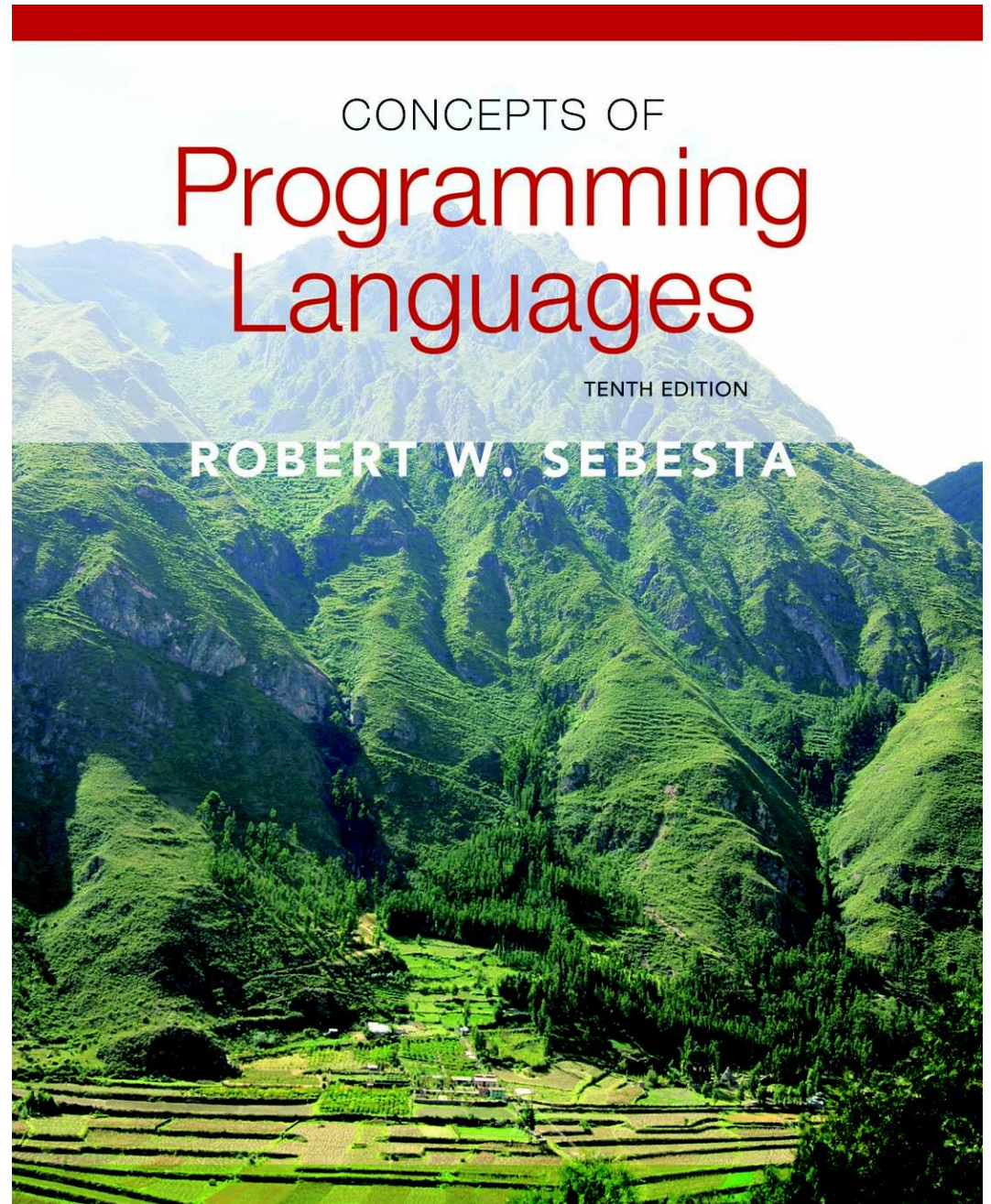


Chapter 3

Describing Syntax and Semantics



Chapter 3 Topics

- Introduction
- The General Problem of Describing Syntax
- Formal Methods of Describing Syntax
- Attribute Grammars
- Describing the Meanings of Programs:
Dynamic Semantics

