
Describing Syntax and Semantics

CS 315 – Programming Languages

Pinar Duygulu

Bilkent University

Introduction

Providing a **precise description** of a programming language is important.

Reasons:

- Diversity of the people who need to understand
- Language implementors must determine how the expressions, statements, etc are formed, and their intended effects – clear description of language make their job easy
- Language users must understand the language by referring to the language manual

ALGOL 60 was the **first** language with a precise description.

Introduction

Syntax of a PL: the **form** of its expressions, statements, and program units.

Semantics of a PL: the **meaning** of those expressions, statements, and program units

e.g: while statement in Java

syntax: while (<boolean_expr>) <statement>

semantics: when boolean_expr is true it will be executed

The meaning of a statement should be clear from its syntax

The general problem of describing syntax

Language: a set of strings of characters from some alphabet.

Natural Languages/ Programming Languages/Formal Languages

Ex: English, Turkish / Pascal, C, FORTRAN / a^*b^* , 0^n1^n

Strings of a language: sentence / program (statement) / word

Alphabet: Σ , All strings: Σ^* , Language: $L \subseteq \Sigma^*$

Syntax rules specify which strings from Σ^* are in the language.

Lexemes

Lower level constructs are given not by the syntax but by lexical specifications. These are called **lexemes**

Examples: identifiers, constants, operators, special words.

`total, sum_of_products, 1254, ++, (:`

So, a language is considered as a set of strings of lexemes rather than strings of chars.

Tokens

- A **token** of a language is a category of its lexemes.
- For example, **identifier** is a token which may have lexemes `sum` `and` `total`

Example in Java language

x = (y+3.1) * z_5;

Lexemes

Tokens

x

identifier

=

equal_sign

(

left_paren

)

right_paren

for

for

y

identifier

+

plus_op

3.1

float_literal

mult_op

z_5

identifier

;

semi_colon

Describing Syntax

- Higher level constructs are given by **syntax rules**.
- Examples: organization of the program, loop structures, assignment, expressions, subprogram definitions, and calls.

Elements of Syntax

- An alphabet of symbols
- Symbols are terminal and non-terminal
 - Terminals cannot be broken down
 - Non-terminals can be broken down further
- Grammar rules that express how symbols are combined to make legal sentences
- Rules are of the general form
non-terminal symbol ::= list of zero or more terminals or non-terminals
- One uses rules to recognize (parse) or generate legal sentences

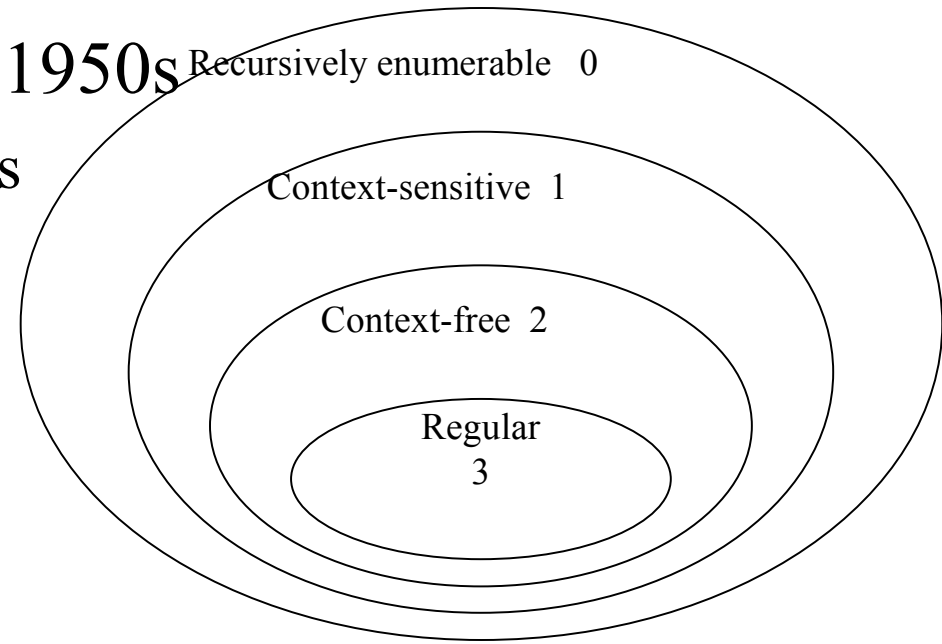
Recognizers vs Generators

- Automata (accept or reject) if input string is in the language
- Grammars (set of rules) easy to understand by humans

Formal Methods for Describing Syntax

- **Noam Chomsky** – linguist - 1950s
 - define four classes of languages

Programming languages
are contained in the class of
CFL's.



ALGOL58 John Backus

ALGOL 60 Peter Naur

Backus-Naur form: A notation to describe
the syntax of programming languages.

