

## BNF Grammars, Parse Trees, Ambiguity

A BNF (Backus-Naur Form) grammar is a triple  $\langle T, N, R \rangle$  where

- $T$  is a finite set of *terminal symbols*;
- $N$  is a finite set of *non-terminal symbols*;
- $R$  is a finite set of *production rules*.

When the BNF grammar is used for a formal description of lexical items ("tokens") and syntax of programming languages, terminal symbols are ASCII or Unicode characters and non-terminal symbols designate *syntactic categories*. In this course, the notation  $\langle X \rangle$  will be used for syntactic categories where  $X$  is a mnemonic identifier or phrase. Terminal symbols will be simply referred to as *terminals*. Each production rule takes the form:

$$\langle X \rangle \rightarrow \alpha_1 \mid \alpha_2 \mid \cdots \mid \alpha_n, \quad n \geq 1$$

where each  $\alpha_i$  is a string of terminals and/or syntactic categories;  $\alpha_i$  may be the empty string  $\epsilon$ .

**Example** Let  $T = \{ 0, 1, \dots, 9, a, b, \dots, z, A, B, \dots, Z, +, -, *, /, (, ), ". \}$ ,  $N = \{ \langle \text{digit} \rangle, \langle \text{letter} \rangle, \langle \text{int} \rangle, \langle \text{id} \rangle, \langle \text{rest} \rangle, \langle \text{float} \rangle, \langle \text{exponent} \rangle, \langle \text{eSign} \rangle, \langle \text{sign} \rangle, \langle E \rangle \}$ , and let  $R$  contain the following production rules:

```
 $\langle \text{digit} \rangle \rightarrow 0 \mid 1 \mid \cdots \mid 9$   
 $\langle \text{letter} \rangle \rightarrow a \mid b \mid \cdots \mid z \mid A \mid B \mid \cdots \mid Z$   
 $\langle \text{int} \rangle \rightarrow \langle \text{digit} \rangle \mid \langle \text{digit} \rangle \langle \text{int} \rangle$   
 $\langle \text{id} \rangle \rightarrow \langle \text{letter} \rangle \langle \text{rest} \rangle$   
 $\langle \text{rest} \rangle \rightarrow \epsilon \mid \langle \text{digit} \rangle \langle \text{rest} \rangle \mid \langle \text{letter} \rangle \langle \text{rest} \rangle$   
 $\langle \text{float} \rangle \rightarrow \langle \text{int} \rangle "." \langle \text{int} \rangle \mid \langle \text{int} \rangle "." \langle \text{int} \rangle \langle \text{exponent} \rangle$   
 $\langle \text{exponent} \rangle \rightarrow \langle \text{eSign} \rangle \langle \text{int} \rangle \mid \langle \text{eSign} \rangle \langle \text{sign} \rangle \langle \text{int} \rangle$   
 $\langle \text{eSign} \rangle \rightarrow "e" \mid "E"$   
 $\langle \text{sign} \rangle \rightarrow + \mid -$   
 $\langle E \rangle \rightarrow \langle \text{id} \rangle \mid \langle \text{int} \rangle \mid \langle \text{float} \rangle \mid \langle E \rangle + \langle E \rangle \mid \langle E \rangle - \langle E \rangle \mid \langle E \rangle * \langle E \rangle \mid \langle E \rangle / \langle E \rangle \mid "(" \langle E \rangle ")"$ 
```

This BNF defines three token categories  $\langle \text{int} \rangle$ ,  $\langle \text{id} \rangle$ ,  $\langle \text{float} \rangle$ ; the categories  $\langle \text{digit} \rangle$ ,  $\langle \text{letter} \rangle$ ,  $\langle \text{rest} \rangle$ ,  $\langle \text{exponent} \rangle$ ,  $\langle \text{eSign} \rangle$ , and  $\langle \text{sign} \rangle$  are not tokens but merely auxiliary categories to assist the definitions of  $\langle \text{int} \rangle$ ,  $\langle \text{id} \rangle$ ,  $\langle \text{float} \rangle$ . The category  $\langle E \rangle$  of arithmetic expressions is then defined from the three token categories, the arithmetic operators  $+$ ,  $-$ ,  $*$ ,  $/$ , and the parentheses.

