

Syntax and Semantics

Hui Chen

Computer Science
Virginia State University
Petersburg, Virginia

January 20, 2016

Outline

- ▶ Backus-Naur Form
 - ▶ derivations, parse trees, ambiguity, descriptions of operator precedence and associativity, and extended Backus-Naur Form.
- ▶ Attribute grammars
- ▶ Operational axiomatic and denotational semantics

Chomsky Hierarchy

- ▶ Also called Chomsky-Schützenberger Hierarchy (Noam Chomsky, 1956)

Class	Grammar	Language	Automaton
Type-0	Unrestricted	Recursively enumerable	Turing machine (TM)
Type-1	Context-sensitive	Context-sensitive	Linear-bounded automaton (LBA)
Type-2	Context-free	Context-free	Pushdown automaton (PDA)
Type-3	Regular	Regular	Deterministic finite automaton (DFA)

- ▶ A strictly nested sets of classes of formal grammars, i.e.,

$$\text{Type-0} \supset \text{Type-1} \supset \text{Type-2} \supset \text{Type-3}$$

- ▶ Context-free and regular grammars are of our primary concern

Context-Free Grammar (CFG)

- ▶ A CFG is a quadruple, $G = (V, T, P, S)$ where
 - ▶ V : the set of variables or non-terminals
 - ▶ T : the set of terminals
 - ▶ P : the set of productions of the form $A \rightarrow \gamma$ where A is a single variable, i.e., $A \in V$ and γ is string of terminals and variables, i.e., $\gamma \in (V \cup T)^*$
 - ▶ S : the start symbol and $S \in V$
- ▶ To describe the grammar of a programming language,
 - ▶ Terminals are lexemes or tokens

Example: A Simple Programming Language¹

- ▶ Operators: $+$ and $*$ represent addition and multiplication, respectively
- ▶ Arguments are identifiers consisting *only* of letters a , b , and digits 0, 1
- ▶ An example statement in the language,

$$(a + b) * (a + b + 1)$$

¹This is an example given in [Hopcroft et al., 2006]

CFG of the Simple Language

- ▶ The language can be specified using a CFG as,

$$G = (\{E, I\}, T, P, E)$$

where

- ▶ E and I are the two variables, and E is the start symbol
- ▶ T , the terminals are the set of symbols $\{+, *, (,), a, b, 0, 1\}$
- ▶ P is the productions, i.e.,

$$1 \quad E \rightarrow I$$

$$2 \quad E \rightarrow E + E$$

$$3 \quad E \rightarrow E * E$$

$$4 \quad E \rightarrow (E)$$

$$5 \quad I \rightarrow a$$

$$6 \quad I \rightarrow a$$

$$7 \quad I \rightarrow Ia$$

$$8 \quad I \rightarrow Ib$$

$$9 \quad I \rightarrow I0$$

$$10 \quad I \rightarrow I1$$

Backus-Naur Form (BNF)

- ▶ John Backus (1959) and Peter Naur (1960) developed to describe syntax of ALGOL 58 and 60
- ▶ BNF is equivalent to context-free grammars
- ▶ Widely used today for describing syntax of programming languages

Production Rules in BNF

- ▶ Nonterminals (or variables in CFG, called *abstractions*) are often enclosed in angle brackets
- ▶ A start symbol is a special element of the nonterminals of a grammar
- ▶ Grammar: a finite non-empty set of rules
- ▶ Examples of BNF rules:

`<ident_list > → identifier`

`<ident_list > → identifier, <ident_list >`

`<if_stmt > → if <logic_expr> then <stmt >`

More than one RHS

- ▶ An abstraction (or a nonterminal symbol) can have more than one right-hand sides
- ▶ Example: applying this rule, we can rewrite,

$$\langle \text{ident_list} \rangle \rightarrow \text{identifier}$$

$$\langle \text{ident_list} \rangle \rightarrow \text{identifier}, \langle \text{ident_list} \rangle$$

as

$$\langle \text{ident_list} \rangle \rightarrow \text{identifier} \mid \text{identifier}, \langle \text{ident_list} \rangle$$

