

## Describing Syntax and Semantics – Chapter 3

### Introduction

The task of providing a concise yet understandable description of a programming language is difficult but essential to the language's success.

The study of programming languages, like the study of natural languages, can be divided into examinations of syntax and semantics.

### Syntax

- The form of a languages expressions, statements, and program units.
- `while ( <boolean_expr> ) <statement>`

### Semantics

- The meaning of those expressions, statements, and program units.
- When the Boolean expression is true the embedded statement is executed.

### The General Problem of Describing Syntax

- A language is a set of strings from some alphabet.
- The strings of a language are called sentences or statements.
- The syntax rules of a language specify which strings of characters from the language's alphabet are in the language.
- The lexemes of a language are the lowest level syntactical unit.
- The tokens of a language are the categories of its lexemes.

`index = 2 * count + 17;`

Lexemes	Tokens
index	identifier
=	assign_op
2	int_literal
*	mul_op
count	identifier
+	add_op
17	int_literal
;	semicolon

## Language Recognizers

In general, languages can be formally defined in two distinct ways: by recognition and by generation.

- A language recognizer reads strings of characters from the alphabet and indicates whether the string is part of the language or not.
- The syntax analysis part of a compiler is a recognizer for the language the compiler translates.

## Language Generators

- A language generator is a device that can be used to generate the sentences of a language.
- It is possible to determine whether the syntax of a particular statement is correct by comparing it with the structure of the generator.

## Formal Methods of Describing Syntax

### Backus-Naur Form and Context-Free Grammars

- Noam Chomsky and John Backus – developed formal syntax description – late 1950s
- Most widely used method for describing programming language syntax

### Context-free Grammars

- Noam Chomsky – linguist – developed two classes of grammars – late 1950s
- Regular grammars – used to describe the tokens of a language
- Context-free grammars – used to describe whole languages
- At the time was not interested in artificial languages

### Origins of Backus-Naur Form

- John Backus – developed to describe ALGOL 58 – 1959
- Peter Naur – revised to describe ALGOL 60 – BNF – 1960
- Nearly identical Chomsky's context-free grammars

## Fundamentals

- A metalanguage is a language used to describe another language.
- BNF is a metalanguage for programming languages.
- BNF uses abstractions for syntactic structures.
  - $\langle \text{assign} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expression} \rangle$
  - LHS:  $\langle \text{assign} \rangle$ 
    - the abstraction
    - non-terminal symbols or nonterminals
  - RHS:  $\langle \text{var} \rangle = \langle \text{expression} \rangle$ 
    - mixture of tokens, lexemes, and references to other abstractions
    - terminal symbols or terminals – tokens and lexemes
  - Referred to as a rule or production
- A Grammar is a finite nonempty set of rules

## Describing Lists

- Syntactic lists are described using recursion.
- A rule is recursive if its LHS appears on its RHS.
  - $\langle \text{ident\_list} \rangle \rightarrow \text{identifier} \mid \text{identifier}, \langle \text{ident\_list} \rangle$

## Grammars and Derivations

- A grammar is a generative device for defining languages.
- The sentences are generated through a sequence of applications of the rules, beginning with a special nonterminal of the grammar called the start symbol, often named  $\langle \text{program} \rangle$ .

### A Grammar for a Small Language

```
 $\langle \text{program} \rangle \rightarrow \text{begin } \langle \text{stmt\_list} \rangle \text{ end}$   
 $\langle \text{stmt\_list} \rangle \rightarrow \langle \text{stmt} \rangle$   
                   $\mid \langle \text{stmt} \rangle ; \langle \text{stmt\_list} \rangle$   
 $\langle \text{stmt} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expression} \rangle$   
 $\langle \text{var} \rangle \rightarrow A \mid B \mid C$   
 $\langle \text{expression} \rangle \rightarrow \langle \text{var} \rangle + \langle \text{var} \rangle$   
                   $\mid \langle \text{var} \rangle - \langle \text{var} \rangle$   
                   $\mid \langle \text{var} \rangle$ 
```

