# Lecture 2

# Describing Syntax

## Chapter 3 – Part 1

First Semester
1447 - 2025

CONCEPTS OF
PROGRAMMING LANGUAGES

ROBERT W. SEBESTA

12/E

# Lecture 2 Topics:

- **Introduction**
- **The General Problem of Describing Syntax:**
  - Language **Recognizers**
  - Language **Generators**
- **Formal Methods of Describing Syntax:**
  - Backus-Naur Form (**BNF**):
    - Extended BNF (**EBNF**)
  - Context-Free **Grammars** (**CFG**):
    - Grammars
    - Derivations
    - Parse Trees
    - Ambiguity
- **Attribute Grammars:**
  - Static Semantics

# Introduction

- The study of <u>programming languages</u>, like the study of <u>natural languages</u>, can be divided into:
  - Examinations of syntax.
  - Examinations of semantics.

- Syntax: the form or structure of the expressions, statements, and program units.

- Semantics: the meaning of the expressions, statements, and program units.

- Syntax and semantics provide <u>a language's definition</u>.

# Introduction: Example

- ## Example:

  - ### The syntax of a Java "**while**" statement is:

    - ➤ while (boolean_expr) statement

  - ### The semantics of the same statement is:

    - ➤ When the current value of the Boolean expression is true, the embedded statement is executed.

    - ➤ Then control implicitly returns to the Boolean expression to repeat the process.

    - ➤ If the Boolean expression is false, control transfers to the statement following the while construct.

# Introduction: Language Users

- Users of a language definition:
  - Other language designers (evaluators).
  - Implementers.
  - Programmers (the users of the language).

# The General Problem of Describing Syntax: Terminology

- A *sentence* (statement) is a string of characters over some alphabet.
- A *language* is a set of sentences (statements).
- A *lexeme* is the lowest level syntactic unit of a language (e.g., `*`, `sum, begin`).
  - The language operators.
  - The language special words.
  - The language numerical literals.
  - Etc.
- A *token* is a category of lexemes (e.g., identifier).

# Example

- **Consider the following Java statement:**
  - ➢ `index = 2 * count + 17;`

| Lexemes | Tokens |
|---------|--------|
| index | identifier |
| = | equal_sign |
| 2 | int_literal |
| * | mult_op |
| count | identifier |
| + | plus_op |
| 17 | int_literal |
| ; | semicolon |

# Formal Definition of Languages

- ## Recognizers:

  - A recognition <u>device</u> reads input strings over the alphabet of the language and decides whether the input strings belong to the language (accept or reject the given input strings).

  - Example: syntax analysis (parsing) part of a compiler.

    - The **syntax analyzer** <u>determines whether the given programs are syntactically correct</u>.

  - Detailed discussion of syntax analysis appears in **Chapter 4.**

- ## Generators:

  - A <u>device</u> that generates sentences of a language.

  - One can determine if the syntax of a particular sentence is syntactically correct by comparing it to the structure of the generator.

# BNF and Context-Free Grammars

- **Context-Free Grammars:**
  - Developed by Noam Chomsky in the mid-1950s.
  - Language generators, meant to <u>describe the syntax of natural languages</u>.
  - Define a class of languages called context-free languages.

- **Backus-Naur Form (1959):**
  - Invented by John Backus to <u>describe the syntax of Algol 58</u>.
  - BNF is <u>equivalent</u> to context-free grammars.

# BNF Fundamentals

- BNF is a natural <u>notation</u> for <u>describing</u> syntax.
- In BNF, <u>abstractions</u> are used to <u>represent classes of syntactic structures</u>.
  - They act like syntactic variables (also called *nonterminal symbols,* or just *terminals*).
- Example:
  - A simple Java <u>assignment statement</u> might be represented by the **abstraction** <**assign**>
  - Pointed brackets are often used to delimit names of abstractions.
  - The actual definition of <**assign**> can be given by:
  - `<assign>` → `<var> = <expression>`

# BNF Fundamentals (continued)

- Terminals are lexemes or tokens.

- A rule has a left-hand side (LHS), which is a nonterminal, and a right-hand side (RHS), which is a string of terminals and/or nonterminals.

- Example:

  - See the next slide!

# $\langle\text{assign}\rangle \rightarrow \langle\text{var}\rangle = \langle\text{expression}\rangle$

- The text on the <u>left side</u> of the <u>arrow</u>, which is aptly called the left-hand side (LHS), is the abstraction being defined.

