actions into the right sides of the productions.

A translation scheme is like a syntax-directed definition except the order of evaluation of the semantic rules is explicitly shown.

$expr \rightarrow expr + term \quad \{print('+')\}$
 $\rightarrow expr - term \quad \{print('-')\}$
 $\rightarrow term$
 $term \rightarrow 0 \quad \{print('0')\}$
 $term \rightarrow 1 \quad \{print('1')\}$
 \dots
 $term \rightarrow 9 \quad \{print('9')\}$