Fundamentals

- A **metalanguage** is a language used to describe another language.
- **BNF** (Backus-Naur Form) is a metalanguage used to describe PL's.
- BNF uses abstractions for syntactic structures.
- $\langle \text{LHS} \rangle \rightarrow \langle \text{RHS} \rangle$
- **LHS**: abstraction being defined
- **RHS**: definition
- Note: Sometimes $::=$ is used for \rightarrow

Fundamentals

- Example, Java assignment statement can be represented by the abstraction `<assign>`. Then the assignment statement of Java can be defined in BNF as
- `<assign> → <var> = <expression>`
- Such a definition is called a **rule** or **production**.
- Here, `<var>` and `<expression>` must also be defined.
- an instance of this abstraction can be
total = sub1 + sub2

Fundamentals

- These abstractions are called **Variables**, or **Nonterminals** of a Grammar.
- Grammar is simply a collection of rules.
- **Lexemes** and **tokens** are the **Terminals** of a grammar.

*An initial example

- Consider the sentence
 - “Marry greets John”
- A simple grammar for it
 - $\langle \text{sentence} \rangle ::= \langle \text{subject} \rangle \langle \text{predicate} \rangle$
 - $\langle \text{subject} \rangle ::= \text{Mary}$
 - $\langle \text{predicate} \rangle ::= \langle \text{verb} \rangle \langle \text{object} \rangle$
 - $\langle \text{verb} \rangle ::= \text{greets}$
 - $\langle \text{object} \rangle ::= \text{John}$

*Alternation

- Multiple definitions can be separated by | to mean **OR**.

`<object> ::= John | Alfred`

This adds “Marry greets Alfred” to legal sentences

`<subject> ::= Marry | John`

`<object> ::= Marry | John`

Alternatively

`<sentence> ::= <subject><predicate>`

`<subject> ::= noun`

`<predicate> ::= <verb><object>`

`<verb> ::= greets`

`<object> ::= <noun>`

`<noun> ::= John | Mary`

*Infinite number of Sentences

$\langle \text{object} \rangle ::=$

 John |

 John again |

 John again and again |

Instead use recursive definition

$\langle \text{object} \rangle ::=$ John |

 John $\langle \text{repeat factor} \rangle$

 $\langle \text{repeat factor} \rangle ::=$ again |

 again and $\langle \text{repeat factor} \rangle$

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

*Simple example for PLs

$$\langle \text{number} \rangle ::= \langle \text{number} \rangle \langle \text{digit} \rangle \mid \langle \text{digit} \rangle$$
$$\langle \text{signed number} \rangle ::= + \langle \text{number} \rangle \mid - \langle \text{number} \rangle$$

*Simple example for PLs

- How you can describe simple arithmetic?
- $\langle \text{expression} \rangle ::= \langle \text{expr} \rangle \langle \text{operator} \rangle \langle \text{expr} \rangle \mid \text{var}$
- $\langle \text{op} \rangle ::= + \mid - \mid * \mid /$
- $\langle \text{var} \rangle ::= a \mid b \mid c \mid \dots$
- $\langle \text{var} \rangle ::= \langle \text{signed number} \rangle$

*All numbers

- $S := '-' FN \mid FN$
- $FN := DL \mid DL '.' DL$
- $DL := D \mid D DL$
- $D := '0' \mid '1' \mid '2' \mid '3' \mid '4' \mid '5' \mid '6' \mid '7' \mid '8' \mid '9'$

*Identifiers

```
<identifier> → <letter>  
| <identifier><letter>  
| <identifier><digit>
```

PASCAL/Ada If Statement

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

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

or

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

Grammars and Derivations

- A grammar is a generative device for defining languages
- The **sentences** of the language are **generated** through a sequence of applications of the rules, starting from the special nonterminal called **start symbol**.
- Such a generation is called a **derivation**.
- **Start symbol** represents a **complete program**. So it is named **<program>**.

Example grammar

```
<program> → begin <stmt_list> end  
<stmt_list> → <stmt>  
             | <stmt> ; <stmt_list>  
<stmt>      → <var> := <expression>  
<var>       → A | B | C  
<expression> → <var>  
             | <var><arith_op> <var>  
<arith_op>  → + | - | * | /
```


Derivations

- In order to **check if** a given string represents a **valid program** in the language, we **try to derive** it in the grammar.
- Example string:
 - **begin A := B; C := A * B end;**
- Derivation **starts** from the **start symbol** **<program>**.
- At each step we replace a nonterminal with its definition (RHS of the rule).

Example

Each of
these strings
is called
**sentential
form**

`<program>` \Rightarrow **begin** `<stmt_list>` **end**
 \Rightarrow **begin** `<stmt>` ; `<stmt_list>` **end**
 \Rightarrow **begin** `<var>` := `<expression>`; `<stmt_list>` **end**
 \Rightarrow **begin** **A** := `<expression>`; `<stmt_list>` **end**
 \Rightarrow **begin** **A** := **B**; `<stmt_list>` **end**
 \Rightarrow **begin** **A** := **B**; `<stmt>` **end**
 \Rightarrow **begin** **A** := **B**; `<var>` := `<expression>` **end**
 \Rightarrow **begin** **A** := **B**; **C** := `<expression>` **end**
 \Rightarrow **begin** **A** := **B**; **C** := `<var>``<arith_op>``<var>` **end**
 \Rightarrow **begin** **A** := **B**; **C** := **A** `<arith_op>` `<var>` **end**
 \Rightarrow **begin** **A** := **B**; **C** := **A** * `<var>` **end**
 \Rightarrow **begin** **A** := **B**; **C** := **A** * **B** **end**

If always the leftmost nonterminal is replaced, then it is called leftmost derivation.

