

# ***Language Description: Syntactic Structure***

**Expression Notations**

**Abstract Syntax Trees**

**Lexical Syntax**

**Context-Free Grammars**

**Grammars for Expressions**

**Variants of Grammars**

# Preface

- clear and complete descriptions
- formal syntax and informal semantics
- *syntax*: how programs in the language are built up
- *semantics*: what programs mean
- The same syntax has different semantics in different parts of the world.
  - Example: *DD/DD/DDDD*
    - 01/02/2001
    - January 2, 2001 or February 1, 2001 ?

# Preface

- abstract syntax: captures intent, independent of notation
  - Example:
    - $a+b$ ,  $(+ a b)$ , ADD a TO b
    - intent: apply the operator  $+$  to the operands  $a$  and  $b$
    - written representation is different in each case
    - abstract syntax is the same

# 2.1 Expression Notations

- Programming languages use a mix of *infix*, *prefix*, and *postfix* notations.
  - Example:
    - infix:  $a+b$
    - prefix:  $+ab$
    - postfix:  $ab+$
- An expression can be enclosed within parentheses without affecting its value.
  - expression  $E$  has the same value as  $(E)$
- *parenthesis-free*
  - *prefix*
  - *postfix*

# 2.1 Expression Notations

- *Prefix* Notation
  - a constant or a variable: itself
  - an operator **op** to subexpressions  $E1$  and  $E2$ :  
**op**  $E1$   $E2$
  - advantage: easy to decode during a left-to-right scan of an expression
  - Example:
    - sum of  $x$  and  $y$ :  $+ x y$
    - product of  $+ x y$  and  $z$ :  $* + x y z$
    - $* + 20 30 60$
    - $* 20 + 30 60$

# 2.1 Expression Notations

- *Prefix Notation* (cont)
  - *arity*: the number of operands of an operator
  - $\mathbf{op}^k$  of arity  $k \geq 0$  to  $E_1, E_2, \dots, E_k$ :  $\mathbf{op}^k E_1 E_2 \dots E_k$
  - during a left-to-right scan, the  $i$ th expression to the right of  $\mathbf{op}^k$  is the  $i$ th operand of  $\mathbf{op}^k$ , for  $1 \leq i \leq k$
  - *read(x)* and *max(x, y)*: a variant of prefix notation
    - allows operators to take a variable number of arguments
    - Example:
      - *write(root)*, *write(root, a, b, c)*
  - Lisp: (*read x*), (*max x y*)

# 2.1 Expression Notations

- *Postfix* Notation
  - a constant or a variable: itself
  - an operator **op** to subexpressions  $E1$  and  $E2$ :  
 $E1 E2 \mathbf{op}$
  - advantage: with the help of a *stack* data structure
  - Example:
    - sum of  $x$  and  $y$ :  $x y +$
    - product of  $x y +$  and  $z$ :  $x y + z *$
    - $20 30 + 60 *$
    - $20 30 60 + *$

# 2.1 Expression Notations

- *Infix* Notation: Precedence and Associativity
  - advantages: familiar, easy to read
  - decode  $a+b*c$ 
    - sum of  $a$  and  $b*c$ ?
    - product of  $a+b$  and  $c$ ?
    - solved by *precedence* and *associativity*
  - an operator at a higher *precedence level* takes its operands before an operator at a lower precedence level
  - precedence:  $* > +$ 
    - the operands of  $*$ :  $b$  and  $c$
    - the operands of  $+$ :  $a$  and  $b*c$



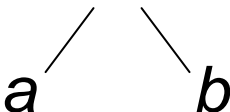
# 2.1 Expression Notations

- *Infix* Notation (cont)
  - precedence:
    - $*, / > +, -$
    - $*$  and  $/$ : at the same level
    - $+$  and  $-$ : at the same level
  - with the same precedence: *from left to right*
  - $4 - 2 - 1 = (4 - 2) - 1 = 2 - 1 = 1$
  - left associative
    - if subexpressions containing multiple occurrences of the operators are grouped *from left to right*
  - $+, -, *, /$ : left associative
  - $b * b - 4 * a * c = (b * b) - ((4 * a) * c)$