Normally tokens are contiguous strings of non-whitespace characters. The usual assumption, therefore, is that the right side of any production rule for tokens may not be broken by whitespaces between any two terminals and/or syntactic categories. For example, we tacitly assume that no whitespace may appear between  $\langle \text{digit} \rangle$  and  $\langle \text{int} \rangle$  for  $\langle \text{int} \rangle$ , between  $\langle \text{letter} \rangle$  and  $\langle \text{rest} \rangle$  for  $\langle \text{id} \rangle$ , etc. (Obvious exceptions are character and string literals which may contain whitespaces, like ' ' and "a b"; in this case explicit syntactic categories are needed for symbols, including whitespaces, that can appear in character and string literals). Non-token syntactic categories, however, may have whitespaces between tokens/syntactic categories on the right sides of production rules. For example, we tacitly assume that there may be whitespaces between  $\langle E \rangle$  and  $+$ ,  $*$  and  $\langle E \rangle$ , "(" and  $\langle E \rangle$ , etc.

**Derivations** Let  $\beta \langle X \rangle \gamma$  be any string with an occurrence of  $\langle X \rangle$  where  $\beta$  and  $\gamma$  are strings of terminals and/or syntactic categories which may be  $\epsilon$ . Let  $\langle X \rangle \rightarrow \alpha_i$  be any choice from a production rule for  $\langle X \rangle$ . Then  $\beta \langle X \rangle \gamma$  is said to *derive*  $\beta \alpha_i \gamma$  in one step, formally denoted  $\beta \langle X \rangle \gamma \Rightarrow \beta \alpha_i \gamma$ .

**Example** The following is a sequence of one-step derivations of "CS316" from  $\langle id \rangle$  (the syntactic category being replaced is underlined):

$\langle id \rangle \Rightarrow$   
 $\langle letter \rangle \langle rest \rangle \Rightarrow$   
 $C \langle rest \rangle \Rightarrow$   
 $C \langle letter \rangle \langle rest \rangle \Rightarrow$   
 $CS \langle rest \rangle \Rightarrow$   
 $CS \langle digit \rangle \langle rest \rangle \Rightarrow$   
 $CS3 \langle rest \rangle \Rightarrow$   
 $CS3 \langle digit \rangle \langle rest \rangle \Rightarrow$   
 $CS31 \langle rest \rangle \Rightarrow$   
 $CS31 \langle digit \rangle \langle rest \rangle \Rightarrow$   
 $CS316 \langle rest \rangle \Rightarrow$   
 $CS316\epsilon =$   
 $CS316$

In this derivation sequence, the leftmost syntactic category is replaced in each step. In general any of the syntactic categories appearing in each step may be replaced. Let the notation  $\alpha \Rightarrow^* \beta$  mean  $\alpha$  derives  $\beta$  in a finite number of steps. Let  $\Gamma$  be a BNF grammar and  $\langle X \rangle$  any syntactic category in  $\Gamma$ . Then the *language defined by  $\langle X \rangle$  in  $\Gamma$*  is defined as the set:

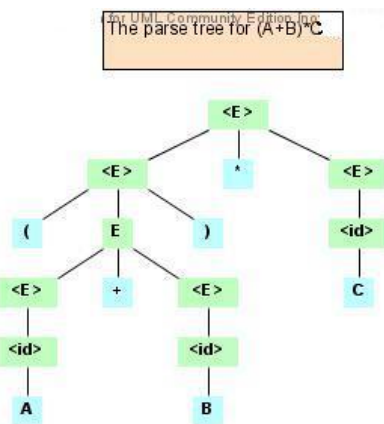
$$\{ \beta \mid \langle X \rangle \Rightarrow^* \beta, \beta \text{ is a string of terminals in } \Gamma \}$$

**Parse Trees** Linear derivation sequences are a useful concept for the formal definition of the languages defined by syntactic categories. They are also useful for deriving token strings, like "CS316", "76453", and "764.34E-5". However, it is very difficult to identify syntactic structures of higher-level, non-token strings in linear derivation sequences, such as the arithmetic expression "CS316+76453-764.34E-5/k+(a\*b)". A much better way of identifying and analyzing syntactic structures of strings is to construct derivations in tree form, called **parse trees**. A parse tree is defined as follows:

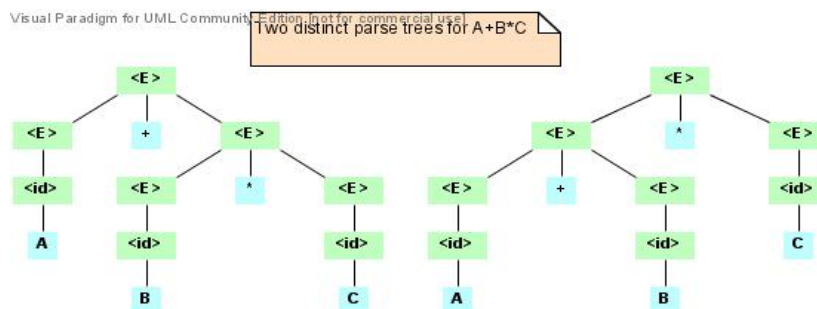
- Each non-leaf node is labeled by a syntactic category;
- Each leaf node is labeled by a terminal;
- Let  $n$  be any non-leaf node,  $\langle X \rangle$  the syntactic category labeling it. Then there must be a production rule  $\langle X \rangle \rightarrow \alpha_i$  such that the children of  $n$  are labeled by the terminals and syntactic categories in  $\alpha_i$  from left to right.

Parse trees can be constructed for tokens or strings of higher-level categories like  $\langle E \rangle$ , and their leaf nodes would be labeled by terminals according to this definition. However, actual compilers' parsers construct parse trees from streams of tokens and do not parse tokens themselves in tree form; here tokens are regarded as "high-level terminal symbols" extracted by lexical analyzers. In accordance with this, whenever constructing parse trees for non-token categories, we will derive tokens directly from their categories. All leaf nodes of such parse trees will then be labeled by tokens rather than terminals.

**Example** The following is the parse tree for "(A+B)\*C" from the category <E>:



**The Ambiguity Problem** A syntactic category <X> in a BNF grammar is said to be *ambiguous* if there exists a string  $\alpha$  of terminals/tokens such that there exist two or more parse trees for  $\alpha$  from <X>. The category <E> defined above is ambiguous as the following distinct parse trees for "A+B\*C" demonstrate:



The problem with ambiguity is that distinct parse trees for the same string do not in general represent equivalent semantic information. The first parse tree above represents  $A+(B*C)$  while the second represents  $(A+B)*C$ . Of course we would tend to assume the first tree is the "right" one based on the conventional precedence rules for the arithmetic operators, but both trees are valid according to the above production rule for <E> since it does not incorporate precedence rules in any way. The problem is clearly aggravated if we are looking at ambiguous syntactic categories for operators, control structures, etc. in a new programming language we are trying to learn – we can never be sure that expressions and statements we write down will be parsed as we intend. The same problem also occurs in compiler construction. The compiler builds parse trees for expressions/statements and then generates intermediate code based on the semantic information encoded in the trees. The first tree above would generate intermediate code to compute  $A+(B*C)$  while the second would generate code to compute  $(A+B)*C$ . So the compiler writer suffers the same problem: Which parse tree is supposed to be built? Clearly, formal descriptions of programming language syntax should use unambiguous BNF grammars or, failing that, ambiguous categories should be provided with separate disambiguation specifications (e.g., concise English comments).

The following is an unambiguous BNF grammar for <E>:

```
<E> → <term> | <term> + <E> | <term> - <E>
<term> → <primary> | <primary> * <term> | <primary> / <term>
<primary> → <id> | <int> | <float> | "(" <E> ")"
```