Another example

$\langle \text{assign} \rangle ::= \langle \text{id} \rangle = \langle \text{exp} \rangle$

$\langle \text{id} \rangle ::= A \mid B \mid C$

$\langle \text{expr} \rangle ::= \langle \text{id} \rangle + \langle \text{exp} \rangle$

$\mid \langle \text{id} \rangle * \langle \text{expr} \rangle$

$\mid (\langle \text{expr} \rangle)$

$\mid \langle \text{id} \rangle$

Derivation

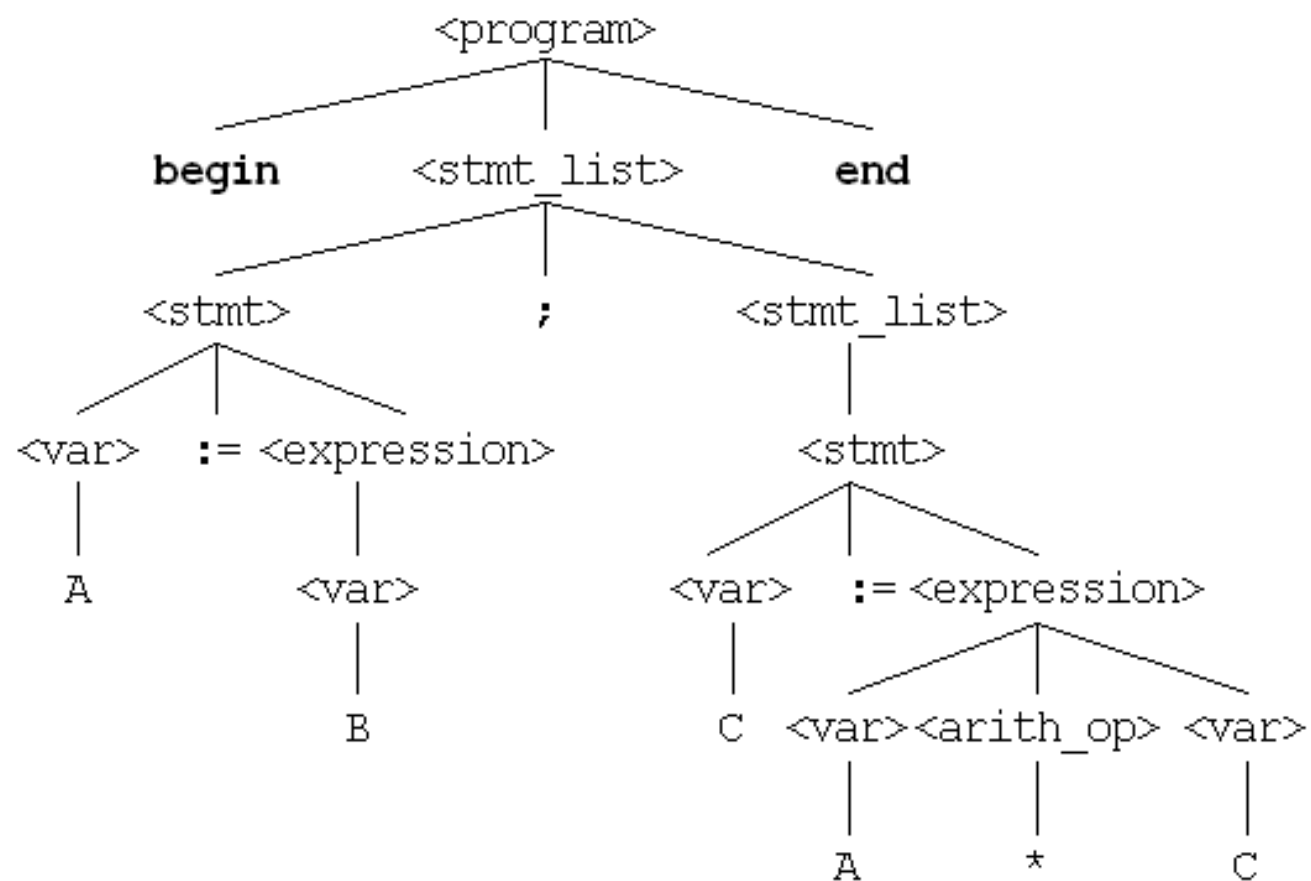
$$A = B * (A + C)$$

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

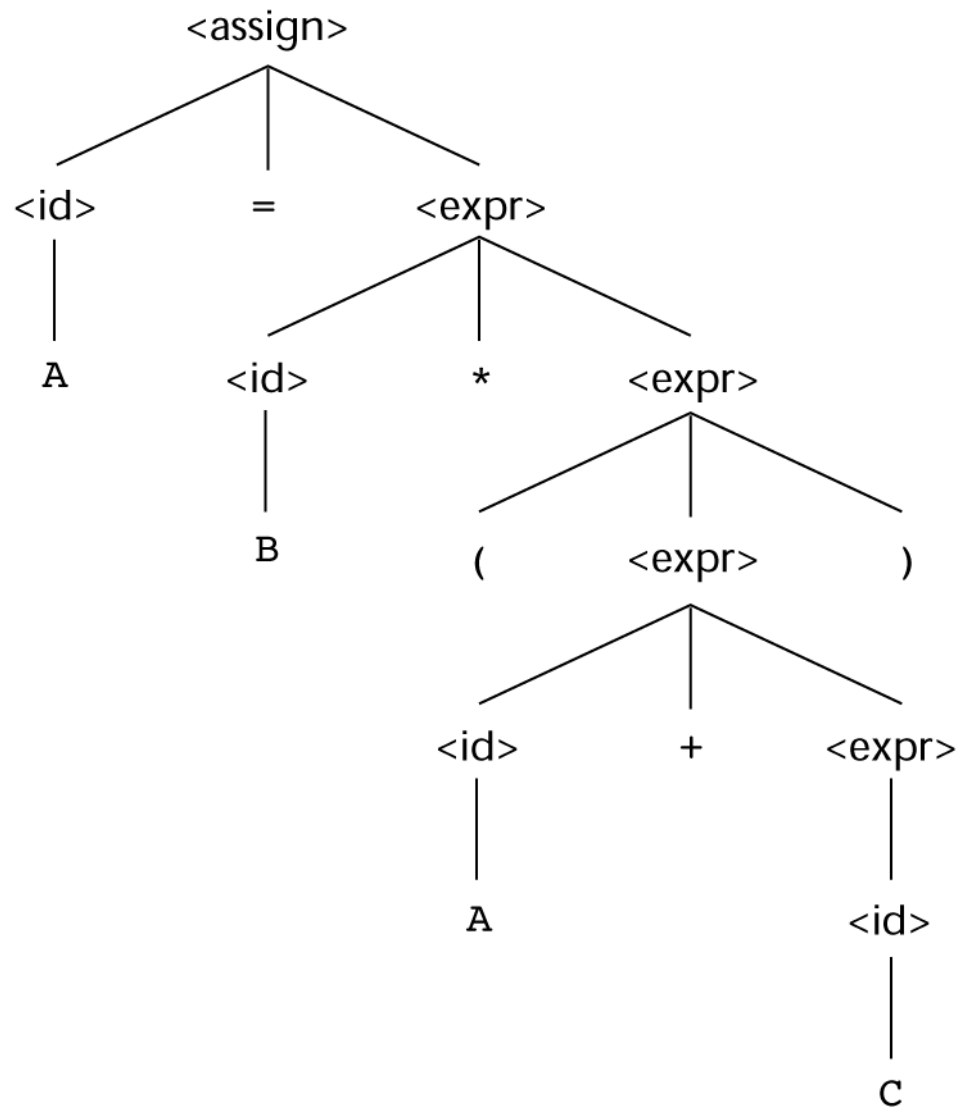
Parse Trees

- Grammars naturally describe the hierarchical syntactic structure of the sentences of the languages that they define
- These hierarchical structures are called parse trees
- Every internal node is a nonterminal, and every leaf is a terminal symbol
- A derivation can also be represented by a parse tree.
- In fact, a parse tree represents many derivations.