# 2.1 Expression Notations

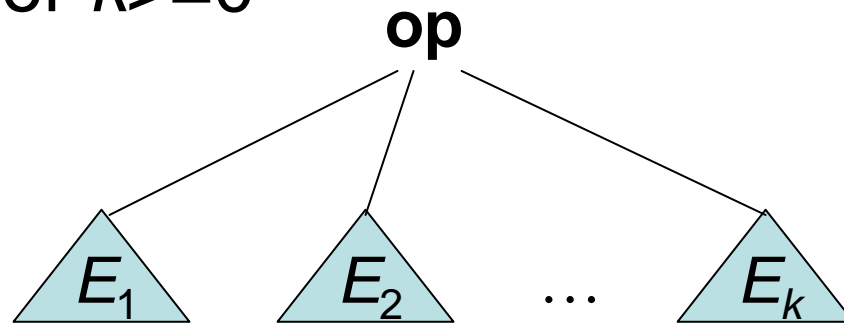
- *Infix* Notation (cont)
  - *right associative*: grouped from right to left
  - Example: exponentiation
    - $2^{3^4} = 2^{(3^4)} = 2^{81}$
- *Mixfix* Notation
  - operations specified by a combination of symbols do not fit neatly into the prefix, infix, postfix notations
  - Example: **if  $a > b$  then  $a$  else  $b$**
  - when symbols or keywords appear interspersed with the components of an expression, the operation will be said to be in *mixfix* notation

## 2.2 Abstract Syntax Trees

- *abstract syntax* of a language
  - identifies the meaningful components of each construct in the language
- prefix:  $+ a b$ , infix:  $a + b$ , postfix:  $a b +$ 
  - all have the same meaningful components
    - the operator:  $+$
    - the subexpressions:  $a$  and  $b$
  - tree representation: 

## 2.2 Abstract Syntax Trees

- Tree Representation of Expressions
  - applying an operator **op** to operands  $E_1, E_2, \dots, E_k$ , for  $k \geq 0$



– if an expression is a constant or a variable, then its tree consists of a leaf

# 2.2 Abstract Syntax Trees

- Tree Representation of Expressions (cont)