Introduction

- **Syntax:** the form or structure of the expressions, statements, and program units
- **Semantics:** the meaning of the expressions, statements, and program units
- Syntax and semantics provide a language's definition
 - Users of a language definition
 - Other language designers
 - Implementers (interpret, compile)
 - Programmers (the users of the language)

The General Problem of Describing Syntax: Terminology

- A *sentence* is a string of characters over some alphabet
- A *language* is a set of sentences
- A *lexeme* is the lowest level syntactic unit of a language (e.g., `*`, `sum`, `begin`)
- A *token* is a category of lexemes (e.g., identifier)

The General Problem of Describing Syntax: Terminology

- `index = 3 * count + 5;`

Lexeme	Token
<code>index, count</code>	identifier
<code>=, *, +</code>	Operator
<code>3, 5</code>	int_literal

Formal Definition of Languages

- How can you define a PL?
- **Recognizers**
 - A **recognition device** reads input strings over the alphabet of the language and decides whether the input strings **belong to the language**
 - Example: syntax analysis part of a compiler
 - Detailed discussion of syntax analysis (**parser**) appears in Chapter 4
 - Limited usage because trial and error mode; generators preferred
- **Generators**
 - A device that **generates sentences of a language**
 - One can determine if the syntax of a particular sentence is **syntactically correct** by **comparing** it to the **structure of the generator**

BNF and Context-Free Grammars

- Context-Free Grammars
 - Developed by Noam Chomsky (linguist) in the mid-1950s
 - Language generators, meant to describe the **syntax of natural languages**
 - Defines a class of languages called context-free languages
- Backus-Naur Form (1959)
 - Invented by John Backus to **describe the syntax of Algol 58**
 - **BNF is equivalent to context-free grammars**

Quiz

- What is the difference between lexeme and token?

BNF Fundamentals

- What is a **metalanguage**?
 - A language that is used to describe other languages.
- **BNF is a metalanguage for programming languages.**

BNF Fundamentals

- In BNF, **abstractions** are used to represent classes of syntactic structures—they act like syntactic variables (also called *nonterminal symbols*, or just *terminals*)
- *Terminals* are
 - lexemes
 - or tokens
- A rule has a left-hand side (LHS), which is a **nonterminal**, and a right-hand side (RHS), which is a string of **terminals** and/or **nonterminals**

BNF Fundamentals (continued)

- Nonterminals are often enclosed in angle brackets
 - Examples of BNF rules:
`<if_stmt> → if <logic_expr> then <stmt>`
`<ident_list> → identifier | identifier, <ident_list>`
- Grammar:
 - a finite non-empty set of rules
- A *start symbol* is a special element of the nonterminals of a grammar

BNF Rules

- An abstraction (or nonterminal symbol) can have more than one RHS

```
<stmt> → <single_stmt>  
        | begin <stmt_list> end
```

Describing Lists

- **Syntactic lists are described using recursion**

$$\begin{aligned} \langle \text{ident_list} \rangle &\rightarrow \text{ident} \\ &\quad | \text{ ident, } \langle \text{ident_list} \rangle \end{aligned}$$

- **A derivation is a repeated application of rules,**
 - starting with the **start symbol** and ending with a **sentence** (all terminal symbols)

A Grammar for a Small Language

$\langle \text{program} \rangle \rightarrow \mathbf{begin} \langle \text{stmts} \rangle \mathbf{end}$

$\langle \text{stmts} \rangle \rightarrow \langle \text{stmt} \rangle \mid \langle \text{stmt} \rangle ; \langle \text{stmts} \rangle$

$\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$

$\langle \text{var} \rangle \rightarrow a \mid b \mid c \mid d$

$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle + \langle \text{term} \rangle \mid \langle \text{term} \rangle - \langle \text{term} \rangle$

$\langle \text{term} \rangle \rightarrow \langle \text{var} \rangle \mid \text{const}$

An Example Derivation

$\langle \text{program} \rangle \Rightarrow \langle \text{stmts} \rangle$
 $\Rightarrow \langle \text{stmt} \rangle$
 $\Rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$
 $\Rightarrow a = \langle \text{expr} \rangle$
 $\Rightarrow a = \langle \text{term} \rangle + \langle \text{term} \rangle$
 $\Rightarrow a = \langle \text{var} \rangle + \langle \text{term} \rangle$
 $\Rightarrow a = b + \langle \text{term} \rangle$
 $\Rightarrow a = b + \text{const}$