- The text to the <u>right of the arrow</u> is the definition of the LHS.

- It is called the right-hand side (RHS) and consists of some mixture of **tokens**, **lexemes**, and **references to** other abstractions. (Actually, **tokens** are also abstractions.)

- Altogether, the definition is called a rule, or production.

- In the example rule just given, the abstractions $\langle\text{var}\rangle$ and $\langle\text{expression}\rangle$ obviously must be defined for the $\langle\text{assign}\rangle$ definition to be useful.

# How a rule can be read?

$$\langle \textbf{assign} \rangle \rightarrow \langle \textbf{var} \rangle = \langle \textbf{expression} \rangle$$

- This particular rule specifies that the abstraction $\langle$assign$\rangle$ is defined as an <u>instance</u> of the abstraction $\langle$var$\rangle$, <u>followed</u> by the lexeme =, <u>followed</u> by an <u>instance</u> of the abstraction $\langle$expression$\rangle$.

- One example sentence whose <u>syntactic structure</u> is described by the rule above is:

  - ```
    total = subtotal1 + subtotal2
    ```

# BNF Fundamentals (continued)

- **Nonterminals** are often enclosed in "*angle brackets*" (abstraction).

  - Examples of BNF rules:

    ```
    <ident_list> → identifier | identifier, <ident_list>

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

- **Grammar**: a finite non-empty set of rules.

- A *start symbol* is a <u>special element</u> of the nonterminals of a grammar.

# BNF Rules

- An abstraction (or nonterminal symbol) can have <u>more than one</u> RHS:

  ```
  <stmt> → <single_stmt>
  <stmt> → begin <stmt_list> end
  ```

- These two rules can be written as:

```
<stmt> → <single_stmt> | begin <stmt_list> end
```

# BNF Rules: More Examples

- a Java if statement can be described with the rules:

<if_stmt> → if ( <logic_expr> ) <stmt>

<if_stmt> → if ( <logic_expr> ) <stmt> else <stmt>

or with the rule

<if_stmt> → if ( <logic_expr> ) <stmt>

    | if ( <logic_expr> ) <stmt> else <stmt>

# Describing Lists

- Example of a list:

  - a list of identifiers appearing on a data declaration statement.

- Syntactic lists are described using recursion:

  ```
  <ident_list> → ident

               | ident, <ident_list>
  ```

- A rule is recursive if its LHS appears in its RHS.

# Derivation

- A grammar is a generative device <u>for</u> <u>defining languages</u>.

- A derivation is a <u>repeated</u> application of rules (grammars), <u>starting</u> with the start symbol and <u>ending</u> with a sentence (all terminal symbols).

- The start symbol represents a complete program and is often named <program>.

# An Example Grammar

`<program>` → `<stmts>`

`<stmts>` → `<stmt>` | `<stmt>` ; `<stmts>`

`<stmt>` → `<var>` **=** `<expr>`

`<var>` → **a** | **b** | **c** | **d**

`<expr>` → `<term>` **+** `<term>` | `<term>` **-** `<term>`

`<term>` → `<var>` | const

- The **language** described by the **grammar** of has only <u>one</u> statement form:
  - assignment.

# An Example Derivation

- How can we **derive** the following **assignment** statement using the previous **grammars** (rules)?

  - `a = b + const`

> The symbol => is read "derives."

```
<program> => <stmts>
            => <stmt>
            => <var> = <expr>
            => a = <expr>
            => a = <term> + <term>
            => a = <var> + <term>
            => a = b + <term>
            => a = b + const
```

# Another Example Grammar

## A Grammar for a Small Language

&lt;program&gt; → **begin** &lt;stmt_list&gt; **end**

&lt;stmt_list&gt; → &lt;stmt&gt;
       | &lt;stmt&gt; ; &lt;stmt_list&gt;

&lt;stmt&gt; → &lt;var&gt; = &lt;expression&gt;

&lt;var&gt; → A | B | C

&lt;expression&gt; → &lt;var&gt; + &lt;var&gt;
       | &lt;var&gt; - &lt;var&gt;
       | &lt;var&gt;

How can we **derive** the following small program?