- Example:  $b*b-4*a*c$

- expression is of the form:  $E1-E2$

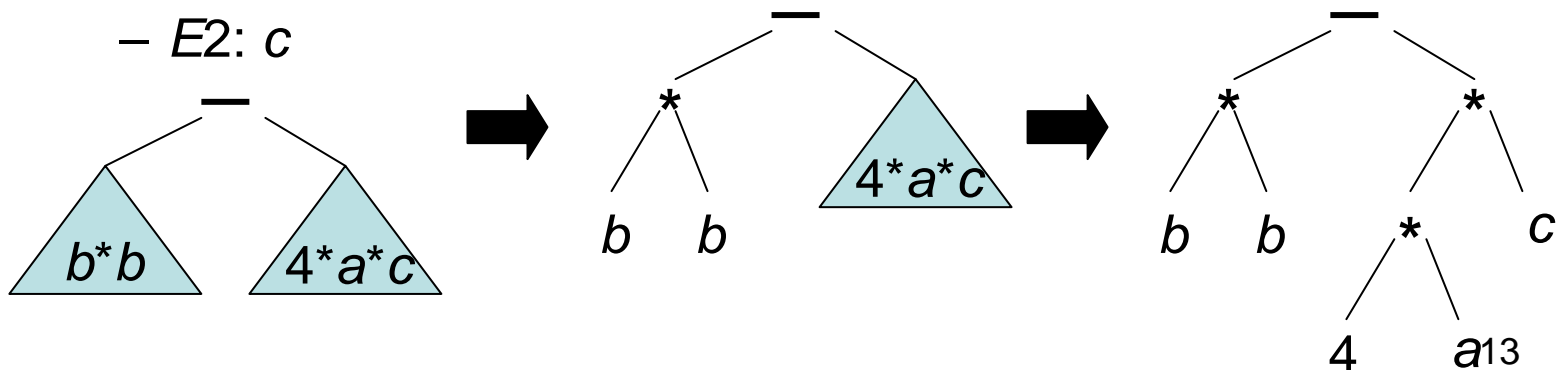
- $E1: b*b$

- $E2: 4*a*c$

- the subexpression has the form  $E1*E2$

- $E1: 4*a$

- $E2: c$

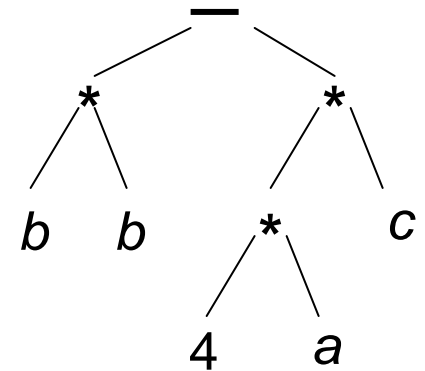


# 2.2 Abstract Syntax Trees

- Abstract Syntax Tree
  - show the operator/operand structure of an expression
  - show the syntactic structure of an expression independent of the notation

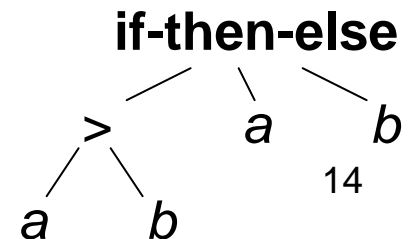
– Example:

- Prefix:  $- * b b * * 4 a c$
- Infix:  $b * b - 4 * a * c$
- Postfix:  $b b * 4 a * c * -$



– can be extended to other constructs by making up suitable operators for the constructs

- Example: **if  $a > b$  then  $a$  else  $b$**



## 2.3 Lexical Syntax

- Tokens and Spellings
  - the syntax of a programming language is specified in terms of units called *tokens* or *terminals*
  - *lexical syntax*
    - specifies the correspondence between the written *representation* of the language and the *tokens* or *terminals* in a grammar for the language
  - *keywords*
    - alphabetic character sequences that are treated as units in a language
    - Example: **if** and **while** are keywords in both Pascal and C
  - *Reserved Words*
    - keywords that cannot be used as names

## 2.3 Lexical Syntax

- Tokens and Spellings (cont)
  - *Spelling*
    - the actual character sequence used to write down an *occurrence* of a token is called the *spelling* of that occurrence
    - Example: the keyword **while** has spelling `while`
  - Using token **name** for names and token **number** for integers, the character sequence
$$b * b - 4 * a * c$$
is represented by the token sequence
$$\mathbf{name}_b * \mathbf{name}_b - \mathbf{number}_4 * \mathbf{name}_a * \mathbf{name}_a$$
  - *White space* in the form of *blank*, *tab*, and *newline* characters can typically be inserted between tokens without changing the meaning of a program



## 2.3 Lexical Syntax

- Tokens and Spellings (cont)
  - *comments* between tokens are ignored
    - (\* and \*) in Pascal
    - /\* and \*/ in C
  - the most complex rules in a lexical syntax are typically the ones describing the syntax of *real numbers*
  - Example:  
314.E-2, 3.14, 0.314E+1,  
0.314E1, .314E1

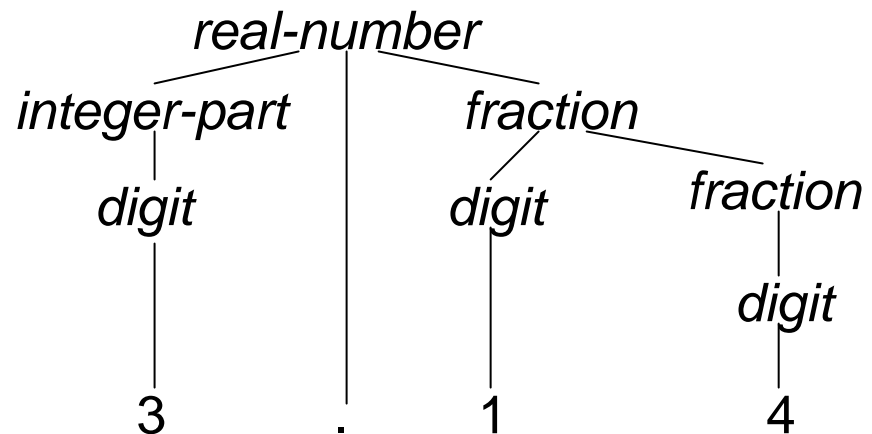
# 2.4 Context-Free Grammar

- *Context-Free Grammar*
  - a set of tokens or terminals
  - a set of nonterminals
  - a set of rules, productions
  - a starting nonterminal
- BNF: Backus-Naur Form
  - BNF rules for real numbers**