Example derivation:

$$\begin{aligned} \langle \text{program} \rangle &\Rightarrow \text{begin } \langle \text{stmt\_list} \rangle \text{ end} \\ &\Rightarrow \text{begin } \langle \text{stmt} \rangle \text{ end} \\ &\Rightarrow \text{begin } \langle \text{var} \rangle = \langle \text{expression} \rangle \text{ end} \\ &\Rightarrow \text{begin } A = \langle \text{expression} \rangle \text{ end} \\ &\Rightarrow \text{begin } A = \langle \text{var} \rangle \text{ end} \\ &\Rightarrow \text{begin } A = B \text{ end} \end{aligned}$$

- The derivation begins with the start symbol –  $\langle \text{program} \rangle$
- The symbol  $\Rightarrow$  is read “derives.”
- Each successive string is the sequence is derived from the previous string by replacing one of the nonterminals with one of its definitions.
- Each string in the derivation is called a sentential form.
- Leftmost derivation – derivation where the leftmost nonterminal is replaced in the previous string.
- Derivation continues until the sentential form consists of no nonterminals.

#### A Grammar for Simple Assignment Statements

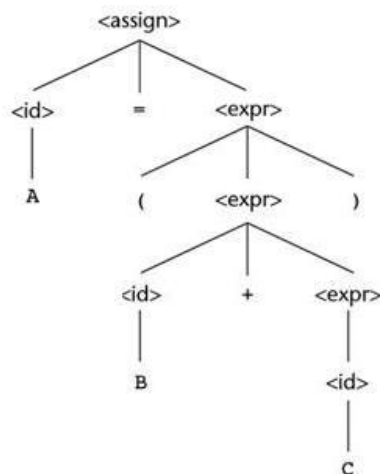
$$\begin{aligned} \langle \text{assign} \rangle &\rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle \\ \langle \text{id} \rangle &\rightarrow \text{A} \mid \text{B} \mid \text{C} \\ \langle \text{expr} \rangle &\rightarrow \langle \text{id} \rangle + \langle \text{expr} \rangle \\ &\mid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

Example derivation:

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

#### Parse Trees

- Grammars naturally describe the hierarchical syntactical structure of the sentences of the language they define.
- These hierarchical structures are called parse trees – leafs are terminal.



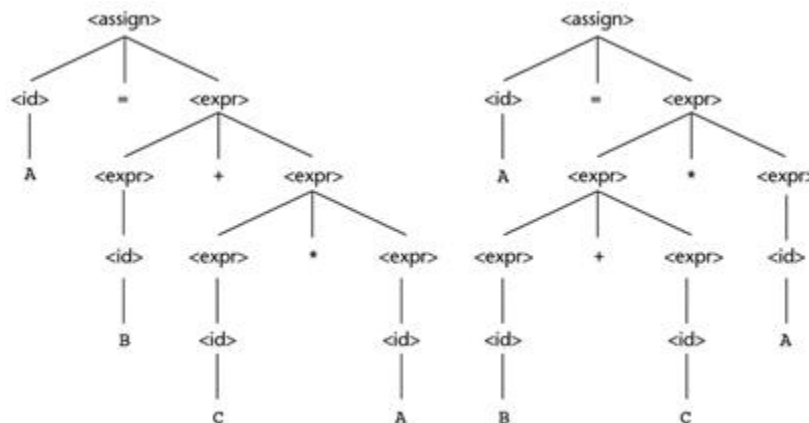
## Ambiguity

A grammar that generates a sentential form for which there are two or more distinct parse trees is said to be ambiguous.

### An Ambiguous Grammar for Simple Assignment Statements

```
<assign> → <id> = <expr>
<id> → A | B | C
<expr> → <expr> + <expr>
        | <expr> * <expr>
        | ( <expr> )
        | <id>
```

The above grammar is ambiguous because there are two distinct parse trees for the following statement:  $A = B + C * A$



- This grammar has less syntactic structure than the previous one.
  - This grammar allows the expression to grow on both sides of the operator where the previous one only allowed it to grow on the right side.
- Syntactic ambiguity is a problem because compilers often base the semantics of structures on their syntactic form.
  - If a language structure has more than one parse tree, then its meaning cannot be determined uniquely.

## Operator Precedence

- Operators generated lower in parse trees are evaluated before those generated higher up in the parse tree. This can be used to support operator precedence.
- The left parse tree shows multiplication lower in the tree than addition which is lower than assignment.
- The right parse tree shows addition lower in the tree than multiplication which is lower than assignment.