- ▶ Another example:

$$\langle \text{stmt} \rangle \rightarrow \langle \text{single_stmt} \rangle \mid \text{begin } \langle \text{stmt_list} \rangle \text{ end}$$

Lists

- ▶ Syntactic lists are described using recursion

$$\langle \text{ident_list} \rangle \rightarrow \text{ident} \mid \text{ident}, \langle \text{ident_list} \rangle$$

Derivation

- ▶ A repeated application of rules, starting with the start symbol and ending with a sentence (all terminal symbols)
 - ▶ Every string of symbols in a 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

An Example of Derivation

- ▶ Given a grammar,

$$\langle \text{program} \rangle \rightarrow \langle \text{stmts} \rangle$$

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

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

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

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

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

- ▶ we can have the following derivation,

$$\langle \text{program} \rangle \Rightarrow \langle \text{stmts} \rangle \Rightarrow \langle \text{stmt} \rangle \Rightarrow \langle \text{var} \rangle = \langle \text{expr} \rangle$$

$$\Rightarrow a = \langle \text{expr} \rangle \Rightarrow a = \langle \text{term} \rangle + \langle \text{term} \rangle$$

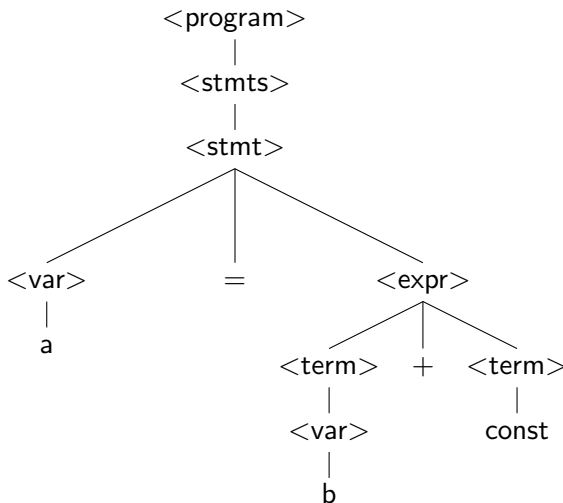
$$\Rightarrow a = \langle \text{var} \rangle + \langle \text{term} \rangle$$

$$\Rightarrow a = b + \langle \text{term} \rangle$$

$$\Rightarrow a = b + \text{const}$$

Parse Tree

- ▶ A parse tree is a hierarchical representation of a derivation
- ▶ Example:

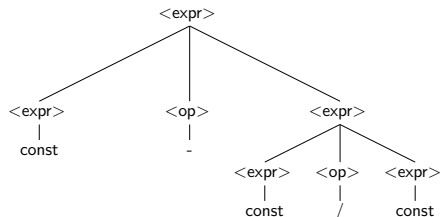
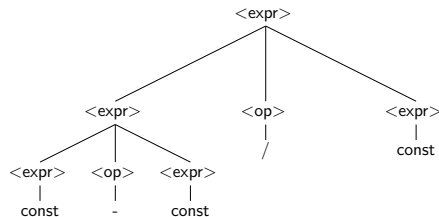


Ambiguity in Grammars

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

Example of Ambiguous Grammar and Parse Trees

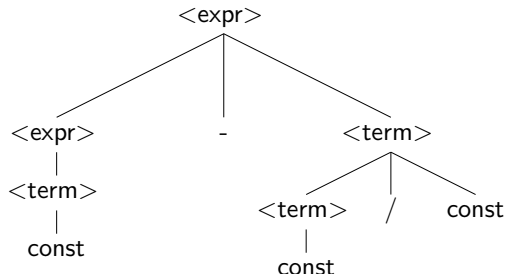
$$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \mid \text{const}$$

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


Unambiguous Grammar

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

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

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


Associativity of Operators

- ▶ Operator associativity can also be indicated by a grammar
- ▶ Example: compare the following two grammars

1. Ambiguous grammar

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

2. Unambiguous grammar

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

Extended BNF (EBNF)

- ▶ The extensions *do not* enhance the descriptive power of BNF; they only increase its *readability* and *writability*
- ▶ Optional parts are placed in brackets [], e.g.,

$$\langle \text{proc_call} \rangle \rightarrow \text{ident} [(\langle \text{expr_list} \rangle)]$$

- ▶ Alternative parts of RHSs are placed inside () and separated via |, e.g.,

$$\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle (+|-) \text{const}$$

- ▶ Repetitions (0 or more times) are placed inside { },

$$\langle \text{ident} \rangle \rightarrow \text{letter} \{ \text{letter} | \text{digit} \}$$

- ▶ Can you rewrite the above examples without using extensions?