<i>&lt;real-number&gt;</i>	::=	<i>&lt;integer-part&gt;</i> . <i>&lt;fraction&gt;</i>
<i>&lt;integer-part&gt;</i>	::=	<i>&lt;digit&gt;</i>   <i>&lt;integer-part&gt;</i> <i>&lt;digit&gt;</i>
<i>&lt;fraction&gt;</i>	::=	<i>&lt;digit&gt;</i>   <i>&lt;digit&gt;</i> <i>&lt;fraction&gt;</i>
<i>&lt;digit&gt;</i>	::=	0   1   2   3   4   5   6   7   8   9

# 2.4 Context-Free Grammar

- *Parse Tree*
  - hierarchical structure
  - nonterminal
  - production
    - a rule that defines a nonterminal in terms of a sequence of terminals and nonterminal
    - each node in the parse tree is based on a production
  - Example: real number

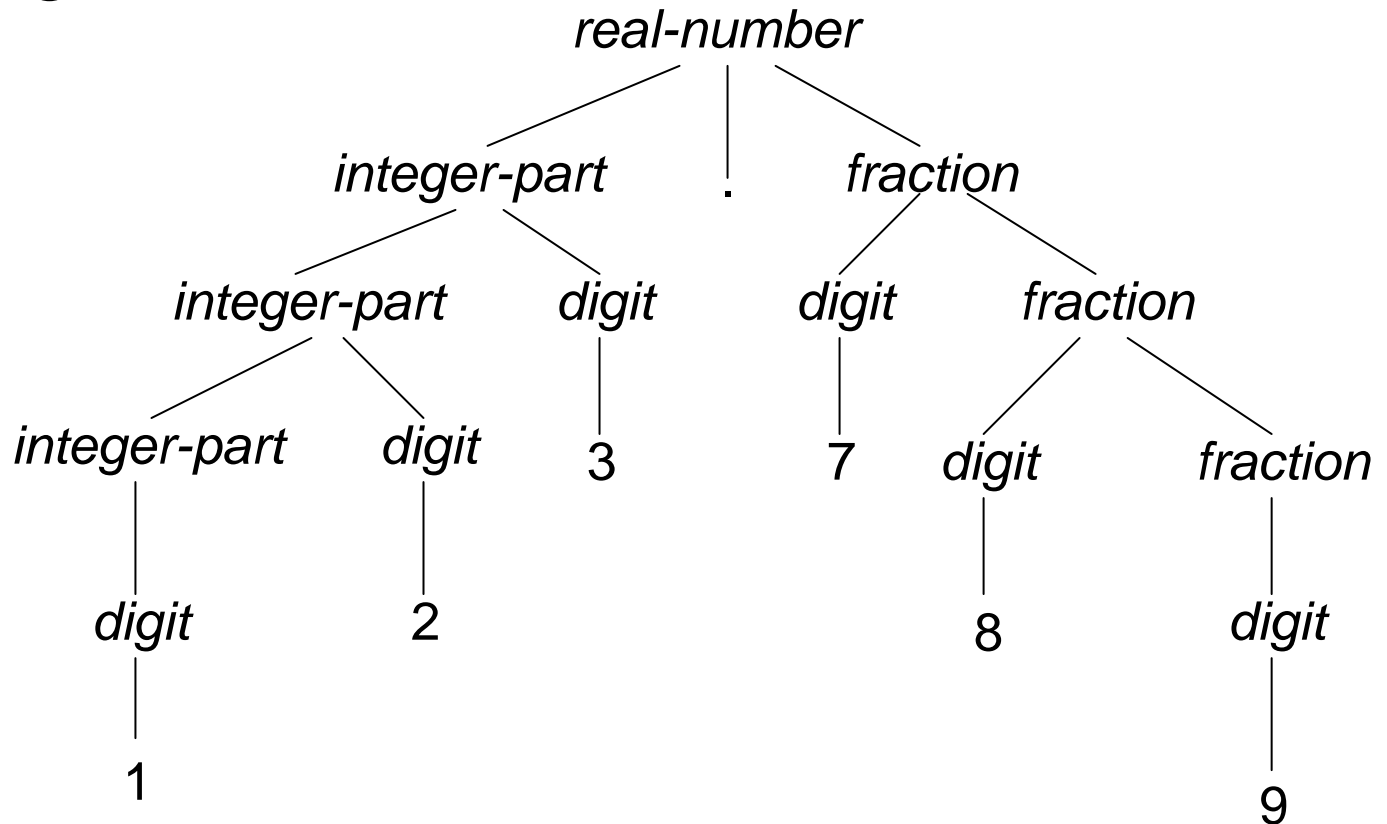