### An Unambiguous Grammar for Expressions

$\langle \text{assign} \rangle \rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$   
 $\langle \text{id} \rangle \rightarrow A \mid B \mid C$   
 $\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle$   
 $\quad \mid \langle \text{term} \rangle$   
 $\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle * \langle \text{factor} \rangle$   
 $\quad \mid \langle \text{factor} \rangle$   
 $\langle \text{factor} \rangle \rightarrow ( \langle \text{expr} \rangle )$   
 $\quad \mid \langle \text{id} \rangle$

Derivation:  $A = B + C * A$

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

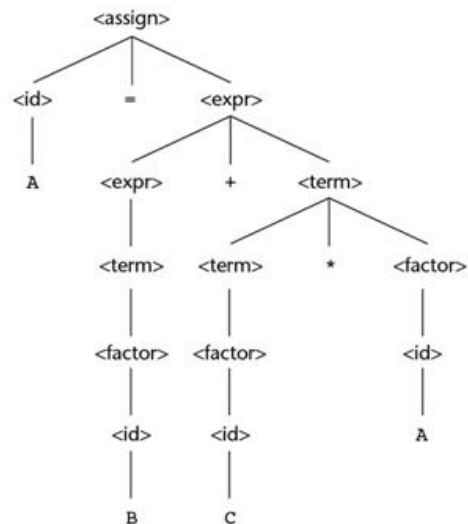
Leftmost

$\Rightarrow A = \langle \text{expr} \rangle$   
 $\Rightarrow A = \langle \text{expr} \rangle + \langle \text{term} \rangle$   
 $\Rightarrow A = \langle \text{term} \rangle + \langle \text{term} \rangle$   
 $\Rightarrow A = \langle \text{factor} \rangle + \langle \text{term} \rangle$   
 $\Rightarrow A = \langle \text{id} \rangle + \langle \text{term} \rangle$   
 $\Rightarrow A = B + \langle \text{term} \rangle$   
 $\Rightarrow A = B + \langle \text{term} \rangle * \langle \text{factor} \rangle$   
 $\Rightarrow A = B + \langle \text{factor} \rangle * \langle \text{factor} \rangle$   
 $\Rightarrow A = B + \langle \text{id} \rangle * \langle \text{factor} \rangle$   
 $\Rightarrow A = B + C * \langle \text{factor} \rangle$   
 $\Rightarrow A = B + C * \langle \text{id} \rangle$   
 $\Rightarrow A = B + C * A$

Rightmost

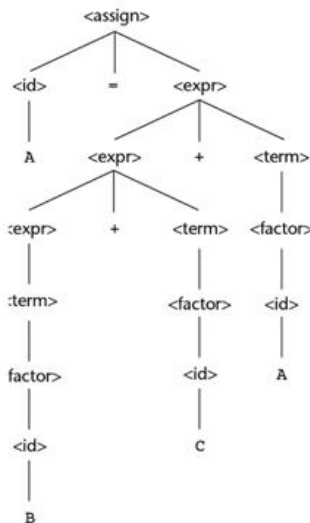
$\Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle + \langle \text{term} \rangle$   
 $\Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle + \langle \text{term} \rangle * \langle \text{factor} \rangle$   
 $\Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle + \langle \text{term} \rangle * \langle \text{id} \rangle$   
 $\Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle + \langle \text{term} \rangle * A$   
 $\Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle + \langle \text{factor} \rangle * A$   
 $\Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle + \langle \text{id} \rangle * A$   
 $\Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle + C * A$   
 $\Rightarrow \langle \text{id} \rangle = \langle \text{term} \rangle + C * A$   
 $\Rightarrow \langle \text{id} \rangle = \langle \text{factor} \rangle + C * A$   
 $\Rightarrow \langle \text{id} \rangle = \langle \text{id} \rangle + C * A$   
 $\Rightarrow \langle \text{id} \rangle = B + C * A$   
 $\Rightarrow A = B + C * A$

- Every derivation with an unambiguous grammar has a unique parse tree.
- Both of the above derivations are represented by the same parse tree.