Recent Variations in EBNF

- ▶ Alternative RHSs are put on separate lines
- ▶ Use of `a :` instead of `→`
- ▶ Use of `opt` for optional parts
- ▶ Use of `oneof` for choices

Static Semantics

- ▶ Context-free grammars (CFGs) has limitations to describe the syntax of programming languages
 - ▶ Some are context-free, but cumbersome to be described in CFGs, e.g., type constraints
 - ▶ Some are non context-free, e.g., variables must be declared before they are used
- ▶ Static semantics rules: checking and analysis of the rules can be done at compile time

Attribute Grammar

- ▶ Formal approach both to describing and checking the correctness of the static semantics rules of a program
- ▶ Additions to CFGs to carry some semantic info on parse tree nodes
 - ▶ Static semantics specification
 - ▶ Static semantics checking

Definition of Attribute Grammar

- ▶ An attribute grammar is a context-free grammar $G = (S, N, T, P)$ with the following additions:
 - ▶ For each grammar symbol x there is a set $A(x)$ of attribute values
 - ▶ Each rule has a set of functions that define certain attributes of the nonterminals in the rule
 - ▶ Each rule has a (possibly empty) set of predicates to check for attribute consistency

Rules in Attribute Grammar

- ▶ Let $X_0 \rightarrow X_1 \dots X_n$ be a rule
- ▶ Functions of the form $S(X_0) = f(A(X_1), \dots, A(X_n))$ define synthesized attributes
- ▶ Functions of the form $I(X_j) = f(A(X_0), \dots, A(X_n))$, for $i \leq j \leq n$, define inherited attributes
- ▶ Initially, there are intrinsic attributes on the leaves

An Example of Attribute Grammars

- ▶ Syntax

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

- ▶ actual_type: synthesized for $\langle \text{var} \rangle$ and $\langle \text{expr} \rangle$
- ▶ expected_type: inherited for $\langle \text{expr} \rangle$

An Example of Attribute Grammars

- ▶ Syntax rule:

$$\langle \text{expr} \rangle \rightarrow \langle \text{var} \rangle [1] + \langle \text{var} \rangle [2]$$

- ▶ Semantic rules:

$$\langle \text{expr} \rangle .\text{actual_type} \rightarrow \langle \text{var} \rangle [1].\text{actual_type}$$

- ▶ Predicate:

$$\langle \text{var} \rangle [1].\text{actual_type} == \langle \text{var} \rangle [2].\text{actual_type}$$

$$\langle \text{expr} \rangle .\text{expected_type} == \langle \text{expr} \rangle .\text{actual_type}$$

- ▶ Syntax rule:

$$\langle \text{var} \rangle \rightarrow \text{id}$$

- ▶ Semantic rule:

$$\langle \text{var} \rangle .\text{actual_type} \leftarrow \text{lookup}(\langle \text{var} \rangle .\text{string})$$

Compute Attribute Values

- ▶ If all attributes were *inherited*, the tree could be decorated in top-down order.
- ▶ If all attributes were *synthesized*, the tree could be decorated in bottom-up order.
- ▶ In many cases, both kinds of attributes are used, and it is some combination of top-down and bottom-up that must be used.