## 2.4 Context-Free Grammar

- *Parse Tree* (cont)
  - *<empty>*: an empty string of length 0
    - useful for specifying optional constructs
    - Example: 0 in 0.5
      - *<integer-part> ::= <empty> | <digit-sequence>*
  - each *leaf* is labeled with a terminal or *<empty>*
  - each *nonleaf* node is labeled with a nonterminal
  - the *root* is labeled with the starting nonterminal

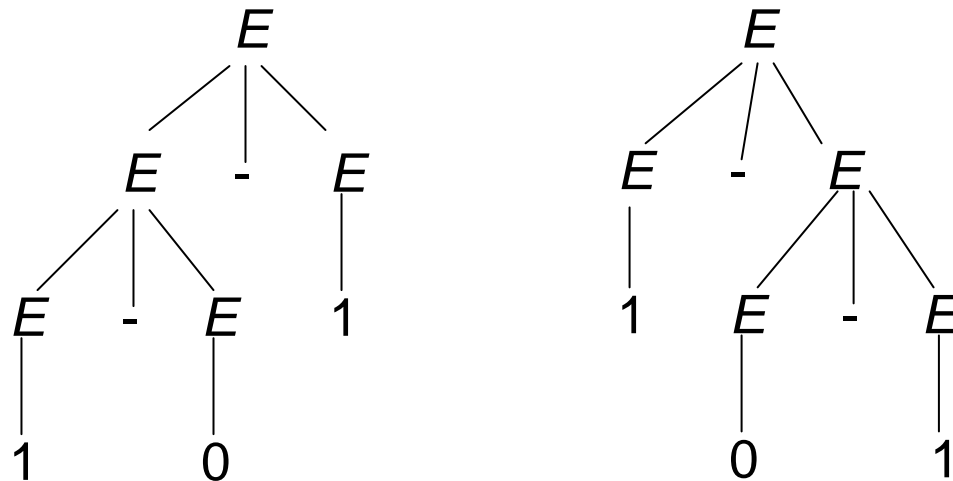
# 2.4 Context-Free Grammar

- parsing 123.789



# 2.4 Context-Free Grammar

- Syntactic Ambiguity
  - if some string in its language has *more than one* parse tree
  - Example.  $E ::= E - E \mid 0 \mid 1$