## Associativity of Operators

- This is the order of evaluation of operators with the same precedence.
- Do parse trees for expressions with two or more adjacent occurrences of equal precedence operators have those occurrences in proper hierarchical order?
  - $A = B + C + A$
  - Parse trees shows the correct order using left associativity.



- Left recursive – a grammar rule that has its LHS as the beginning of its RHS
  - $\langle \text{epxr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle$
  - Supports left associativity
- Right recursive – a grammar rule that has its LHS as the ending of its RHS
  - $\langle \text{factor} \rangle \rightarrow \langle \text{exp} \rangle ** \langle \text{factor} \rangle$
  - Supports right associativity

## Extended BNF – EBNF

- Deals with minor inconveniences in BNF
- Improves readability and writability of BNF
- Added optional RHS component representation
  - $\langle \text{selection} \rangle \rightarrow \text{if} ( \langle \text{expression} \rangle ) \langle \text{statement} \rangle [ \text{else} \langle \text{statement} \rangle ]$
- Added repeated RHS component representation – replaces recursive form
  - $\langle \text{ident\_list} \rangle \rightarrow \langle \text{identifier} \rangle \{ , \langle \text{identifier} \rangle \}$
- Added multiple choice RHS component representation – replaces multiple rules
  - $\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle ( * | / | \% ) \langle \text{factor} \rangle$

## BNF and EBNF Versions of an Expression Grammar

BNF:

```
<expr> → <expr> + <term>
        | <expr> - <term>
        | <term>
<term> → <term> * <factor>
        | <term> / <factor>
        | <factor>
<factor> → <exp> ** <factor>
          | <exp>
<exp> → ( <expr> )
        | id
```

EBNF:

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

## Recent variations on BNF and EBNF

- Colon replaces the arrow – RHS placed on next line
- Removal of | for alternative options – each placed on a separate line
- Subscript<sub>opt</sub> used instead of the [] for optional components
  - ConstructorDeclarator -> SimpleName ( FormalParameterList<sub>opt</sub> )
- Use “one of” instead of ( | ) to indicate a choice
  - AssignmentOperator -> one of = \*= /= %= += -= <<= >>= &= ^= |=

## Attribute Grammars

- An attribute grammar is a device used to describe more of the structure of a programming language than can be described with a context-free grammar.
- An attribute grammar is an extension of a context-free grammar.
  - Allows for certain language rules to be described, like type compatibility.

## Static Semantics

- There are some characteristics of the structure of programming languages that are difficult to describe with BNF, and some that are impossible.
  - Java type compatibility rules:
    - integer assigned to float-point allowed – reverse not allowed
  - Could be specified in BNF – grammar would become too large
- The static semantics of a language is only indirectly related to the meaning of programs during execution; rather, it has to do with the legal forms of a program (syntax rather than semantics).
- Static semantics is so named because the analysis required to check these specifications can be done at compile time.



## Basic Concepts

- Attribute grammars are context-free grammars with attributes, attribute computation functions, and predicate functions.
- Attributes – are similar to variables – can have values assigned to them.
- Attribute computation functions (semantic functions) – associated with grammar rules.
- Predicate functions – state semantic rules – associated with grammar rules.

## Attribute Grammars Defined

- Each grammar symbol  $X$  has a set of attributes values  $A(X)$  consisting of:
  - Synthesized attributes  $S(X)$  – used to pass semantic information up a parse tree.
  - Inherited attributes  $I(X)$  – used to pass semantic information down and across a parse tree.
- Each grammar rule has a set of semantic functions that define certain attributes of the nonterminals in the rule.
  - For rule  $X_0 \rightarrow X_1 \dots X_n$ 
    - $S(X_0) = f(A(X_1) \dots A(X_n))$  – functions define synthesized attributes
    - $I(X_0) = f(A(X_1) \dots A(X_n))$  – functions define inherited attributes
- Each grammar rule has a set of predicates to check for attribute consistency. This set may be empty.
  - Predicate functions have the form of a Boolean expression on the union of the attribute set  $\{A(X_0) \dots A(X_n)\}$  and a set of literal attribute values.
  - With all allowed derivations every predicate associated with every nonterminal is true.
  - A false predicate function value indicates a syntax or static semantics rule violation.