$\langle \text{program} \rangle \rightarrow \langle \text{stmts} \rangle$
 $\langle \text{stmts} \rangle \rightarrow \langle \text{stmt} \rangle \mid \langle \text{stmt} \rangle ; \langle \text{stmts} \rangle$
 $\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$
 $\langle \text{var} \rangle \rightarrow a \mid b \mid c \mid d$
 $\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle + \langle \text{term} \rangle \mid \langle \text{term} \rangle - \langle \text{term} \rangle$
 $\langle \text{term} \rangle \rightarrow \langle \text{var} \rangle \mid \text{const}$

Derivations

- Every string of symbols in a derivation is a *sentential form*
- A *sentence* is a sentential form that has only terminal symbols
- *Leftmost derivation*
 - leftmost nonterminal in each sentential form is the one that is expanded

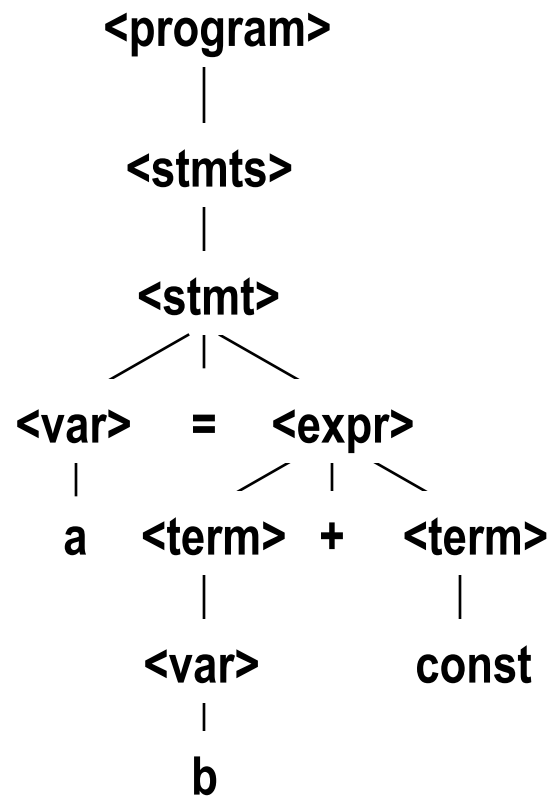
An Example Derivation

$\langle \text{program} \rangle \Rightarrow \langle \text{stmts} \rangle$
 $\Rightarrow \langle \text{stmt} \rangle$
 $\Rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$
 $\Rightarrow a = \langle \text{expr} \rangle$
 $\Rightarrow a = \langle \text{term} \rangle + \langle \text{term} \rangle$
 $\Rightarrow a = \langle \text{var} \rangle + \langle \text{term} \rangle$
 $\Rightarrow a = b + \langle \text{term} \rangle$
 $\Rightarrow a = b + \text{const}$

$\langle \text{program} \rangle \rightarrow \langle \text{stmts} \rangle$
 $\langle \text{stmts} \rangle \rightarrow \langle \text{stmt} \rangle \mid \langle \text{stmt} \rangle ; \langle \text{stmts} \rangle$
 $\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$
 $\langle \text{var} \rangle \rightarrow a \mid b \mid c \mid d$
 $\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle + \langle \text{term} \rangle \mid \langle \text{term} \rangle - \langle \text{term} \rangle$
 $\langle \text{term} \rangle \rightarrow \langle \text{var} \rangle \mid \text{const}$