# 2.4 Context-Free Grammar

- *Dangling-Else Ambiguity*

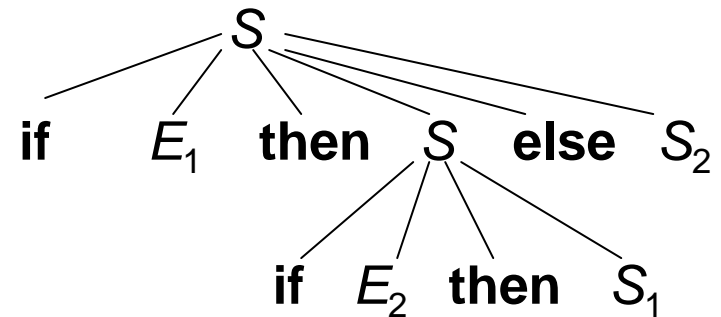
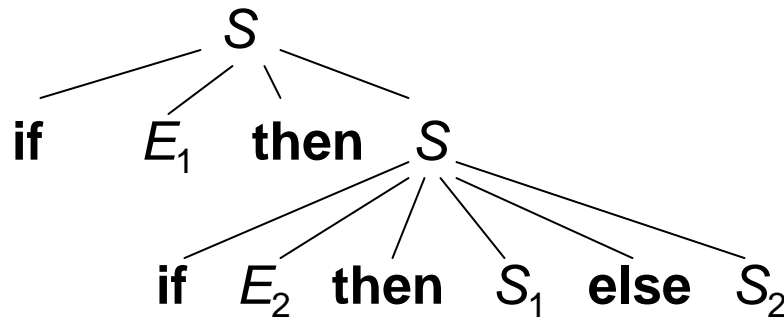
– production:

$S ::= \text{if } E \text{ then } S$

$S ::= \text{if } E \text{ then } S \text{ else } S$

– parsing:

$\text{if } E_1 \text{ then if } E_2 \text{ then } S_1 \text{ else } S_2$



resolved by matching an **else** with the nearest unmatched **if**

# 2.4 Context-Free Grammar

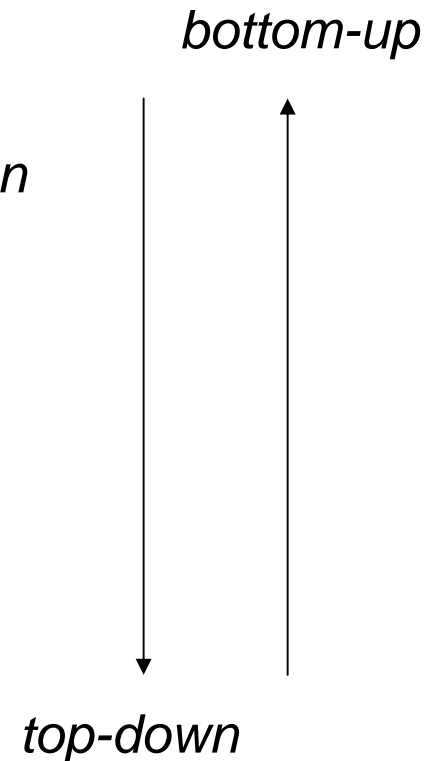
- *Derivations*
  - *top-down* parser
    - from the *root* of a parse tree toward the *leaves*
  - *bottom-up* parser
    - from the *leaves* of a parse tree toward the *root*
  - a derivation consists of a sequence of strings, beginning with the starting nonterminal
  - each successive string is obtained by *replacing a nonterminal by the right side of one of its production*
  - a derivation ends with a string consisting *entirely of terminals*



# 2.4 Context-Free Grammar

- *Derivations* (cont)
  - Example: parsing 21.89

*real-number ::= integer-part . fraction*  
*::= integer-part digit . fraction*  
*::= digit digit . fraction*  
*::= 2 digit . fraction*  
*::= 2 1 . fraction*  
*::= 2 1 . digit fraction*  
*::= 2 1 . 8 fraction*  
*::= 2 1 . 8 digit*  
*::= 2 1 . 8 9*