Parse trees



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

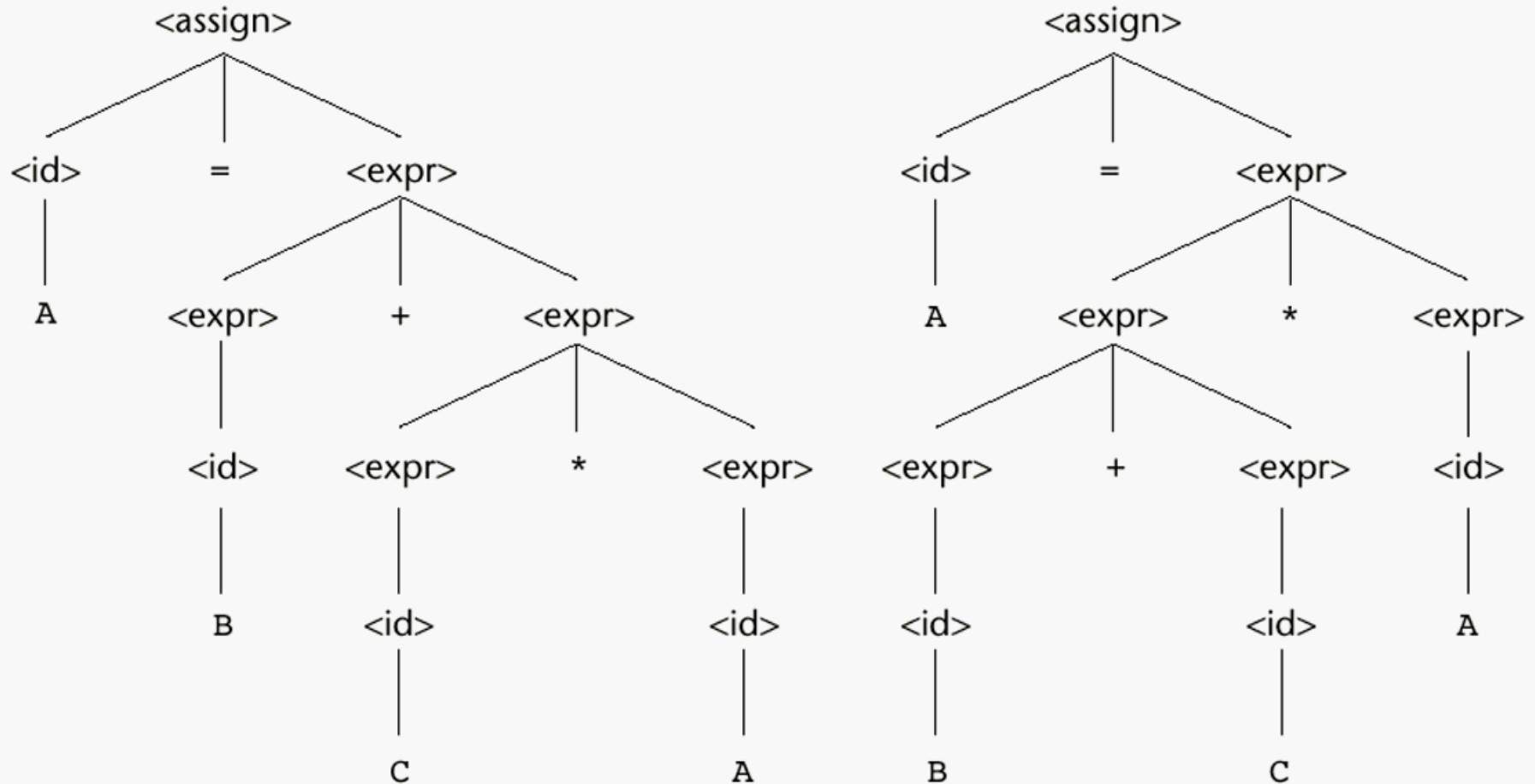


Ambiguous Grammar

A grammar that generates a sentential form for which there are two or more distinct parse trees is called as ambiguous

```
<assign> ::= <id> = <expr>
<id> ::= A | B | C
<expr> ::= <expr> + <expr>
          | <expr> * <expr>
          | (<expr>)
          | <id>
```


Two distinct parse trees for the same sentence, $A = B + C * A$



