

Sebesta: Concepts of Programming Languages Chapter 3

Describing Syntax and Semantics

CS 4100
March 2, 2011

Some Chapter 3 Topics

- Introduction
- The General Problem of Describing Syntax
- Formal Methods of Describing Syntax

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
 - 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)

Formal Definition of Languages

- **Recognizers**
 - A recognition device reads input strings of the language and decides whether the input strings belong to the language
 - Example: syntax analysis part of a compiler
- **Generators**
 - A device that generates sentences of a language
 - One can determine if the syntax of a particular sentence is correct by comparing it to the structure of the generator

Formal Methods of Describing Syntax

- **Backus-Naur Form and Context-Free Grammars**
 - Most widely known method for describing programming language syntax
- **Extended BNF**
 - Improves readability and writability of BNF
- **Grammars and Recognizers**

BNF and Context-Free Grammars

- **Context-Free Grammars**
 - Developed by Noam Chomsky in the mid-1950s
 - Language generators, meant to describe the syntax of natural languages
 - Define a class of languages called context-free languages

Chomsky Hierarchy

- **Type-0** Recursively enumerable
 - Turing machine
 - any string of non-terminals ::= any other string of non-terminals and terminals
- **Type-1** Context-sensitive
 - Linear-bounded non-deterministic Turing machine
 - any string of non-terminals ::= any other string of non-terminals and terminals
- **Type-2** Context-free
 - Non-deterministic pushdown automaton
 - $\langle nt \rangle ::=$ any string of terminal and non-terminal symbols
- **Type-3** Regular
 - Finite state automaton
 - $\langle nt \rangle ::= k \langle nt \rangle$ or $\langle nt \rangle ::= k$

Backus-Naur Form (BNF)

- Backus-Naur Form (1959)
 - Invented by John Backus to describe Algol 58
 - BNF is equivalent to context-free grammars
 - BNF is a *metalanguage* used to describe another language
 - In BNF, abstractions are used to represent classes of syntactic structures--they act like syntactic variables (also called *nonterminal symbols*)

BNF Fundamentals

- Non-terminals: BNF abstractions
- Terminals: lexemes and tokens
- Grammar: a collection of rules
 - Examples of BNF rules:
 - `<ident_list> → identifier | identifier, <ident_list>`
 - `<if_stmt> → if <logic_expr> then <stmt>`

BNF Rules

- A rule has a left-hand side (LHS) and a right-hand side (RHS), and consists of *terminal* and *nonterminal* symbols
- A grammar is a finite nonempty set of 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
 - `<ident_list> → ident`
 - `| ident, <ident_list>`
- A derivation is a repeated application of rules, starting with the start symbol and ending with a sentence (all terminal symbols)

An Example Grammar

```

<program> → <stmts>
<stmts> → <stmt> | <stmt> ; <stmts>
<stmt> → <var> = <expr>
<expr> → <term> + <term> | <term> - <term>
<term> → <var> | const
<var> → a | b | c | d
    
```

An Example Derivation

```

<program> => <stmts>
=> <stmt>
=> <var> = <expr>
=> a = <expr>
=> a = <term> + <term>
=> a = <var> + <term>
=> a = b + <term>
=> a = b + const
    
```

```

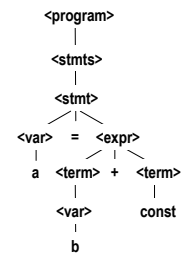
<program> → <stmts>
<stmts> → <stmt> | <stmt> ; <stmts>
<stmt> → <var> = <expr>
<expr> → <term> + <term> | <term> - <term>
<term> → <var> | const
<var> → a | b | c | d
    
```

Derivation

- Every string of symbols in the derivation is a sentential form
- A sentence is a sentential form that has only terminal symbols
- A leftmost derivation is one in which the leftmost nonterminal in each sentential form is the one that is expanded
- A derivation may be neither leftmost nor rightmost

Parse Tree

- A hierarchical representation of a derivation

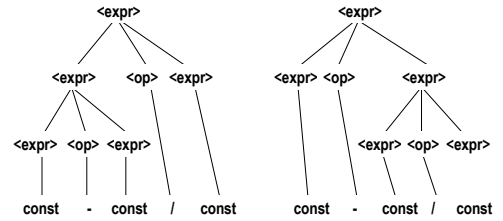


Ambiguity in Grammars

- 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

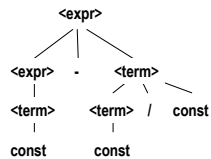
$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \mid \text{const}$
 $\langle \text{op} \rangle \rightarrow / \mid -$



An Unambiguous Expression Grammar

- If we use the parse tree to indicate precedence levels of the operators, we cannot have ambiguity

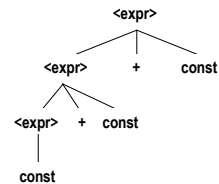
$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle - \langle \text{term} \rangle \mid \langle \text{term} \rangle$
 $\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle / \text{const} \mid \text{const}$



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)



Extended BNF

- Optional parts are placed in brackets []
`<proc_call> -> ident [(<expr_list>)]`
- Alternative parts of RHSs are placed inside parentheses and separated via vertical bars
`<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> }
```