

COMP 356
Programming Language Structures
Notes for Classes 1–3 (Chapter 3 of *Concepts of Programming Languages*)
Syntax and Semantics

Acknowledgment: The majority of these notes were written by Prof. Tim Wahls; however, some significant changes have been made.

1 Introduction: why study the structure of programming languages?

The theory of computability tells us that all modern programming languages are *Turing complete*. In other words, these languages can be used to solve exactly the same set of problems (the set of *computable* problems). In fact, some surprisingly minimal subsets of a language are Turing complete (e.g. `if` and `goto` are sufficient; `if` and `while` are sufficient).

So, (1) why would we ever need more than one programming language? And (2) why would we bother to study the implementation of and differences between programming languages?

Possible answers to (1) include:

- different languages make it easier or harder to achieve certain tasks. Examples?
- different languages achieve different performance for certain tasks. Examples?

Possible answers to (2) include:

- Understanding implementation of a language can help improve the quality of the programs you write in that language. Examples?
- Understanding different languages broadens your selection of problem-solving techniques. Examples?
- Various practical advantages e.g. Knowing several languages makes you more attractive to an employer; you can select an appropriate language for a given task.
- Some fundamental ideas in the theory of computer science (e.g. recursion, abstraction) are reflected in the design of some languages. Examples?

See Section 1.1 of the textbook for some more detailed discussion.

2 Syntax and semantics overview

Definitions:

- **syntax** is the form (structure, grammar) of a language
- **semantics** is the meaning of a language

Example:

```
if (a > b) a = a + 1;  
else b = b + 1;
```

- syntax: `if-else` is an operator that takes three operands—a condition and two statements
- semantics: if the value of `a` is greater than the value of `b`, then increment `a`. Otherwise, increment `b`.

Both the syntax and semantics of a programming language must be carefully defined so that:

- language implementors can implement the language (correctly), so that programs developed with one implementation run correctly under another (portability)
- programmers can use the language (correctly)

3 Describing Syntax

A **language** is a set of strings of characters from some alphabet.

Examples:

- binary numbers (using the alphabet $\{0, 1\}$)
- Java identifiers
- Java programs that compile with no errors
- English words (assuming we agree on a specific dictionary)
- English sentences (assuming we agree on a dictionary and all the grammatical rules of English, which probably isn't possible in practice)

The syntax rules of a language determine whether or not arbitrary strings belong to the language. The first step in specifying syntax is describing the basic units or “words” of the language, called **lexemes**. For example, some typical Java lexemes include:

- `if`
- `++`
- `+`
- `3.27`
- `count`

Lexemes are grouped into categories called **tokens**. Each token has one or more lexemes.

token	sample lexemes
identifier	<code>count</code> , <code>i</code> , <code>x2</code> , ...
<code>if</code>	<code>if</code>
<code>add_op</code>	<code>+</code> , <code>-</code>
<code>semi</code>	<code>;</code>
<code>int_lit</code>	<code>0</code> , <code>10</code>
<code>double_lit</code>	<code>-2.4</code>

Table 1: Example tokens and some associated lexemes.

Tokens are specified using regular expressions or finite automata.

The scanner/lexical analyzer of a compiler processes the character strings in the source program and determines the tokens that they represent.

Once the tokens of a language are defined, the next step is to determine which sequences of tokens are in the language (syntax).

3.1 Using BNF to Describe Syntax

BNF \equiv Backus-Naur Form

BNF is:

- a **metalinguage**—a language used to describe other languages
- the standard way to describe programming language syntax
- often used in language reference manuals

The class (set) of languages that can be described using BNF is called the **context-free languages**, and BNF descriptions are also called **context-free grammars** or just **grammars**.

symbol	meaning
\rightarrow	is defined as
	or (alternatives)
\langle something \rangle	a nonterminal —replace by the definition of something
something	(with no \langle \rangle) a token or terminal —a “word” in the language being defined (not replaced)
ϵ	the empty string (nothing)

Table 2: BNF notation.