```
begin
A = B + C ;
B = C
end
```

# Another Example Derivation

```
<program> => begin <stmt_list> end
        => begin <stmt> ; <stmt_list> end
        => begin <var> = <expression> ; <stmt_list> end
        => begin A = <expression> ; <stmt_list> end
        => begin A = <var> + <var> ; <stmt_list> end
        => begin A = B + <var> ; <stmt_list> end
        => begin A = B + C ; <stmt_list> end
        => begin A = B + ; <stmt> end
        => begin A = B + C ; <var> = <expression> end
        => begin A = B + C ; B = <expression> end
        => begin A = B + C ; B = <var> end
        => begin A = B + C ; B = C end
```

# Derivations

- Every string of symbols in a derivation is a *sentential form.*

- A *sentence* (statement) 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 <u>neither</u> leftmost <u>nor</u> rightmost.

# Yet Another Example Grammar

## A Grammar for Simple Assignment Statements

$$\langle assign \rangle \rightarrow \langle id \rangle = \langle expr \rangle$$
$$\langle id \rangle \rightarrow A \mid B \mid C$$
$$\langle expr \rangle \rightarrow \langle id \rangle + \langle expr \rangle$$
$$\mid \langle id \rangle * \langle expr \rangle$$
$$\mid ( \langle expr \rangle )$$
$$\mid \langle id \rangle$$

How can we **derive** the following **assignment** statement?

$$A = B * ( A + C )$$

# Yet Another Example Derivation

- The statement A = B * ( A + C ) is generated by the following leftmost derivation:

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

# Parse Tree

- A hierarchical representation of a derivation.

A **parse tree** for the simple statement
`a = b + const`

See slide 20

```
                      <program>
                          |
                      <stmts>
                          |
                       <stmt>
                      /   |   \
                 <var>    =    <expr>
                   |          /   |   \
                   a    <term>  +  <term>
                          |            |
                        <var>        const
                          |
                          b
```

# Parse Tree



A **parse tree** for the simple statement
`A = B * (A + C)`

# Ambiguity in Grammars

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

- This type of grammar allows the parse tree of an expression to grow on both left and right.
    - It should allow the tree to grow on the right only in such cases.

- How can this be a problem?
    - It confuses compilers during syntax analysis as compilers use the parse tree to generate code.
    - So, the meaning of the structure cannot be determined uniquely.

# An Ambiguous Expression Grammar

<expr> → <expr> <op> <expr>  |  const

<op> → /  |  -

# Another Ambiguous Grammar

## An Ambiguous Grammar for Simple Assignment Statements

$\langle assign \rangle \rightarrow \langle id \rangle = \langle expr \rangle$

$\langle id \rangle \rightarrow A \mid B \mid C$

$\langle expr \rangle \rightarrow \langle expr \rangle + \langle expr \rangle$
$\quad\quad\quad \mid \langle expr \rangle * \langle expr \rangle$
$\quad\quad\quad \mid ( \langle expr \rangle )$
$\quad\quad\quad \mid \langle id \rangle$

- This grammar is ambiguous because the sentence A = B + C * A has <u>two</u> *distinct* **parse trees**.

  - See the figure in the next slide.

**Two** *distinct* **parse trees** for the <u>same</u> sentence, A = B + C * A

# An Unambiguous Expression Grammar

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

