

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:

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

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

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

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,

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

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

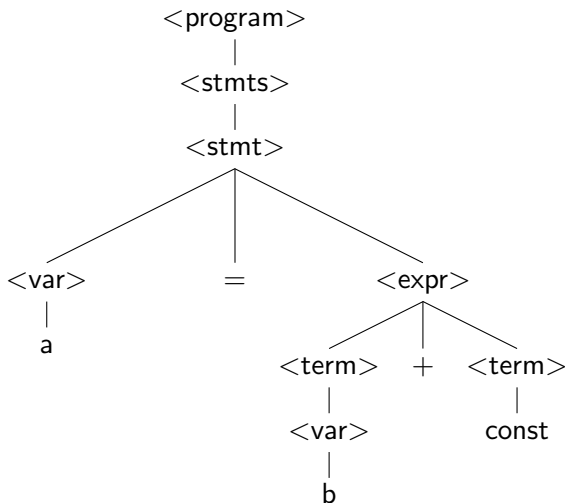
$$\Rightarrow a = \langle \text{var} \rangle + \langle \text{term} \rangle$$

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

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

Parse Tree

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

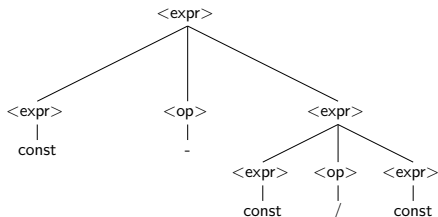
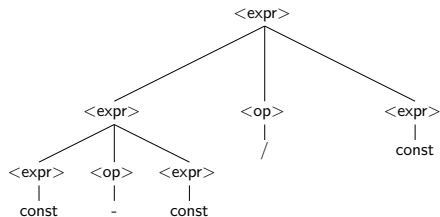


Ambiguity in Grammars

- ▶ 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

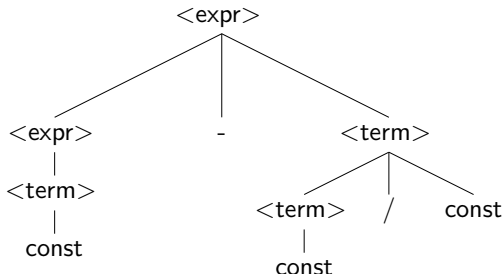
$$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \mid \text{const}$$

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


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 `<var>` and `<expr>`
- ▶ `expected_type`: inherited for `<expr>`

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:

$$\begin{aligned} \langle \text{var} \rangle [1].\text{actual_type} &== \langle \text{var} \rangle [2].\text{actual_type} \\ \langle \text{expr} \rangle .\text{expected_type} &== \langle \text{expr} \rangle .\text{actual_type} \end{aligned}$$

- ▶ 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

Constructing Denotational Semantics Specification

- ▶ define *syntactic domain*: mathematical objects for language entities
- ▶ define *semantic domain*: function that maps language entities onto mathematical objects
- ▶ syntactic domain (domain \mathbb{D}): collection of values that are legitimate parameters to the function
- ▶ semantic domain (range \mathbb{R}): collection of objects to which the parameters are mapped

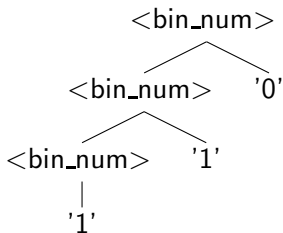
$$f : \mathbb{D} \mapsto \mathbb{R}$$

An Example: Binary Numbers

- ▶ Grammar for binary numbers

$$\begin{aligned}
 \langle \text{bin_num} \rangle &\rightarrow '0' \\
 &| '1' \\
 &| \langle \text{bin_num} \rangle '0' \\
 &| \langle \text{bin_num} \rangle '1'
 \end{aligned}$$

- ▶ Parse Tree for 110



An Example: Binary Numbers

- ▶ Now need to define the *meaning* of binary numbers
- ▶ *syntactic domain* (domain):

$$\mathbb{D} = \{'0', '1', \langle \text{bin_num} \rangle '0', \langle \text{bin_num} \rangle '1'\}$$

- ▶ *semantic domain* (range):

$$\mathbb{R} = \{0, 1, 2 \cdot M_{bin}(\langle \text{bin_num} \rangle), 2 \cdot M_{bin}(\langle \text{bin_num} \rangle) + 1\}$$

- ▶ mapping from domain onto range $M_{bin} : \mathbb{D} \mapsto \mathbb{R}$

$$M_{bin}('0') = 0$$

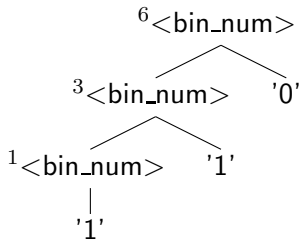
$$M_{bin}('1') = 1$$

$$M_{bin}(\langle \text{bin_num} \rangle '0') = 2 \cdot M_{bin}(\langle \text{bin_num} \rangle)$$

$$M_{bin}(\langle \text{bin_num} \rangle '1') = 2 \cdot M_{bin}(\langle \text{bin_num} \rangle) + 1$$

An Example: Binary Numbers

- ▶ Decorated Parse Tree for 110



An Example: Decimal Numbers

- ▶ Grammar for decimal numbers (in EBNF)

$$\begin{aligned} \langle \text{dec_num} \rangle &\rightarrow '0'|'1'|'2'|'3'|'4'|'5'|'6'|'7'|'8'|'9' \\ &|\langle \text{dec_num} \rangle ('0'|'1'|'2'|'3'|'4'|'5'|'6'|'7'|'8'|'9') \end{aligned}$$

- ▶ What are the mapping function and its syntactic and semantic domains?
- ▶ Can you provide an example of decorated parse tree for 3231?

