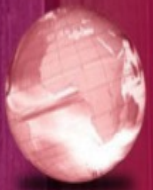


GLOBAL
EDITION



Chapter 3

Syntax and Semantics

Concepts of Programming Languages

ELEVENTH EDITION

Robert W. Sebesta

Introduction

- Syntax
 - General Problem of Describing Syntax
 - Formal Methods of Describing Syntax
- Semantics
 - Attribute Grammars
 - Formal Methods of Describing Semantics

Introduction

- A concise, yet, understandable description of a programming language.
- A standard, formal way of authoring this description.
- Making the description hospitable to a verity of users (programmers)
- To enable a programmer to encode software from information provided in a reference manual.

-
- **What is Syntax, Semantics?**
 - The ***syntax*** of a programming language is the form of its expressions, statements, and program units: How things appear
 - Its ***semantics*** is the meaning of those expressions, statements, and program units: What things mean
 - **Example from natural language**
 - A valid English statement is correct both syntactically and semantically: e.g. The wind subsided.
 - The statement follows a syntactic rule:
 - ARTICLE SUBJECT VERB, and a meaning is “attached” to it.
 - Consider the sentence: **Colorless green ideas sleep furiously.**
 - This is correct syntactically, but nonsensical semantically.
 - We evaluate the “syntactic correctness” of a language based on it’s grammar, which all (almost) English speakers agree upon.

Syntax

- Example from a Programming Language
 - The **syntax** (grammar) of a C **if** statement is:
`if (<expr>) <statement>`
 - **Semantics** for this statement: **if** the current value of the expression is true, the embedded statement is selected for execution.
 - In a well-designed programming language syntax should suggest semantics, i.e., the form of a statement should strongly suggest its logical meaning.

Describing Syntax

- **The General Problem**

- Languages (natural or artificial) are made of **sentences**, or strings of words, and the aim of a syntactic description is to tell which sentences belong to the language and which don't.
- The lowest level syntactic units which are not described by a syntactic description are called **lexemes**. e.g. a word in an English sentence.
 - Lexemes in a programming language may be: identifiers, operators, literals and special words.

-
- A **token** in a programming language is a ‘class of lexemes’: an identifier is a token that can have lexemes such as `isFull`, and `leftIndex` etc.
 - A token may have only one lexeme: e.g. the token `plus_op` as one lexeme `+`

Example: What are the tokens and the corresponding lexemes in the following statement:

```
index = 2 * count + 17;
```

```
index = 2 * count + 17;
```

<i>Lexemes</i>	<i>Tokens</i>
index	identifier
=	equal_sign
2	int_literal
*	mult_op
count	identifier
+	plus_op
17	int_literal
;	semicolon

Language Recognizers

- Recognizers are detectors: which either accept or reject sentences, depending upon whether they belong to a language or not.
- Suppose we have a language L that uses an alphabet Σ of characters
- To define L formally using the recognition method, we would need to construct a mechanism R , called a recognition device, capable of reading strings of characters from the alphabet Σ
- R would indicate whether a given input string was or was not in L . In effect, R would either accept or reject the given string
- The syntax analysis part of a compiler is a recognizer for the language
- It need only determine whether given programs are in the language.

Language Generators

- Generators have the ability to generate sentences of a language, and maybe used to describe the language.
- We can think of the generator as having a button that produces a sentence of the language every time it is pushed
- A generator seems to be a device of limited usefulness as a language descriptor.

Formal Methods for describing syntax

- Context Free Grammars (CFGs)
- A powerful generative description of a programming language's syntax.
- It is a meta-language, as it is used to describe another language – the programming language under consideration.

-
- A statement (also called a production, or rule) in a CFG, is a mix of abstractions and concrete symbols. e.g. following is the description of an assignment statement in C

$\text{<assign>} \rightarrow \text{<var>} = \text{<expression>}$

$\text{LHS} \rightarrow \text{RHS}$

- LHS : the abstraction to be defined
- RHS: the text (mix of tokens, lexemes and further abstractions to be defined) is the definition

-
- The abstractions in a production are called non-terminals, and the concrete symbols constitute the terminals.
 - A grammar is a collection of productions
 - Non-terminals may have more than one definitions. For example, the if statement in Pascal is described as

$\langle \text{if_stmt} \rangle \rightarrow \text{if} (\langle \text{logic_expr} \rangle) \langle \text{stmt} \rangle$

$\langle \text{if_stmt} \rangle \rightarrow \text{if} (\langle \text{logic_expr} \rangle) \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle$

or with the rule

$\langle \text{if_stmt} \rangle \rightarrow \text{if} (\langle \text{logic_expr} \rangle) \langle \text{stmt} \rangle$

$\quad | \text{ if} (\langle \text{logic_expr} \rangle) \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle$

-
- Some productions may be recursive. For example, the production describing data declaration, where a list of identifiers (separated by commas) may occur:

$$\begin{aligned} \langle \text{ident_list} \rangle &\rightarrow \text{identifier} \\ &\mid \text{identifier}, \langle \text{ident_list} \rangle \end{aligned}$$

-
- How can a complete program be described syntactically, using CFG?