```
<expr> → <expr> - <term>  |  <term>
<term> → <term> / const | const
```

# Operators Precedence

- Given the expression `x + y * z`, one obvious semantic_issue is the order of evaluation of the two operators.

  - Is it add and then multiply, or vice versa?

- This semantic question can be answered by <u>assigning different precedence levels to operators</u>.

- As grammar can describe a <u>certain syntactic structure</u> so that part of the meaning of the structure can be <u>determined from</u> its parse tree.

# Operators Precedence

- It is a fact that an operator in an arithmetic expression is generated lower in the parse tree must be evaluated first.

- This fact can be used to indicate that a lower operator in a parse tree has precedence over an operator produced higher up in the tree.

- However, *is this fact sufficient to solve the operator precedence problem?*

  - Not always! See the unambiguous grammar in slide 24.

# Operators Precedence

- For example, using the grammar in slide 24, try to sketch the parse trees for these two expressions:
  - `A + B * C`
  - `A * B + C`
- **What have you noticed?**
- You will see that:
  - For `A + B * C`, the (`*`) operator is the <u>lowest</u> in the tree, which will lead to a correct evaluation.
  - However, for `A * B + C` instead, the (`+`) operator is the lowest (indicating it is to be done first), which will lead to an incorrect evaluation.
- So, the grammar (slide 24) is sensitive to the order of the operators in the expressions.

# Operators Precedence

- So, how this problem can be solved?
  - Simply, take the order into consideration when designing the grammar by:
    - Use separate <u>nonterminal</u> symbols to represent the operands of the operators that have different precedence.
    - This requires additional <u>nonterminals</u> and some new rules.
  - For example, to correct the grammar in slide 24, we could use <u>three nonterminals</u> to represent operands, which <u>allows the grammar to force different operators to different levels</u> in the parse tree.
  - But, how? See the next slide!

# Operators Precedence

- If \<expr> is the <u>root symbol</u> for expressions, + can be forced to the <u>top</u> of the parse tree by having \<expr> directly generate only + operators, using the <u>new nonterminal</u>, \<term>, as the <u>right operand</u> of +.

- Next, we can define \<term> to generate * operators, using \<term> as the <u>left operand</u> and a <u>new nonterminal</u>, \<factor>, as its <u>right operand</u>.

- Now, * will always be <u>lower</u> in the parse tree, simply <u>because it is farther from the start symbol</u> than + in every derivation.

## A Grammar for Simple Assignment Statements

<assign> → <id> = <expr>

<id> → A | B | C

<expr> → <id> + <expr>

      | <id> * <expr>

      | ( <expr>)

      | <id>

## An Unambiguous Grammar for Expressions

<assign> → <id> = <expr>

<id> → A | B | C

<expr> → <expr> + <term>

      | <term>

<term> → <term> * <factor>

      | <factor>

<factor> → ( <expr> )

      | <id>

This grammar generates the same language as the above grammar. It is unambiguous and it specifies the usual precedence order of multiplication and addition operators.

# Operators Precedence: Example (Leftmost Derivation)

A = B + C * A



```
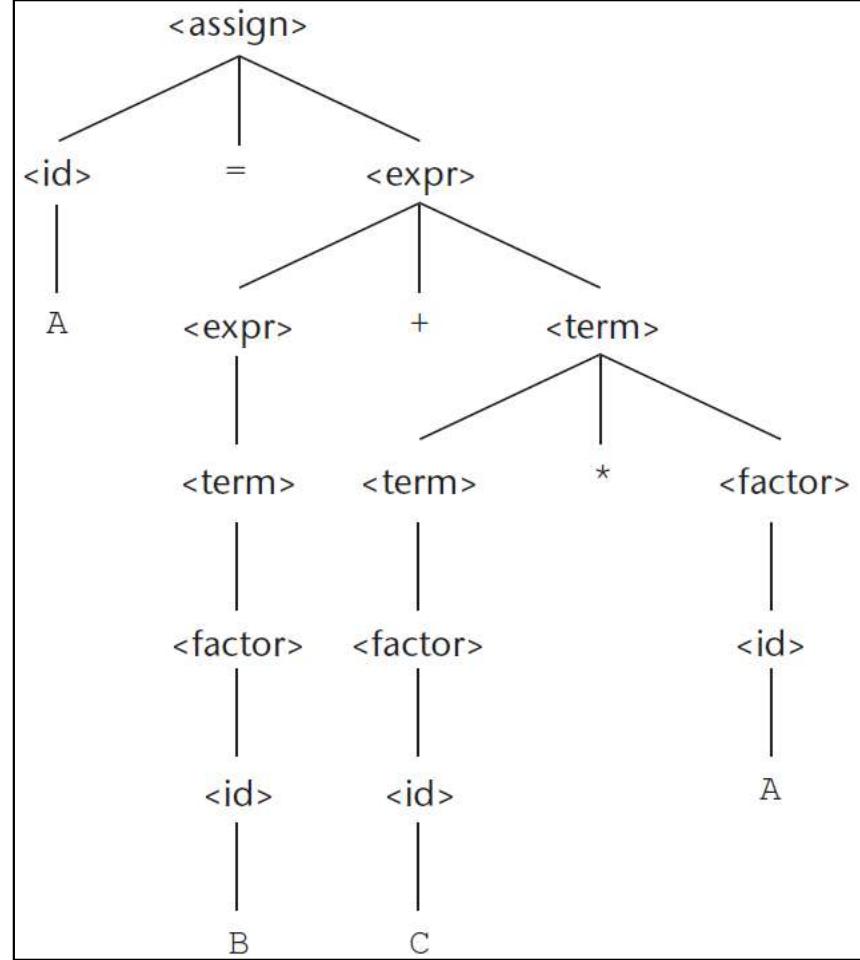<assign> => <id> = <expr>
        => A = <expr>
        => A = <expr> + <term>
        => A = <term> + <term>
        => A = <factor> + <term>
        => A = <id> + <term>
        => A = B + <term>
        => A = B + <term> * <factor>
        => A = B + <factor> * <factor>
        => A = B + <id> * <factor>
        => A = B + C * <factor>
        => A = B + C * <id>
        => A = B + C * A
```

The unique **parse tree** for A = B + C * A using an unambiguous grammar

# Operators Precedence: Example (Rightmost Derivation)

A = B + C * A