An Example of Computing Attribute Values

$\langle \text{expr} \rangle . \text{expected_type} \leftarrow \text{inherited from parent}$

$\langle \text{var} \rangle [1]. \text{actual_type} \leftarrow \text{lookup}(A)$

$\langle \text{var} \rangle [2]. \text{actual_type} \leftarrow \text{lookup}(B)$

$\langle \text{var} \rangle [1]. \text{actual_type} == \langle \text{var} \rangle [2]. \text{actual_type}$

$\langle \text{expr} \rangle . \text{actual_type} \leftarrow \langle \text{var} \rangle [1]. \text{actual_type}$

$\langle \text{expr} \rangle . \text{actual_type} == \langle \text{expr} \rangle . \text{expected_type}$

Dynamic Semantics

- ▶ meaning, of the expressions, statements, and program units of a programming language
- ▶ need for a methodology and notation for describing semantics.
 - ▶ Programmers need to know what statements mean
 - ▶ Compiler writers must know exactly what language constructs do
 - ▶ Correctness proofs would be possible
 - ▶ Compiler generators would be possible
 - ▶ Designers could detect ambiguities and inconsistencies

Describing Semantics

- ▶ no universally accepted notation or approach has been devised for dynamic semantics
- ▶ briefly describe several of the methods that have been developed
 - ▶ Operational Semantics
 - ▶ Denotational Semantics
 - ▶ Axiomatic Semantics

Operational Semantics

- ▶ To describe the meaning of a statement or program by specifying the effects of running it on a machine.
- ▶ The effects on the machine are viewed as the sequence of changes in its state (memory, registers, etc.)
- ▶ To use operational semantics for a high-level language, a *virtual machine* or an idealized computers is used

Applications of Operational Semantics

- ▶ A complete computer simulation
- ▶ The process:
 - ▶ Build a translator (translates source code to the machine code of an idealized computer)
 - ▶ Build a simulator for the idealized computer
- ▶ Evaluation of operational semantics:
 - ▶ Good if used informally (language manuals, etc.)
 - ▶ Extremely complex if used formally (e.g., VDL), it was used for describing semantics of PL/I.

Evaluation

- ▶ good if used informally (e.g., in programming language manuals)
- ▶ extremely complex if used formally (e.g., VDL)

Denotational Semantics

- ▶ Originally developed
in [Strachey and Scott, 1970, Scott and Strachey, 1971]
- ▶ The most rigorous and most widely known formal method for
describing the meaning of programs
- ▶ Based on recursive function theory

Evaluation

- ▶ Can be used to prove the correctness of programs
- ▶ Provides a rigorous way to think about programs
- ▶ Can be an aid to language design
- ▶ Has been used in compiler generation systems
- ▶ Because of its complexity, it are of little use to language users

Axiomatic Semantics

- ▶ Based on formal logic (predicate calculus)
- ▶ Original purpose: formal program verification
- ▶ Axioms or inference rules are defined for each statement type in the language (to allow transformations of logic expressions into more formal logic expressions)
- ▶ The logic expressions are called *assertions*

Assertions in Axiomatic Semantics

- ▶ An assertion before a statement (a precondition) states the relationships and constraints among variables that are true at that point in execution
- ▶ An assertion following a statement is a postcondition
- ▶ A weakest precondition is the least restrictive precondition that will guarantee the postcondition

Evaluation

- ▶ Developing axioms or inference rules for all of the statements in a language is difficult
- ▶ It is a good tool for correctness proofs, and an excellent framework for reasoning about programs, but it is not as useful for language users and compiler writers
- ▶ Its usefulness in describing the meaning of a programming language is limited for language users or compiler writers

Denotation and Operational Semantics

- ▶ In operational semantics, the state changes are defined by coded algorithms
- ▶ In denotational semantics, the state changes are defined by rigorous mathematical functions

Summary

- ▶ BNF and context-free grammars are equivalent meta-languages
 - ▶ Well-suited for describing the syntax of programming languages
- ▶ An attribute grammar is a descriptive formalism that can describe both the syntax and the semantics of a language
- ▶ Three primary methods of semantics description
 - ▶ Operation, axiomatic, denotational

References I



Hopcroft, J. E., Motwani, R., and Ullman, J. D. (2006).
Introduction to Automata Theory, Languages, and Computation (3rd Edition).
Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA.



Scott, D. S. and Strachey, C. (1971).
Toward a mathematical semantics for computer languages, volume 1.
Oxford University Computing Laboratory, Programming Research Group.



Strachey, C. and Scott, D. (1970).
Mathematical semantics for two simple languages.
Princeton Univ.