Ambiguity

The grammar of a PL must not be ambiguous

There are solutions for correcting the ambiguity

- Operator precedence
- Associativity rules

Operator precedence

In mathematics $*$ operation has a higher precedence than $+$
This can be implemented with extra nonterminals

$\langle \text{assign} \rangle ::= \langle \text{id} \rangle = \langle \text{expr} \rangle$

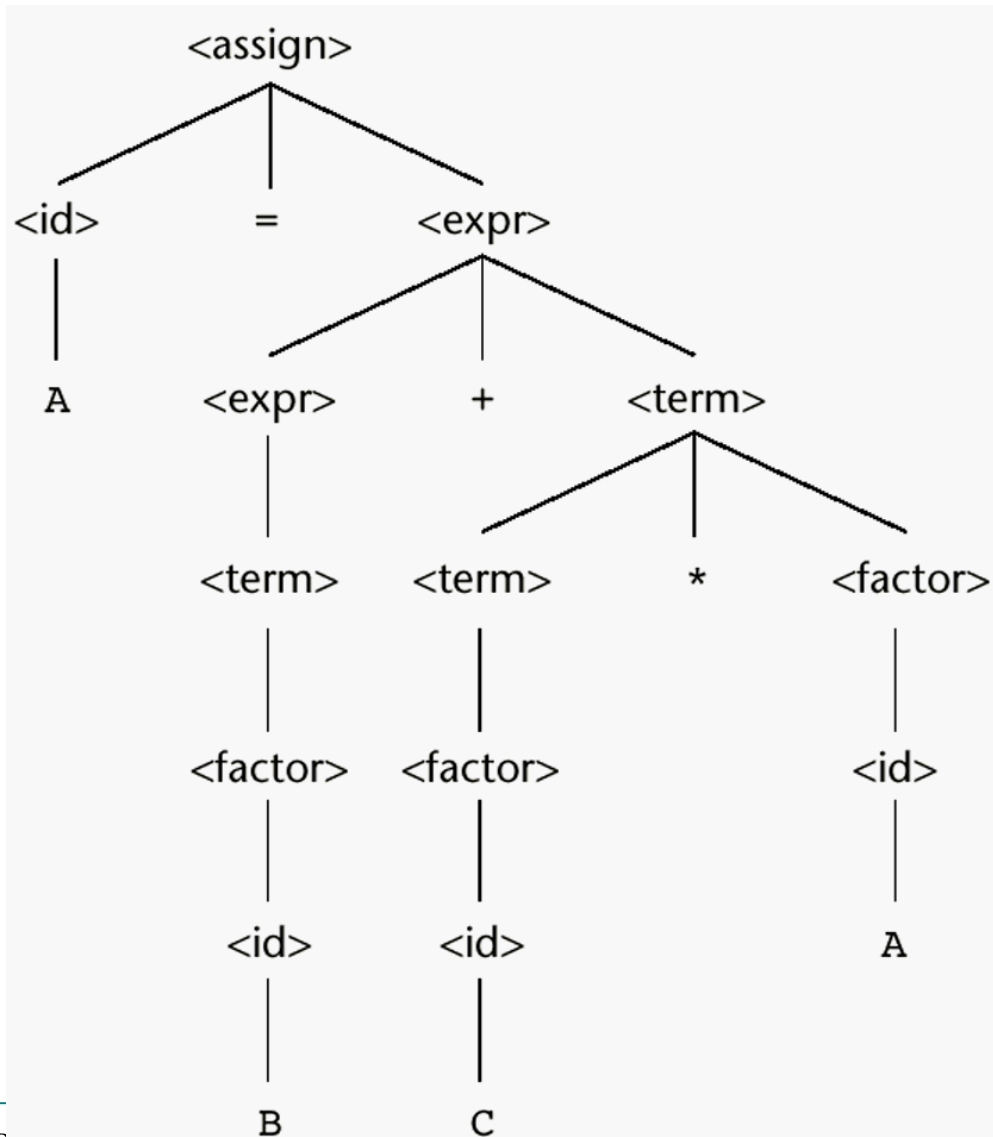
$\langle \text{id} \rangle ::= A \mid B \mid C$