Parse Tree

- A hierarchical representation of a derivation



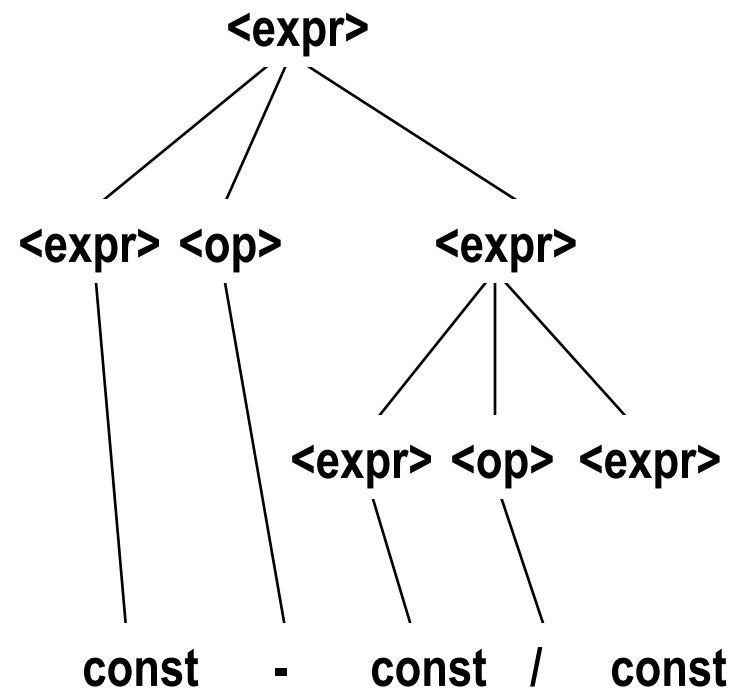
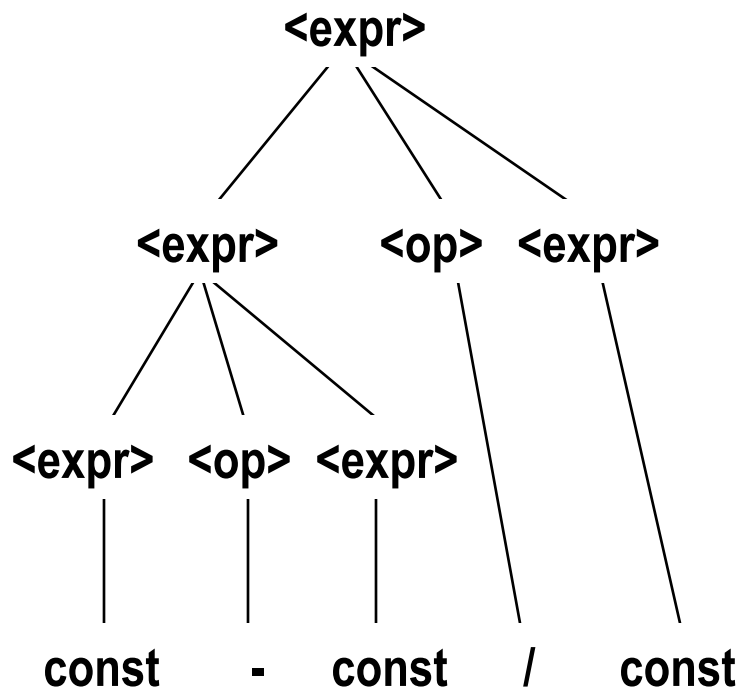
Ambiguity in Grammars

- A grammar is *ambiguous*
 - if and only if it generates a sentential form that has two or more distinct parse trees

An Ambiguous Expression Grammar

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \mid \text{const}$

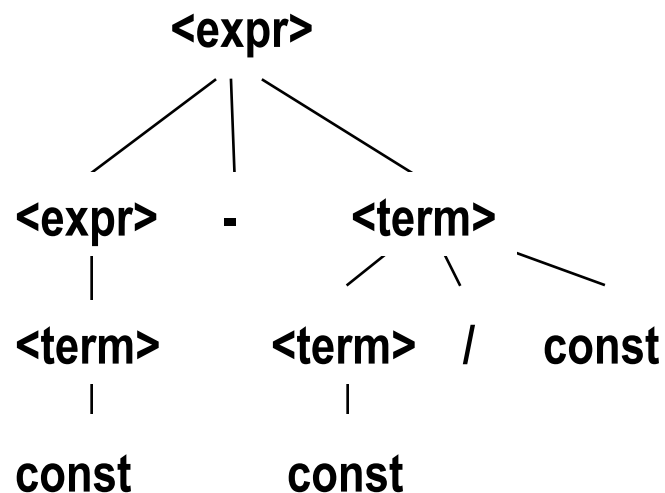
$\langle \text{op} \rangle \rightarrow / \mid -$