$\langle \text{program} \rangle \rightarrow \text{begin } \langle \text{stmt_list} \rangle \text{ end}$

$\langle \text{stmt_list} \rangle \rightarrow \langle \text{stmt} \rangle$
 $\quad \quad \quad | \langle \text{stmt} \rangle ; \langle \text{stmt_list} \rangle$

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

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

$\langle \text{expression} \rangle \rightarrow \langle \text{var} \rangle + \langle \text{var} \rangle$
 $\quad \quad \quad | \langle \text{var} \rangle - \langle \text{var} \rangle$
 $\quad \quad \quad | \langle \text{var} \rangle$

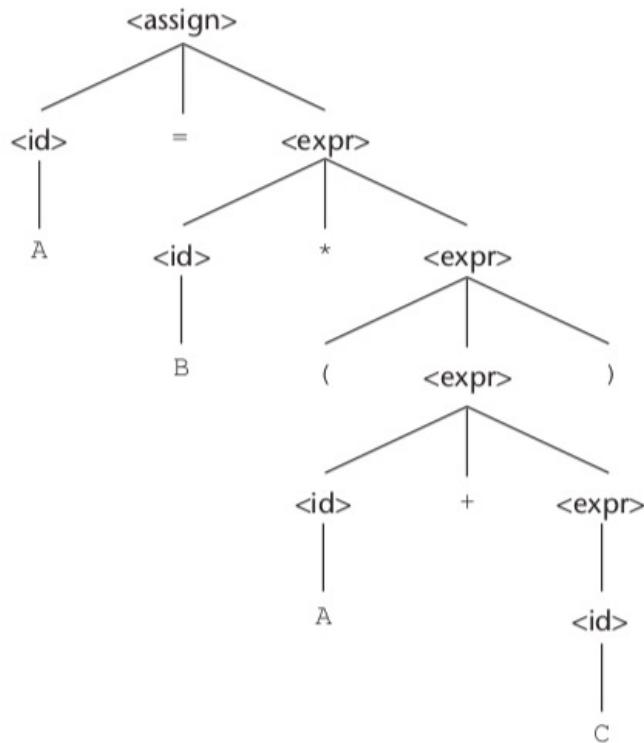
<program> => **begin** <stmt_list> **end**
=> **begin** <stmt> ; <stmt_list> **end**
=> **begin** <var> = <expression> ; <stmt_list> **end**
=> **begin** A = <expression> ; <stmt_list> **end**
=> **begin** A = <var> + <var> ; <stmt_list> **end**
=> **begin** A = B + <var> ; <stmt_list> **end**
=> **begin** A = B + C ; <stmt_list> **end**
=> **begin** A = B + ; <stmt> **end**
=> **begin** A = B + C ; <var> = <expression> **end**
=> **begin** A = B + C ; B = <expression> **end**
=> **begin** A = B + C ; B = <var> **end**
=> **begin** A = B + C ; B = C **end**

$$\langle \text{assign} \rangle \rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$$
$$\langle \text{id} \rangle \rightarrow A \mid B \mid C$$
$$\begin{aligned} \langle \text{expr} \rangle \rightarrow & \langle \text{id} \rangle + \langle \text{expr} \rangle \\ & \mid \langle \text{id} \rangle * \langle \text{expr} \rangle \\ & \mid (\langle \text{expr} \rangle) \\ & \mid \langle \text{id} \rangle \end{aligned}$$
$$A = B * (A + C)$$

$\langle \text{assign} \rangle \Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$
 $\Rightarrow A = \langle \text{expr} \rangle$
 $\Rightarrow A = \langle \text{id} \rangle * \langle \text{expr} \rangle$
 $\Rightarrow A = B * \langle \text{expr} \rangle$
 $\Rightarrow A = B * (\langle \text{expr} \rangle)$
 $\Rightarrow A = B * (\langle \text{id} \rangle + \langle \text{expr} \rangle)$
 $\Rightarrow A = B * (A + \langle \text{expr} \rangle)$
 $\Rightarrow A = B * (A + \langle \text{id} \rangle)$
 $\Rightarrow A = B * (A + C)$

Parse tree

- Parse Trees represent the inherent hierarchical structures in a language sentence.