Example (if statements in C):

\langle if-stmt $\rangle \rightarrow$ if (\langle expr \rangle) \langle stmt \rangle
 \langle if-stmt $\rangle \rightarrow$ if (\langle expr \rangle) \langle stmt \rangle else \langle stmt \rangle

Each definition above is called a **rule** or **production**. A full BNF definition would include definitions for all the nonterminals used— \langle expr \rangle and \langle stmt \rangle are not defined above.

The two rules above can be abbreviated using | as follows:

\langle if-stmt $\rangle \rightarrow$ if (\langle expr \rangle) \langle stmt \rangle
| if (\langle expr \rangle) \langle stmt \rangle else \langle stmt \rangle

Example (expressions in C):

\langle expr $\rangle \rightarrow$ id | num | (\langle expr \rangle) | \langle expr \rangle \langle op \rangle \langle expr \rangle
 \langle op $\rangle \rightarrow$ + | - | * | / | == | < | <= | > | >=

Note that BNF definitions can be recursive.

3.2 Derivations

A **derivation** is a sequence of replacements using the rules of a grammar. Derivations:

- are often used to show that a particular sequence of tokens belongs to the language defined by the grammar
- always begin with the **start symbol** for the grammar
 - by tradition, the nonterminal on the LHS of the first rule is the start symbol
 - o.w. the start symbol could have a special name such as $\langle \text{start} \rangle$ or $\langle \text{program} \rangle$

Example: A derivation to show that:

$\text{id} + \text{num}$

is a valid expression ($\langle \text{expr} \rangle$):

$$\begin{aligned}\langle \text{expr} \rangle &\Rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \\ &\Rightarrow \text{id} \langle \text{op} \rangle \langle \text{expr} \rangle \\ &\Rightarrow \text{id} + \langle \text{expr} \rangle \\ &\Rightarrow \text{id} + \text{num}\end{aligned}$$

Example: A derivation to show that:

$\text{id} < (\text{num} / \text{num})$

is a valid expression ($\langle \text{expr} \rangle$):

$$\begin{aligned}\langle \text{expr} \rangle &\Rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \\ &\Rightarrow \text{id} \langle \text{op} \rangle \langle \text{expr} \rangle \\ &\Rightarrow \text{id} < \langle \text{expr} \rangle \\ &\Rightarrow \text{id} < (\langle \text{expr} \rangle) \\ &\Rightarrow \text{id} < (\langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle) \\ &\Rightarrow \text{id} < (\text{num} \langle \text{op} \rangle \langle \text{expr} \rangle) \\ &\Rightarrow \text{id} < (\text{num} / \langle \text{expr} \rangle) \\ &\Rightarrow \text{id} < (\text{num} / \text{num})\end{aligned}$$

Definitions:

- the symbol \Rightarrow is called **derives**
- each string of symbols derived from the start symbol (including the start symbol itself) is called a **sentential form**
- a **leftmost derivation** is a derivation in which the leftmost nonterminal is always chosen for replacement
- a **rightmost derivation** is a derivation in which the rightmost nonterminal is always chosen for replacement

Derivation order has no effect on the set of strings that can be derived.

BNF example: A subset of statements in C:

$\langle \text{stmt} \rangle \rightarrow \langle \text{if-stmt} \rangle \mid \langle \text{loop-stmt} \rangle \mid \langle \text{assign-stmt} \rangle \mid \langle \text{cmpd-stmt} \rangle$
 $\langle \text{if-stmt} \rangle \rightarrow \text{if } (\langle \text{expr} \rangle) \langle \text{stmt} \rangle$
 $\quad \mid \text{if } (\langle \text{expr} \rangle) \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle$
 $\langle \text{loop-stmt} \rangle \rightarrow \text{while } (\langle \text{expr} \rangle) \langle \text{stmt} \rangle$
 $\langle \text{assign-stmt} \rangle \rightarrow \text{id} = \langle \text{expr} \rangle ;$
 $\langle \text{cmpd-stmt} \rangle \rightarrow \{ \langle \text{stmt-list} \rangle \}$
 $\langle \text{stmt-list} \rangle \rightarrow \epsilon \mid \langle \text{stmt} \rangle \mid \langle \text{stmt} \rangle \langle \text{stmt-list} \rangle$
 $\langle \text{expr} \rangle \rightarrow \text{id} \mid \text{num} \mid (\langle \text{expr} \rangle) \mid \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle$
 $\langle \text{op} \rangle \rightarrow + \mid - \mid * \mid / \mid == \mid < \mid <= \mid > \mid >=$