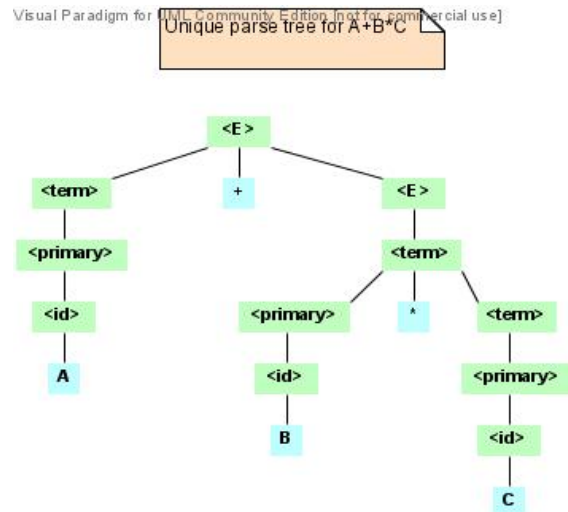
We can read these production rules informally as follows:

- An expression is a sequence of one or more terms separated by + or –.
- A term is a sequence of one or more primaries separated by \* or /.
- A primary is an id, int, float, or a parenthesized expression.

This grammar disambiguates expressions by two factors:

- it introduces a 3-level hierarchy that incorporates the usual precedence rule: \* and / have higher precedence than + and −;
- at the same precedence level, the operators associate to right.

A term cannot contain + or − unless it is enclosed in (…), and a primary cannot contain any operator unless it is enclosed in (…).



**Extended BNF (EBNF) Grammars** The extended BNF grammar provides concise notations for multiple choices and iterations that can be used on the right-hand sides of production rules. It does not, however, provide additional definitional power since any EBNF grammar can always be translated into an equivalent one in the original BNF.

- **Mandatory Multiple Choice:**  $\{\alpha_1 \mid \cdots \mid \alpha_n\}$  denotes a mandatory choice of exactly one  $\alpha_i$ .
- **Optional Multiple Choice:**  $[\alpha_1 \mid \cdots \mid \alpha_n]$  denotes an optional choice of at most one  $\alpha_i$ .
- **Zero or More Iterations:**  $\{\alpha\}^*$  denotes zero or more iterations of  $\alpha$ , i.e.,  $\epsilon, \alpha, \alpha\alpha, \alpha\alpha\alpha, \dots, \alpha^i, \dots$
- **One or More Iterations:**  $\{\alpha\}^+$  denotes one or more iterations of  $\alpha$ , i.e.,  $\alpha, \alpha\alpha, \alpha\alpha\alpha, \dots, \alpha^i, \dots$

$\{\alpha_1 \mid \cdots \mid \alpha_n\}$  and  $\{\alpha_1 \mid \cdots \mid \alpha_n\}^+$  may be abbreviated, respectively, to  $\{\alpha_1 \mid \cdots \mid \alpha_n\}$  and  $\{\alpha_1 \mid \cdots \mid \alpha_n\}^+$ .

### Example

```

<int> → { <digit> }+
<id> → <letter> { <letter> | <digit> }
<float> → { <digit> }+ "." { <digit> }+ [ (E|e) [+|-] { <digit> }+ ]
<E> → <term> { (+|-) <term> }
<term> → <primary> { (*|/) <primary> }
<primary> → <id> | <int> | <float> | "(" <E> ")"
  
```

Here the left or right associativity of the arithmetic operators must be commented separately, since the iterative definitions of <E> and <term>, unlike the recursive definitions, do not incorporate operator associativity.

**Meta Symbols** The following characters are called meta symbols (or special symbols) of the BNF/EBNF grammar since they have predefined, fixed meanings throughout all BNF/EBNF grammars: |, <, >, (, ), [, ], {, }. If any of these needs to be used as a terminal symbol, it must be enclosed in double quotes "..." to avoid confusion. Double quotes may also be used for any terminal symbols — especially easily missed terminals like "." and "," — to make them stand out clearly.

### Example

$\langle S \rangle \rightarrow \langle \text{id} \rangle = \langle E \rangle ; \mid \text{if } "(" \langle B \rangle ")" \langle S \rangle [ \text{else } \langle S \rangle ] \mid \text{while } "(" \langle B \rangle ")" \langle S \rangle \mid \text{do } \langle S \rangle \text{ while } "(" \langle B \rangle "$   
 $" ) " \mid \{ " \langle S \text{ List} \rangle " \}$   
 $\langle S \text{ List} \rangle \rightarrow \{ \langle S \rangle \}^+$

Note: In the 2-branch conditionals, each **else** matches with the closest preceding unmatched **if**.