A = B * (A + C)

$\langle \text{assign} \rangle \Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$
 $\Rightarrow A = \langle \text{expr} \rangle$
 $\Rightarrow A = \langle \text{id} \rangle * \langle \text{expr} \rangle$
 $\Rightarrow A = B * \langle \text{expr} \rangle$
 $\Rightarrow A = B * (\langle \text{expr} \rangle)$
 $\Rightarrow A = B * (\langle \text{id} \rangle + \langle \text{expr} \rangle)$
 $\Rightarrow A = B * (A + \langle \text{expr} \rangle)$
 $\Rightarrow A = B * (A + \langle \text{id} \rangle)$
 $\Rightarrow A = B * (A + C)$

Ambiguity

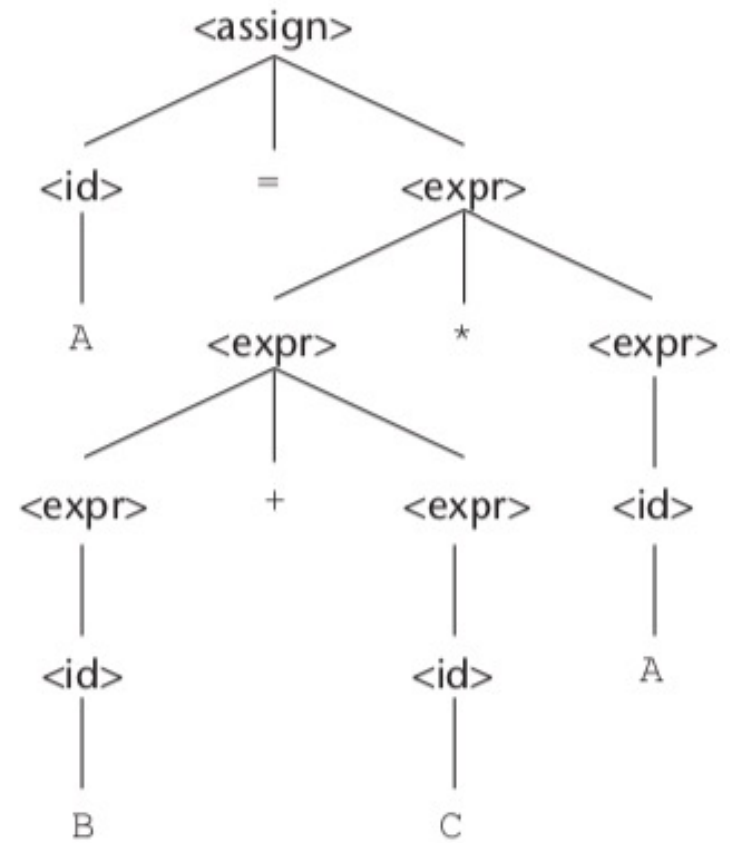
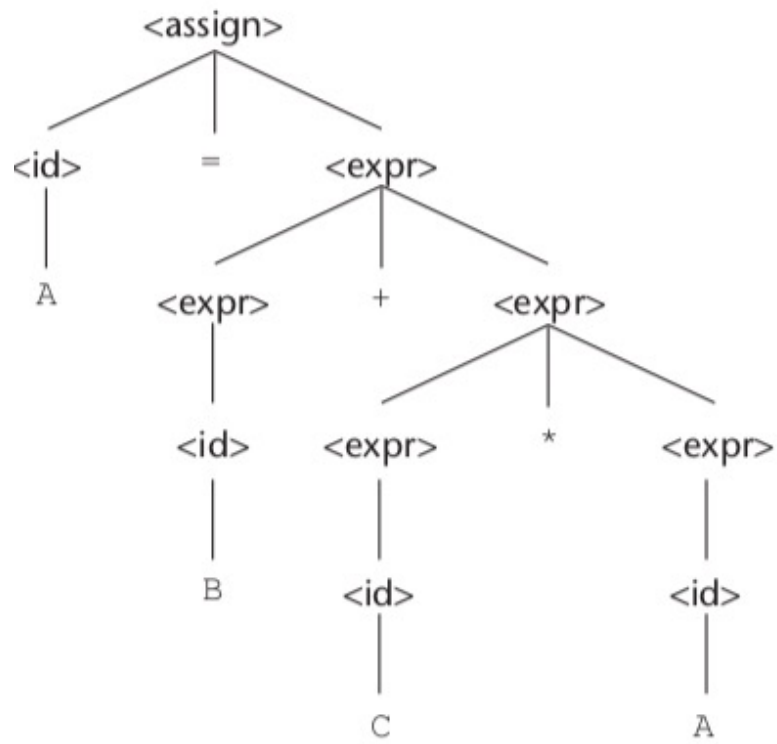
- A grammar which generates a sentences for which there are two or more distinct parse trees is ambiguous.