Even more definitions:

- a BNF rule in which the nonterminal on the LHS is also the first (leftmost) symbol on the RHS is **left recursive**
- a BNF rule in which the nonterminal on the LHS is also the last (rightmost) symbol on the RHS is **right recursive**

For example:

- $\langle \text{stmt-list} \rangle \rightarrow \langle \text{stmt} \rangle \langle \text{stmt-list} \rangle$
is right recursive
- $\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle$
is both left and right recursive

3.3 Parse Trees

A **parse tree** is a graphical way of representing a derivation.

- the root of the parse tree is always the start symbol
- each interior node is a nonterminal
- each leaf node is a token
- the children of a nonterminal (interior node) are the RHS of some rule whose LHS is the nonterminal

For example, a parse tree for:

if (id > num) id = num; else { id = id + num; id = id; }

using the previous grammar is:

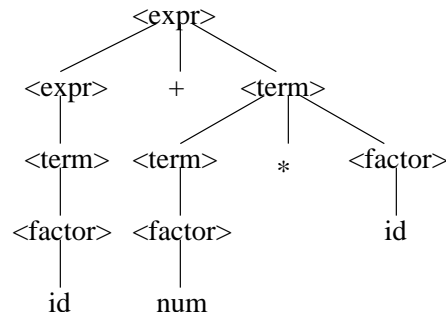
$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle \mid \langle \text{term} \rangle$
 $\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle * \langle \text{factor} \rangle \mid \langle \text{factor} \rangle$
 $\langle \text{factor} \rangle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

The intuition behind this approach is to try to force the + to occur higher in the parse tree.

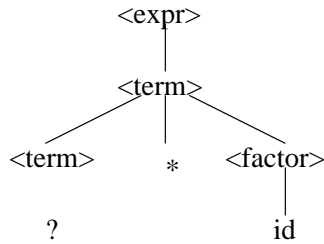
The parse tree for :

id + num * id

is:



Trying to find another tree:



This grammar modification gives rise to three proof obligations:

- the two grammars define the same language
- the second grammar always gives correct associativity and precedence
- the second grammar is not ambiguous

These proofs are omitted.

3.4 Extended BNF (EBNF)

EBNF provides some additional notation, but does not allow any additional languages to be defined (as compared to BNF). Hence, it is only a notational convenience.

notation	meaning
$[x]$	optional (0 or 1 occurrence of x)
$\{x\}$	repetition (0 or more occurrences of x)
$(x \mid y)$	choice (x or y)

Table 3: Additional EBNF notation.

Examples:

BNF: $\langle \text{stmt-list} \rangle \rightarrow \epsilon \mid \langle \text{stmt} \rangle \langle \text{stmt-list} \rangle$

EBNF: $\langle \text{stmt-list} \rangle \rightarrow \{ \langle \text{stmt} \rangle \}$

BNF: $\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid \langle \text{expr} \rangle * \langle \text{expr} \rangle$

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

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

EBNF: $\langle \text{expr} \rangle \rightarrow [\langle \text{expr} \rangle +] \langle \text{term} \rangle$

4 Static Semantics and Attribute Grammars

Static semantics is any aspect of semantics (meaning) that can be checked at compile time (statically). Examples:

- checking that variables are declared before they are used
- type checking

A more specific example: consider the Java code

```
String x = 5;
```

Question (1): is this an element of the language generated by the Java BNF?

Question (2): is this legal Java code?

As we can see from this example, many language properties can be checked statically, but are impossible or difficult to specify with a BNF definition.

An **attribute grammar** is a grammar with attributes (properties) attached to each symbol, plus rules for computing and checking the values of attributes. Attribute grammars are useful for checking static semantic properties.

If parse trees are implemented with a record or class instance for each node, then the attribute is just an additional field in this representation.