```
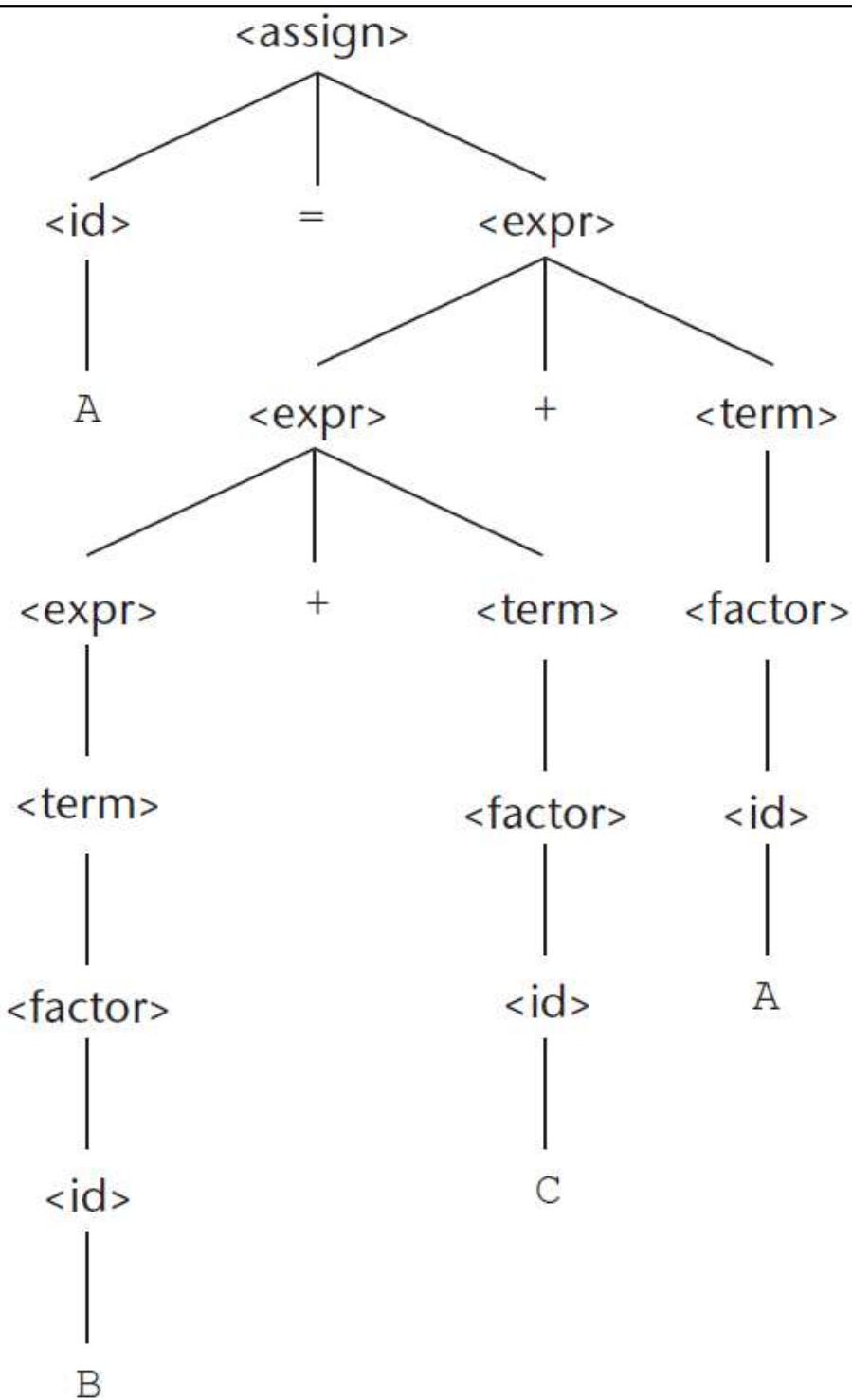<assign> => <id> = <expr>
        => <id> = <expr> + <term>
        => <id> = <expr> + <term> * <factor>
        => <id> = <expr> + <term> * <id>
        => <id> = <expr> + <term> * A
        => <id> = <expr> + <factor> * A
        => <id> = <expr> + <id> * A
        => <id> = <expr> + C * A
        => <id> = <term> + C * A
        => <id> = <factor> + C * A
        => <id> = <id> + C * A
        => <id> = B + C * A
        => A = B + C * A
```

# Associativity of Operators

- When an <u>expression</u> includes two operators that have the <u>same precedence</u> (as `*` and `/` usually have)—for example, "`A / B * C`", then a <u>semantic rule</u> is required to specify which should have precedence.

- This rule is named associativity.

- A grammar for <u>expressions</u> may correctly imply operator associativity.

- Consider the following example of an assignment statement:

  - `A = B + C + A`

  - After using the grammar in the <u>next slide</u> for the derivation of this statement, then its parse tree will look like:

An Unambiguous Grammar for Expressions

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

A **parse tree** for
`A = B + C + A`
illustrating the
associativity of addition

# Associativity of Operators

- The previous parse tree shows the left addition operator lower than the right addition operator.
  - This is the correct order if addition is meant to be left associative, which is typical.
  - In most cases, the associativity of addition in a computer is irrelevant.
- In mathematics, addition is associative, which means that left and right associative orders of evaluation mean the same thing.
  - That is, `(A + B) + C = A + (B + C)`
- Subtraction and division are not associative, whether in mathematics or in a computer.
  - Therefore, correct associativity may be essential for an expression that contains either of them.

# Associativity of Operators

- When a grammar rule has <u>its LHS also appearing at the <u>beginning</u> of its RHS</u>, the rule is said to be left recursive.