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

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

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

A = B + C * A



BNF:

$$\begin{aligned}\langle \text{expr} \rangle &\rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle \\ &\quad | \langle \text{expr} \rangle - \langle \text{term} \rangle \\ &\quad | \langle \text{term} \rangle \\ \langle \text{term} \rangle &\rightarrow \langle \text{term} \rangle * \langle \text{factor} \rangle \\ &\quad | \langle \text{term} \rangle / \langle \text{factor} \rangle \\ &\quad | \langle \text{factor} \rangle \\ \langle \text{factor} \rangle &\rightarrow \langle \text{exp} \rangle ** \langle \text{factor} \rangle \\ &\quad \langle \text{exp} \rangle \\ \langle \text{exp} \rangle &\rightarrow (\langle \text{expr} \rangle) \\ &\quad | \text{id}\end{aligned}$$

EBNF:

$$\begin{aligned}\langle \text{expr} \rangle &\rightarrow \langle \text{term} \rangle \{ (+ \mid -) \langle \text{term} \rangle \} \\ \langle \text{term} \rangle &\rightarrow \langle \text{factor} \rangle \{ (* \mid /) \langle \text{factor} \rangle \} \\ \langle \text{factor} \rangle &\rightarrow \langle \text{exp} \rangle \{ ** \langle \text{exp} \rangle \} \\ \langle \text{exp} \rangle &\rightarrow (\langle \text{expr} \rangle) \\ &\quad | \text{id}\end{aligned}$$

Attribute Grammars

- An attribute grammar is an extension of a CFG, designed to enforce the syntax plus the static semantics of a programming language.
- Attribute grammars are CFGS with the following additions:
 - Attributes
 - Attribute computation functions
 - Predicate functions

Static Semantics Vs Dynamic Semantics

- The static semantics of a language is only indirectly related to the meaning of programs during execution. These can be checked compile time.
- Dynamic semantics, which is the meaning of expressions, statements, and program units. These cannot be checked at runtime and can be only be checked at runtime.

Attribute Grammars

- **Attribute** : these are ‘variables’ associated with grammar symbols, and can have various values attached to them.
- **Attribute computation functions**: these are functions associated with grammar rules (productions) and specify how attributed values are computed.
- **Predicate functions**: These are associated with grammar rules (productions) and state some of the syntax and static semantic rules of the language

Attribute Grammars

- Synthesized and Inherited attributes – for each grammar symbol X , these are represented by sets $S(X)$ and $I(X)$
- Synthesized attributes are used to pass semantic information up a parse tree.
- Inherited attributes are used to pass semantic attributes down a parse tree.

Attribute Grammars

- **Syntactic Rules**

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

$\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle + \langle \text{var} \rangle \mid \langle \text{var} \rangle$

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

- **Semantic Rules**

- **actual_type**: a synthetic attribute for non-terminals, storing their actual type: int or real.
 - In case of a variable intrinsic-attributes are read-off.
 - In case of an expression the value of the attributed is determined as a function of its children nodes.
- **expected_type**: an inherited attribute associated with the non-terminal $\langle \text{expr} \rangle$
 - stores the type value expected of the $\langle \text{expr} \rangle$ given the type value of the variable on the left hand side.

Attribute Grammars

1. Syntax rule: $\langle \text{assign} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$
Semantic rule: $\langle \text{expr} \rangle.\text{expected_type} \leftarrow \langle \text{var} \rangle.\text{actual_type}$
2. Syntax rule: $\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle[2] + \langle \text{var} \rangle[3]$
Semantic rule: $\langle \text{expr} \rangle.\text{actual_type} \leftarrow$
 if $(\langle \text{var} \rangle[2].\text{actual_type} = \text{int})$ and
 $(\langle \text{var} \rangle[3].\text{actual_type} = \text{int})$
 then int
 else real
 end if

Predicate: $\langle \text{expr} \rangle.\text{actual_type} == \langle \text{expr} \rangle.\text{expected_type}$
3. Syntax rule: $\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle$
Semantic rule: $\langle \text{expr} \rangle.\text{actual_type} \leftarrow \langle \text{var} \rangle.\text{actual_type}$
Predicate: $\langle \text{expr} \rangle.\text{actual_type} == \langle \text{expr} \rangle.\text{expected_type}$
4. Syntax rule: $\langle \text{var} \rangle \rightarrow A \mid B \mid C$
Semantic rule: $\langle \text{var} \rangle.\text{actual_type} \leftarrow \text{look-up}(\langle \text{var} \rangle.\text{string})$

The look-up function looks up a given variable name in the symbol table and returns the variable's type.

Attribute Grammars