A parse tree of an attribute grammar is the parse tree based on its underlying BNF grammar, with a possible empty set of attribute values attached to each node. If all the attributes values in a parse tree have been computed, the tree is said to be fully attributed. This is not always done in practice.

## Intrinsic Attributes

- These are synthesized attributes of the leaf nodes whose values are determined outside of the parse tree.
  - Variable data type could come from the symbol table.
- Given the intrinsic attribute values on a parse tree, the semantic functions can be used to compute the remaining attributes.

## Examples of Attribute Grammars

$\langle \text{assign} \rangle \rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$   
 $\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle + \langle \text{var} \rangle \mid \langle \text{var} \rangle$   
 $\langle \text{var} \rangle \rightarrow A \mid B \mid C$

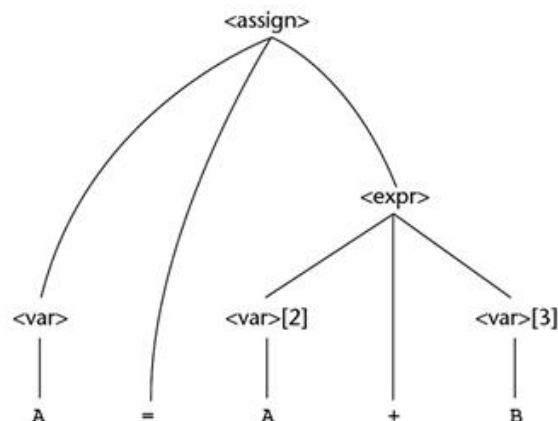
- **actual\_type**
  - Synthesized attribute for  $\langle \text{var} \rangle$  and  $\langle \text{expr} \rangle$  – stores the actual type
    - variable type is intrinsic
    - expression – determined from the actual types of its children
- **expected\_type**
  - Inherited attribute  $\langle \text{expr} \rangle$  – stores the expected type for the expression
  - type determined by variable on left side of assignment statement

### An Attribute Grammar for Simple Assignment Statements

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.

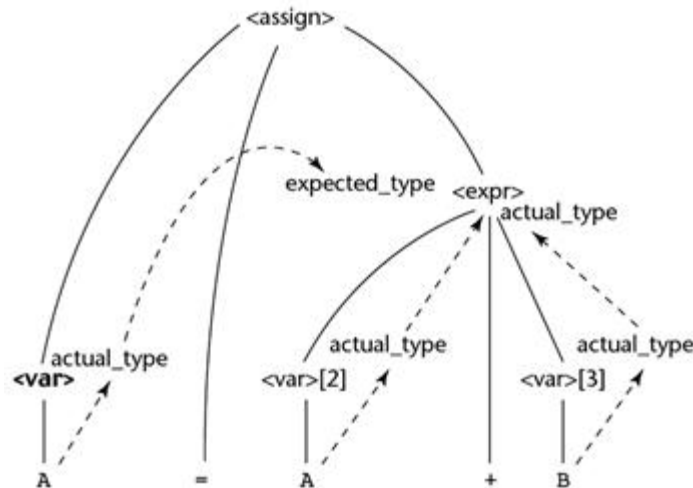
Parse tree using the above grammar for the statement  $A = A + B$



## Computing Attribute Values

- Top-down if all attributes were inherited.
- Bottom-up if all attributes were synthesized.
- Both directions because there are both inherited and synthesized attributes

1.  $\langle \text{var} \rangle.\text{actual\_type} \leftarrow \text{look-up}(A)$  – Rule 4
2.  $\langle \text{expr} \rangle.\text{expected\_type} \leftarrow \langle \text{var} \rangle.\text{actual\_type}$  – Rule 1
3.  $\langle \text{var} \rangle[2].\text{actual\_type} \leftarrow \text{look-up}(A)$  – Rule 4  
 $\langle \text{var} \rangle[3].\text{actual\_type} \leftarrow \text{look-up}(B)$  – Rule 4
4.  $\langle \text{expr} \rangle.\text{actual\_type} \leftarrow \text{either int or real}$  – Rule 2
5.  $\langle \text{expr} \rangle.\text{expected\_type} == \langle \text{expr} \rangle.\text{actual\_type}$  TRUE or FALSE – Rule 2



Assume that the data type of A is real and that the data type of B is int.