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

$\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{factor} \rangle$
 $\mid \langle \text{factor} \rangle$

$\langle \text{factor} \rangle ::= (\langle \text{expr} \rangle)$
 $\mid \langle \text{id} \rangle$

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



Leftmost derivation using unambiguous grammar

$\langle \text{assign} \rangle \Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$
 $\Rightarrow A = \langle \text{expr} \rangle$
 $\Rightarrow A = \langle \text{expr} \rangle + \langle \text{term} \rangle$
 $\Rightarrow A = \langle \text{term} \rangle + \langle \text{term} \rangle$
 $\Rightarrow A = \langle \text{factor} \rangle + \langle \text{term} \rangle$
 $\Rightarrow A = \langle \text{id} \rangle + \langle \text{term} \rangle$
 $\Rightarrow A = B + \langle \text{term} \rangle$
 $\Rightarrow A = B + \langle \text{term} \rangle * \langle \text{factor} \rangle$
 $\Rightarrow A = B + \langle \text{factor} \rangle * \langle \text{factor} \rangle$
 $\Rightarrow A = B + \langle \text{id} \rangle * \langle \text{factor} \rangle$
 $\Rightarrow A = B + C * \langle \text{factor} \rangle$
 $\Rightarrow A = B + C * \langle \text{id} \rangle$
 $\Rightarrow A = B + C * A$

Rightmost derivation using unambiguous grammar

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

Associativity of Operators

What about equal precedence operators?

In math addition and multiplication are associative

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

However computer arithmetic may not be associative

e.g: for floating point addition where floating points values store 7 digits of accuracy, adding eleven numbers together where one of the numbers is 10^7 and the others are 1 result would be $1.000001 * 10^7$ only if the ten 1s are added first

Subtraction and division are not associative

$$A/B/C/D = ? \quad ((A/B)/C)/D \neq A/(B/(C/D))$$

Associativity

In a BNF rule, if the LHS appears at the beginning of the RHS, the rule is said to be left recursive

Left recursion specifies left associativity

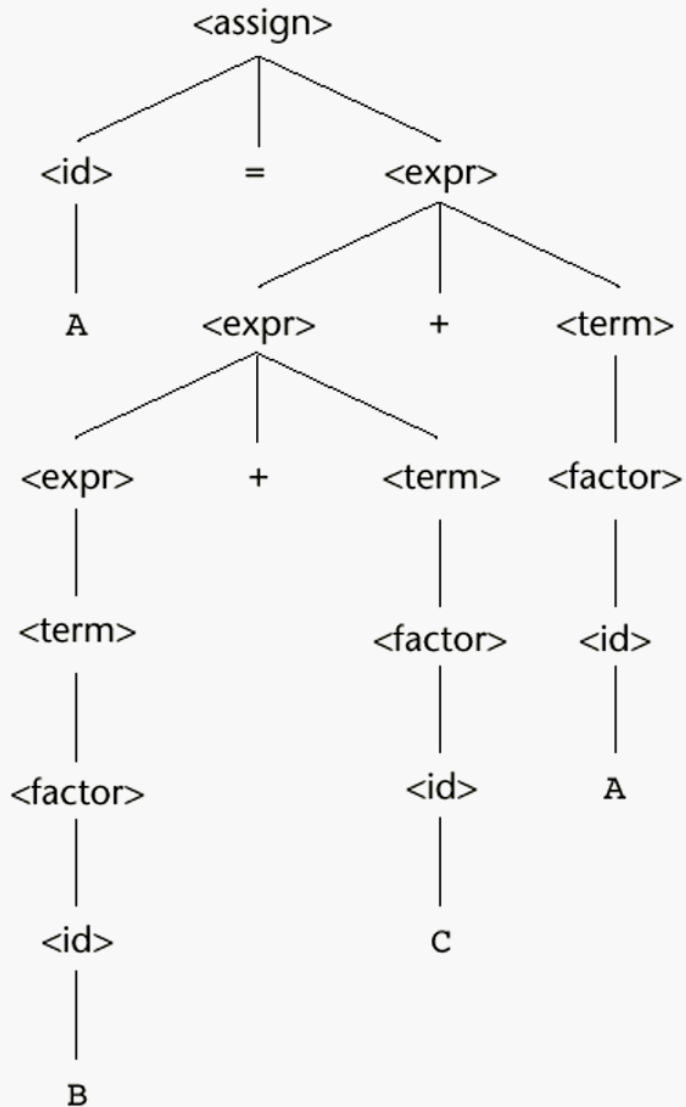
$$\begin{aligned} \langle \text{expr} \rangle &::= \langle \text{expr} \rangle + \langle \text{term} \rangle \\ &\quad | \langle \text{term} \rangle \end{aligned}$$