Types of attributes:

- a **synthesized attribute** is an attribute whose value is computed based on the attribute values of child nodes (in a parse tree). Synthesized attributes pass information up the parse tree.
- an **inherited attribute** is an attribute whose value is computed based on the attribute values of parent and sibling nodes (in a parse tree). Inherited attributes pass information down the parse tree.
- an **intrinsic attribute** is an attribute whose value is not based on other attribute values.

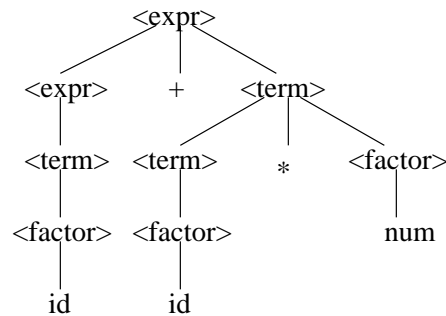
For example, consider the grammar:

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

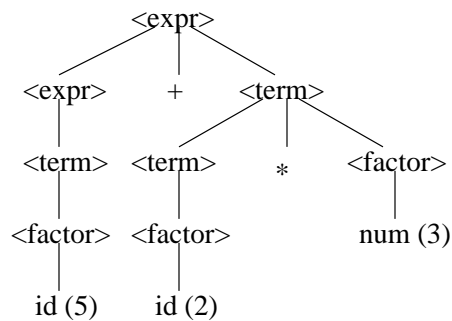
and expression:

$x + y * 3$

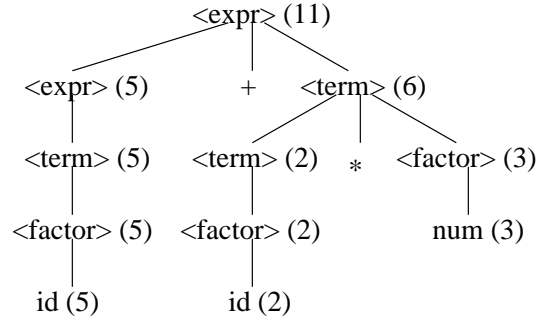
with parse tree:



Consider using an attribute named value to represent the mathematical result of an expression. If the value of x is 5 and y is 2, then the parse tree can be decorated/annotated with these intrinsic attribute values as follows:



Now synthesized attribute values can be computed for each of the nonterminals in the parse tree:



More practically, attributes can be used for type checking by adding a type attribute to each symbol. For example, consider the grammar:

```

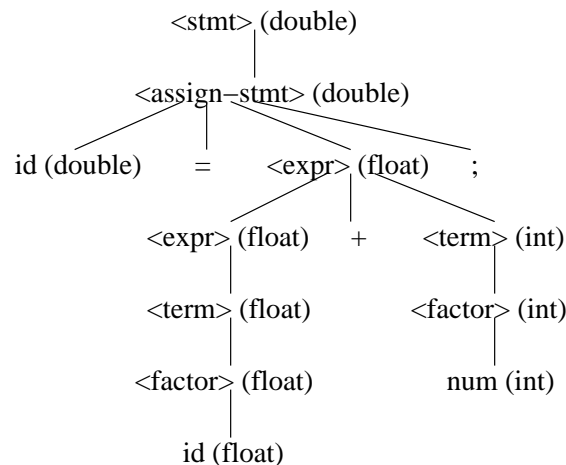
<stmt> → <assign-stmt>
<assign-stmt> → id = <expr> ;
<expr> → <expr> + <term> | <term>
<term> → <term> * <factor> | <factor>
<factor> → id | num | (<expr>)

```

and assignment statement:

```
x = y + 5;
```

and suppose that y is of type float and x is of type double. The annotated parse tree is:



This simple example shows only the *actual* types of the nodes in the parse tree, inferred as synthesized attributes. A compiler would also use *expected* types—a completely separate attribute, inferred as an inherited attribute. The textbook example 3.6 goes over this in detail. Once both kinds of types have been inferred (actual and expected), a compiler can check the resulting parse tree to see if the actual types are the same as the expected types. Of course, it will report an error if there are any invalid mismatches. But note that some kinds of mismatches are okay, because many languages allow *coercion* of types: for example, in Java, an int value can be coerced into a float or double without causing a compiler error.

Parser generators such as YACC and CUP generate parsers from attribute grammar specifications of the language being parsed.