# 2.5 Grammars for Expressions

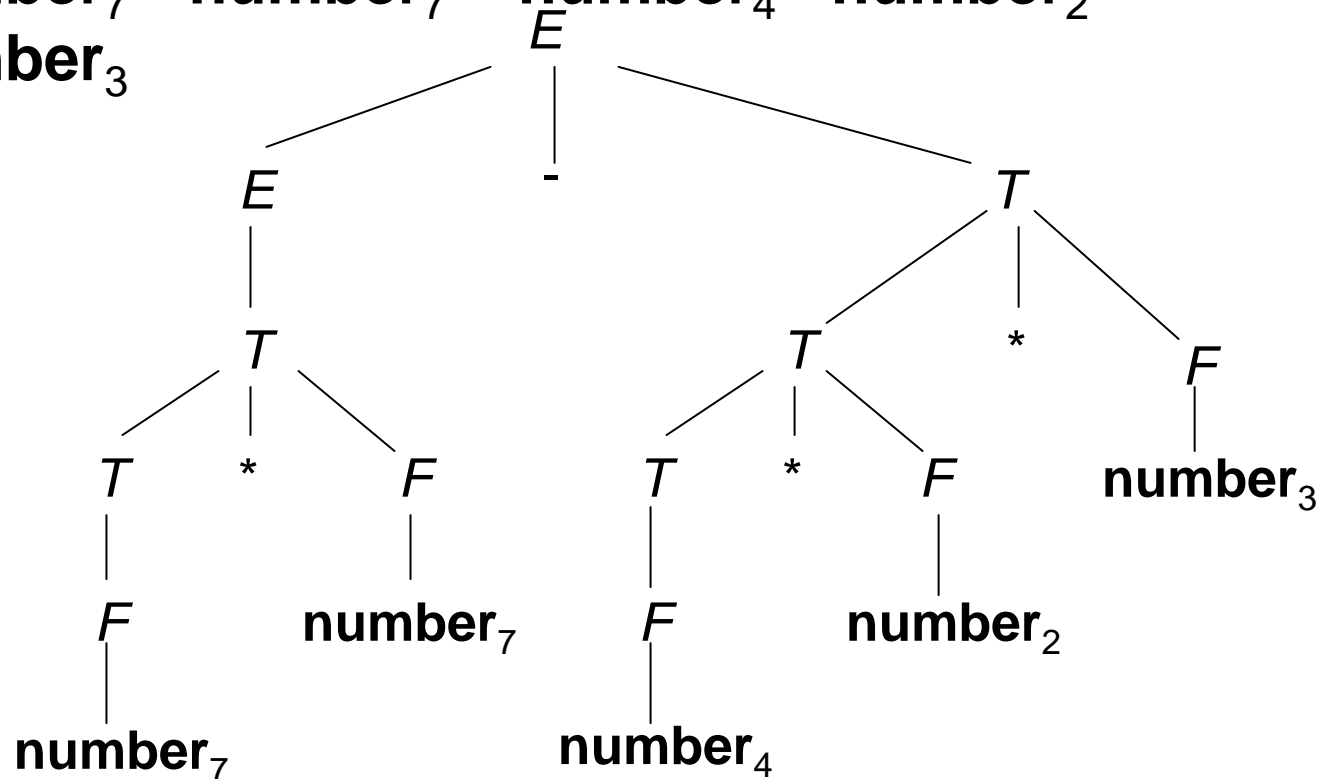
- Lists in Infix Expression
  - terms and factors
    - Example:  $a*b+c*d+e$ 
      - terms:  $a*b$ ,  $c*d$ ,  $e$
      - factors:  $a$ ,  $b$  for  $a*b$
  - a grammar for arithmetic expressions

$$E ::= E + T \mid E - T \mid T$$
$$T ::= T * F \mid T / F \mid F$$
$$F ::= \text{number} \mid \text{name} \mid ( E )$$

# 2.5 Grammars for Expressions

– parse tree:

- $\text{number}_7 * \text{number}_7 - \text{number}_4 * \text{number}_2 * \text{number}_3$