Similar for the right recursion

In most of the languages exponentiation is defined as a right associative operation

$$\begin{aligned} \langle \text{factor} \rangle &::= \langle \text{expr} \rangle ** \langle \text{factor} \rangle \\ &\quad | \langle \text{expr} \rangle \\ \langle \text{expr} \rangle &::= (\langle \text{expr} \rangle) \\ &\quad | \langle \text{id} \rangle \end{aligned}$$

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



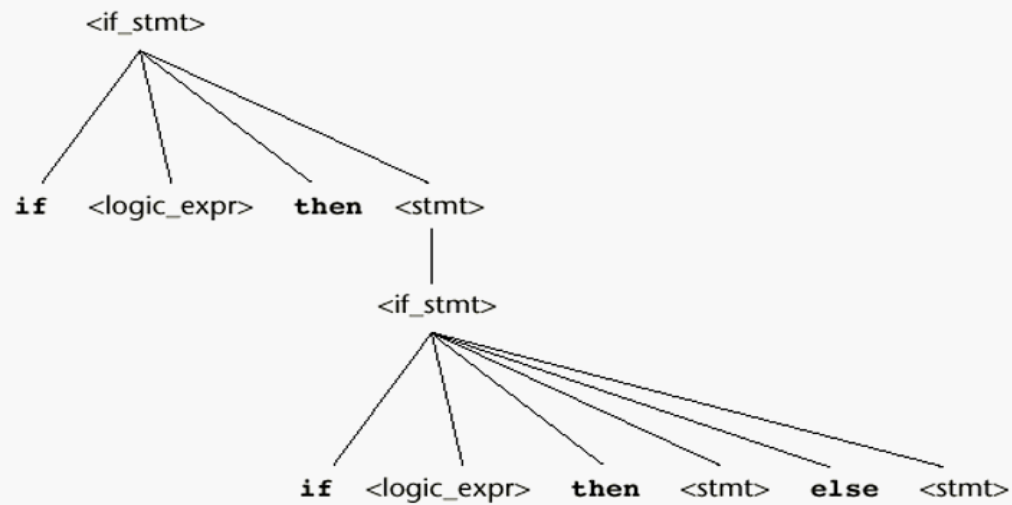
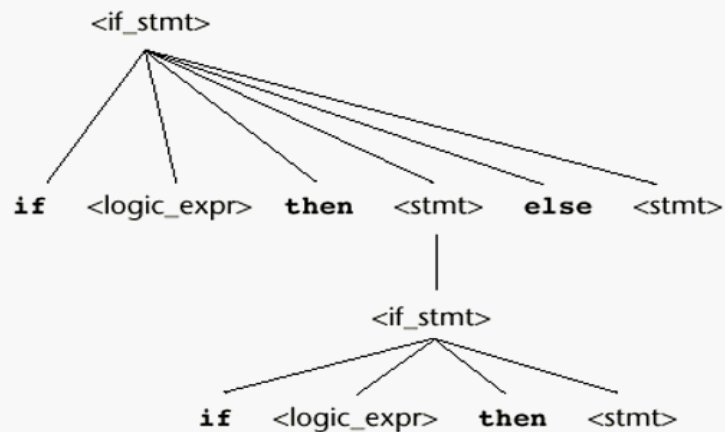
Left associativity

Left addition is lower than the right addition

Two distinct parse trees for the same sentential form

$\langle \text{if_stmt} \rangle ::= \text{if } \langle \text{logic_expr} \rangle \text{ then } \langle \text{stmt} \rangle$
 $\quad \mid \text{if } \langle \text{logic_expr} \rangle \text{ then } \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle$

If C1 then if C2 then A else B



An Unambiguous grammar for if then else

To design an unambiguous if-then-else statement we have to decide which if a dangling else belongs to

Dangling else problem: there are more if then else

Most PL adopt the following rule:

“an else is matched with the closest previous unmatched if statement”
(unmatched if = else-less if)

$\langle \text{stmt} \rangle ::= \langle \text{matched} \rangle \mid \langle \text{unmatched} \rangle$

$\langle \text{matched} \rangle ::= \text{if } \langle \text{logic_expr} \rangle \text{ then } \langle \text{matched} \rangle \text{ else } \langle \text{matched} \rangle$
 \mid any non-if-statement

$\langle \text{unmatched} \rangle ::= \text{if } \langle \text{logic_expr} \rangle \text{ then } \langle \text{stmt} \rangle$
 $\mid \text{if } \langle \text{logic_expr} \rangle \text{ then } \langle \text{matched} \rangle \text{ else } \langle \text{unmatched} \rangle$

there is a unique parse tree for this if statement

*BNF and Extended BNF

EBNF: same power but more convenient

[X] : X is optional (0 or 1 occurrence)