The State of a Program

- ▶ Let the state s of a program be represented as a set of ordered pairs

$$s = \{ \langle i_1, v_1 \rangle, \langle i_2, v_2 \rangle, \dots, \langle i_n, v_n \rangle \}$$

where

i_j is the name of j -th variable and v_j is the current value of variable i_j , $1 \leq j \leq n$. The value of v_j can be the special value *undef*, which indicates that its associated variable is currently undefined.

- ▶ Let *VARMAP* be a function of two parameters: a variable name and the program state. The value of *VARMAP*(i_j, s) is v_j (the value paired with i_j in state s).
- ▶ See the binary numbers and decimal numbers examples

Three Language Constructs

- ▶ Expressions
- ▶ Assignment Statements
- ▶ Logical Pretest Loops

Expressions

- ▶ Map expressions onto $\mathbb{Z} \cup \{error\}$
- ▶ Assume
 - ▶ Expressions have no side effects
 - ▶ Operators are $+$ and $*$
 - ▶ Expression can have at most one operator
 - ▶ Only operands are scalar integer variables and integer literals
 - ▶ There are no parentheses
 - ▶ The value of an expression is an integer.

Grammar of Expressions

▶ Grammar

$\langle \text{expr} \rangle \rightarrow \langle \text{dec_num} \rangle \mid \langle \text{var} \rangle \mid \langle \text{binary_expr} \rangle$

$\langle \text{binary_expr} \rangle \rightarrow \langle \text{left_expr} \rangle \langle \text{operator} \rangle \langle \text{right_expr} \rangle$

$\langle \text{left_expr} \rangle \rightarrow \langle \text{dec_num} \rangle \mid \langle \text{var} \rangle$

$\langle \text{right_expr} \rangle \rightarrow \langle \text{dec_num} \rangle \mid \langle \text{var} \rangle$

$\langle \text{operator} \rangle \rightarrow + \mid *$

Mapping Function for the Expressions

- ▶ use $\Delta =$ to define mathematical functions

$$\begin{aligned}
 M_e(\langle \text{expr} \rangle, s) \Delta = & \text{case } \langle \text{expr} \rangle \text{ of} \\
 & \langle \text{dec_num} \rangle \Rightarrow M_{\text{dec}}(\langle \text{dec_num} \rangle, s) \\
 & \langle \text{var} \rangle \Rightarrow \text{if VARMAP}(\langle \text{var} \rangle, s) == \mathbf{undef} \\
 & \quad \text{then } \mathbf{error} \\
 & \quad \text{else VARMAP}(\langle \text{var} \rangle, s) \\
 & \langle \text{binary_expr} \rangle \Rightarrow \\
 & \quad \text{if}(M_e(\langle \text{binary_expr} \rangle.\langle \text{left_expr} \rangle, s) == \mathbf{undef} \text{ OR} \\
 & \quad \quad M_e(\langle \text{binary_expr} \rangle.\langle \text{right_expr} \rangle, s) == \mathbf{undef}) \\
 & \quad \text{then } \mathbf{error} \\
 & \quad \text{else if } (\langle \text{binary_expr} \rangle.\langle \text{operator} \rangle == '+') \\
 & \quad \quad \text{then } M_e(\langle \text{binary_expr} \rangle.\langle \text{left_expr} \rangle, s) + \\
 & \quad \quad \quad M_e(\langle \text{binary_expr} \rangle.\langle \text{right_expr} \rangle, s) \\
 & \quad \quad \text{else } M_e(\langle \text{binary_expr} \rangle.\langle \text{left_expr} \rangle, s) * \\
 & \quad \quad \quad M_e(\langle \text{binary_expr} \rangle.\langle \text{right_expr} \rangle, s)
 \end{aligned}$$

Assignment Statement

- Maps state sets to state sets $s \cup \{error\}$

$M_a(x = E, s) \Delta=$ if $M_e(E, s) == \mathbf{error}$

then **error**

else $s' = \{ \langle i_1, v_1' \rangle, \langle i_2, v_2' \rangle, \dots, \langle i_n, v_n' \rangle \}$, where

for $j = 1, 2, \dots, n$

if $i_j == x$

then $v_j' = M_e(E, s)$

else $v_j' = \text{VARMAP}(i_j, s)$

Logical Pretest Loops

- ▶ Maps state sets to state sets $s \cup \{error\}$

$$M_l(\mathbf{while} B \mathbf{do} L, s) \Delta= \begin{array}{l} \text{if } M_b(B, s) == \mathbf{undef} \\ \text{then } \mathbf{error} \\ \text{else if } M_b(B, s) == \mathbf{false} \\ \text{then } s \\ \text{else if } M_{sl}(L, s) == \mathbf{error} \\ \text{then } \mathbf{error} \\ \text{else } M_l(\mathbf{while} B \mathbf{do} L, M_{sl}(L, s)) \end{array}$$

Meaning of Loops

- ▶ The value of the program variables after the statements in the loop have been executed the prescribed number of times, assuming there have been no errors.
- ▶ The loop has been converted from *iteration* to *recursion*, where the recursion control is mathematically defined by other recursive state mapping functions
- ▶ *Recursion* is easier to describe with mathematical rigor than *iteration*.
- ▶ Observation: according to the definition, like actual program loops, may compute nothing because of nontermination

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.