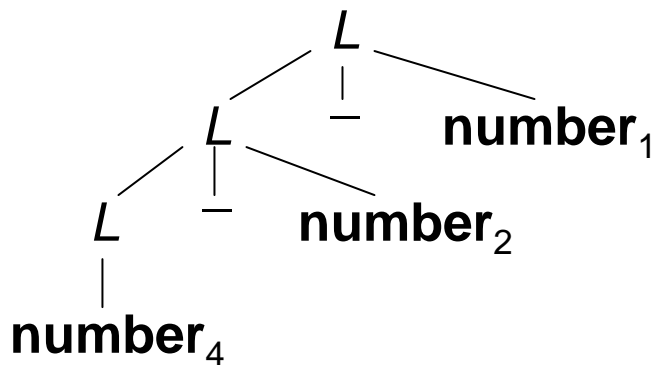
# 2.5 Grammars for Expressions

- Handling Associativity

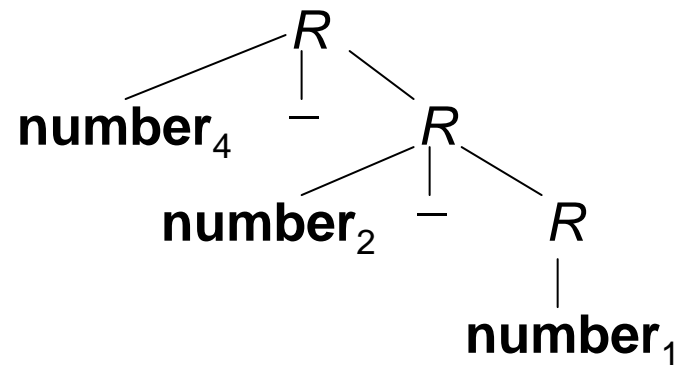
- Example:  $4 - 2 - 1$

$L ::= L + \text{number}$   
|  $L - \text{number}$   
|  $\text{number}$

$R ::= \text{number} + R$   
|  $\text{number} - R$   
|  $\text{number}$



left associative



right associative

# 2.5 Grammars for Expressions

- Handling Associativity and Precedence

$A ::= E := A \mid E$

$E ::= E + T \mid E - T \mid T$

$T ::= T * F \mid T / F \mid F$

$F ::= ( E ) \mid \text{name} \mid \text{number}$

## 2.6 Variants of Grammars

- BNF

$$\begin{aligned} \langle \textit{expression} \rangle &::= \langle \textit{expression} \rangle + \langle \textit{term} \rangle \\ &| \langle \textit{expression} \rangle - \langle \textit{term} \rangle \\ &| \langle \textit{term} \rangle \\ \langle \textit{term} \rangle &::= \langle \textit{term} \rangle * \langle \textit{factor} \rangle \\ &| \langle \textit{term} \rangle / \langle \textit{factor} \rangle \\ &| \langle \textit{factor} \rangle \\ \langle \textit{factor} \rangle &::= \mathbf{number} \\ &| \mathbf{name} \\ &| ( \langle \textit{expression} \rangle ) \end{aligned}$$

# 2.6 Variants of Grammars

- EBNF

- extension of BNF
- allows lists and optional elements to be specified
- for convenience, not for additional capability
- equivalent to BNF
- { and } : zero or more

- Example: { *<statements>* ; }

- *<statement list>* ::= { *<statement>* ; }

*<statement list>* ::= *<empty>*

≡ | *<statement>* ; *<statement list>*

- reduce the number of productions and nonterminal

# 2.6 Variants of Grammars

- EBNF (cont)
  - [ and ] : optional
    - Example:
      - $\langle real\ number \rangle ::= [ \langle integer\ part \rangle ] . \langle fraction \rangle$
      - $\langle \overline{real\ number} \rangle ::= \langle integer\ part \rangle . \langle fraction \rangle$   
| .  $\langle fraction \rangle$
  - | : choice
  - ( and ) : grouping
  - tokens & metasympols
    - enclosing tokens within single quotes: ‘(



## 2.6 Variants of Grammars

- EBNF:  $\langle expression \rangle ::= \langle term \rangle \{ ( + | - ) \langle term \rangle \}$   
 $\langle term \rangle ::= \langle factor \rangle \{ ( * | / ) \langle factor \rangle \}$   
 $\langle factor \rangle ::= '( \langle expression \rangle )' | \mathbf{name} | \mathbf{number}$
- BNF:  $\langle expression \rangle ::= \langle expression \rangle + \langle term \rangle$   
|  $\langle expression \rangle - \langle term \rangle$   
|  $\langle term \rangle$   
 $\langle term \rangle ::= \langle term \rangle * \langle factor \rangle$   
|  $\langle term \rangle / \langle factor \rangle$   
|  $\langle factor \rangle$   
 $\langle factor \rangle ::= \mathbf{number}$   
|  $\mathbf{name}$   
|  $( \langle expression \rangle )$

## 2.6 Variants of Grammars

- Syntax Chart (syntax graph, syntax diagram)
  - graphical notation
  - equivalent BNF

# 2.6 Variants of Grammars

- Syntax Chart