Equivalent to $X|\text{empty}$

$\langle \text{writeln} \rangle ::= \text{WRITELN } [(\langle \text{item_list} \rangle)]$

$\langle \text{selection} \rangle ::= \text{if } (\langle \text{expression} \rangle) \langle \text{statement} \rangle [\text{else} \langle \text{statement} \rangle]$

{X}: 0 or more occurrences

$A ::= \{X\}$ is equivalent to $A ::= XA|\text{empty}$

$\langle \text{identlist} \rangle = \langle \text{identifier} \rangle \{, \langle \text{identifier} \rangle\}$

$\{X1|X2|X3\}$: choose X1 or X2 or X3

$A ::= B(X|Y)C$ is equal to $A ::= AXC | AYC$

$\langle \text{for_stmt} \rangle ::= \text{for } \langle \text{var} \rangle := \langle \text{exp} \rangle (\text{to}|\text{downto}) \langle \text{exp} \rangle \text{ do } \langle \text{stmt} \rangle$

$\langle \text{term} \rangle ::= \langle \text{term} \rangle (*|/|\%) \langle \text{factor} \rangle$

BNF and Extended BNF

BNF:

$$\begin{aligned} \langle \text{expr} \rangle &::= \langle \text{expr} \rangle + \langle \text{term} \rangle \\ &\quad | \langle \text{expr} \rangle - \langle \text{term} \rangle \\ &\quad | \langle \text{term} \rangle \end{aligned}$$
$$\begin{aligned} \langle \text{term} \rangle &::= \langle \text{term} \rangle * \langle \text{factor} \rangle \\ &\quad | \langle \text{term} \rangle / \langle \text{factor} \rangle \\ &\quad | \langle \text{factor} \rangle \end{aligned}$$
$$\begin{aligned} \langle \text{factor} \rangle &::= \langle \text{expr} \rangle ** \langle \text{factor} \rangle \\ &\quad | \langle \text{expr} \rangle \end{aligned}$$
$$\begin{aligned} \langle \text{expr} \rangle &::= (\langle \text{expr} \rangle) \\ &\quad | \langle \text{id} \rangle \end{aligned}$$

BNF and Extended BNF

$\langle \text{expr} \rangle ::= \langle \text{term} \rangle \{ (+ \mid -) \langle \text{term} \rangle \}$
 $\langle \text{term} \rangle ::= \langle \text{factor} \rangle \{ (* \mid /) \langle \text{factor} \rangle \}$
 $\langle \text{factor} \rangle ::= \langle \text{expr} \rangle \{ ** \langle \text{expr} \rangle \}$
 $\langle \text{expr} \rangle ::= (\langle \text{expr} \rangle)$
 $\quad \mid \text{id}$

*Extended BNF

- $\langle \text{number} \rangle ::= \{ \langle \text{digit} \rangle \}$
- $\langle \text{signed number} \rangle ::= [+ \mid -] \langle \text{number} \rangle$

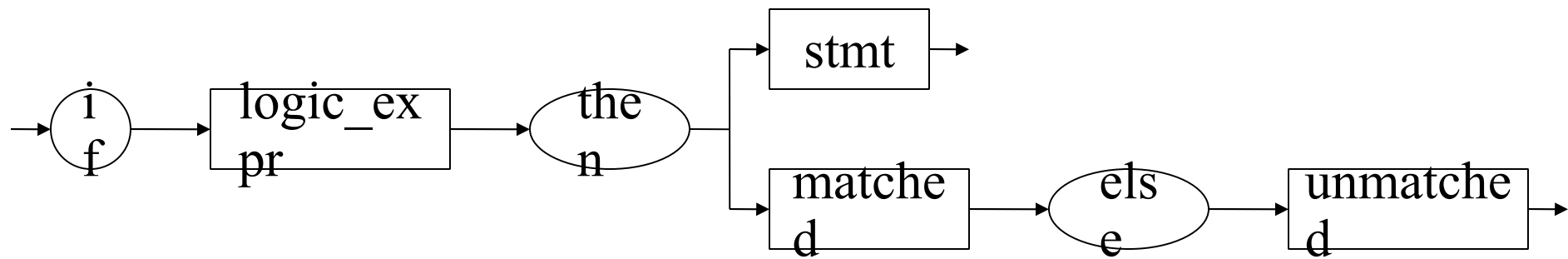
*Syntax Graphs

- Another form for representation of PL syntax
- Syntax graphs = syntax diagrams = syntax charts
- Equivalent to BNF in the power of representation, but easier to understand
- A separate graph is given for each syntactic unit (for each nonterminal)
- A rectangle represents a nonterminal, contains the name of the syntactic unit
- A circle (ellipse) represents a terminal

*Example

$\langle \text{unmatched} \rangle ::= \text{if } \langle \text{logic_expr} \rangle \text{ then } \langle \text{stmt} \rangle$
 $\quad \quad \quad | \text{if } \langle \text{logic_expr} \rangle \text{ then } \langle \text{matched} \rangle \text{ else } \langle \text{unmatched} \rangle$

unmatched



A syntax graph consists of one entry and one or more exit points
 If there exists a path from the input entry to any of the exit points
 corresponding to the string, then the string represents a valid instance
 of that unit. There may exist loops in the path