An Unambiguous Expression Grammar

- If parse tree to indicate precedence levels of the operators

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle - \langle \text{term} \rangle \mid \langle \text{term} \rangle$
 $\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle / \text{const} \mid \text{const}$



Quiz

- Prove that the following grammar is ambiguous:

$\langle A \rangle \rightarrow \langle C \rangle$

- $\langle C \rangle \rightarrow \langle C \rangle \langle \text{op} \rangle \langle C \rangle \mid \langle \text{id} \rangle$

- $\langle \text{op} \rangle \rightarrow * \mid / \mid + \mid -$

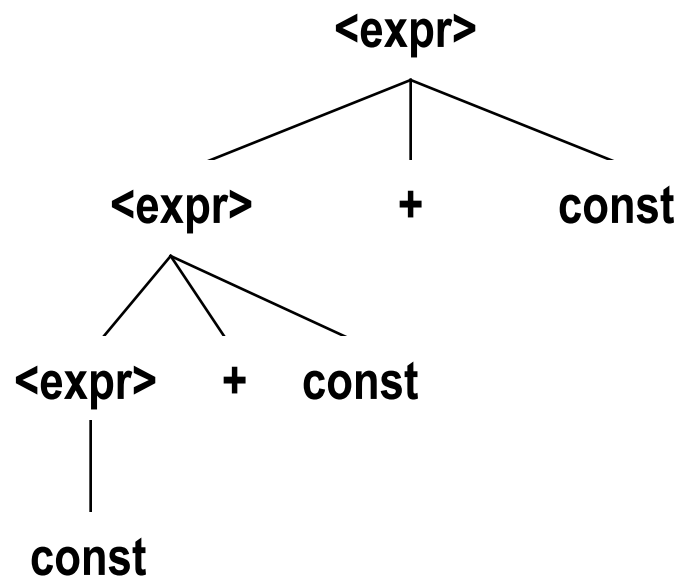
- $\langle \text{id} \rangle \rightarrow a \mid b \mid c$

Associativity of Operators

- Operator associativity can also be indicated by a grammar

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid \text{const}$ (ambiguous)

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \text{const} \mid \text{const}$ (unambiguous)



Associativity of Operators

- Integer addition

$$(3+4) + 5 = 3 + (4+5)$$

- Floating point addition

$$(10^7 + 1) + \dots + 1 \stackrel{?}{=} 10^7 + (1 + \dots + 1)$$

$$- 1.000000 \times 10^7 \stackrel{?}{=} 1.000001 \times 10^7$$

Quiz

- Rewrite the following BNF to give / precedence over * .

$\langle \text{assign} \rangle \rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$

$\langle \text{id} \rangle \rightarrow A \mid B \mid C$

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle / \langle \text{term} \rangle \mid \langle \text{term} \rangle$

$\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle * \langle \text{factor} \rangle \mid \langle \text{factor} \rangle$

$\langle \text{factor} \rangle \rightarrow (\langle \text{expr} \rangle) \mid \langle \text{id} \rangle$

$B = A / (A / B) * C$

Quiz

- Using the grammar above, show a parse tree and a leftmost derivation for the following statement:

$$B = A / (A / B) * C$$

Quiz

- Describe, in English, the language defined by the following grammar:

$$S \rightarrow a S \mid S b \mid a \mid b$$

Extended BNF

- **Optional parts** are placed in brackets []

`<proc_call> -> ident [(<expr_list>)]`

- **Alternative parts:**

- parentheses

`<term> → <term> (+|-) const`

- **Repetitions**

- (0 or more) are placed inside braces { }

`<ident> → letter {letter|digit}`

BNF and EBNF

- BNF

```
<expr> → <expr> + <term>
        | <expr> - <term>
        | <term>
<term>  → <term> * <factor>
        | <term> / <factor>
        | <factor>
```

- EBNF

```
<expr> → <term> { (+ | -) <term> }
<term> → <factor> { (* | /) <factor> }
```

Static Semantics

- Nothing to do with the meaning
- Context-free grammars (CFGs) cannot describe all of the syntax of programming languages
 - (e.g., variables must be declared before they are used)