  - This left recursion specifies <u>left associativity</u>.
- For **example**, the left recursion of the rules of the grammar below <u>causes it to make both addition and multiplication left associative</u>.

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

# Associativity of Operators

- The exponentiation operator is <u>right associative</u> in most languages that provide it.
- To indicate right associativity, <u>right recursion</u> can be used.
- A grammar rule is right recursive if the LHS appears at the right end of the RHS.
- Rules such as:

$$\begin{aligned}
\langle factor \rangle &\rightarrow \langle exp \rangle ** \langle factor \rangle \\
&\mid \langle exp \rangle \\
\langle exp \rangle &\rightarrow ( \langle expr \rangle ) \\
&\mid id
\end{aligned}$$

could be used to describe exponentiation as a right-associative operator.

# Extended BNF

- **Optional** parts are placed in <u>brackets</u> [ ]

  `<proc_call> → ident [(<expr_list>)]`

- **Alternative** parts of RHSs are placed <u>inside</u> <u>parentheses</u> and <u>separated</u> via <u>vertical bars</u>.

  `<term> → <term> (+|-) const`

- **Repetitions** (zero or more) are placed inside braces { }

  `<ident> → letter {letter|digit}`

# BNF and EBNF

- **BNF**:

  <expr> → <expr> + <term>

      | <expr> - <term>

      | <term>

  <term> → <term> * <factor>

      | <term> / <factor>

      | <factor>

- **EBNF**:

  <expr> → <term> {(+ | -) <term>}

  <term> → <factor> {(* | /) <factor>}

# Recent Variations in EBNF

- **Alternative** RHSs are put on <u>separate lines</u>.

- Use of a **colon** instead of `=>`

- Use of `opt` for **optional** parts.

- Use of `oneof` for **choices**.

# Any Questions?

- Please, read chapter 3 (*first 2 sections*)
- I hope you were taking some notes!
- To test your understanding of this lecture, have a go with the "Review Questions" in page 156 of the textbook.
- We will do more exercises later on.
- Please, keep reviewing this lecture regularly.
- We may have a quiz (lecture 1) next week!
- Please, start doing your assignment.