

# Formal Languages and Grammars

---

Chapter 2: Sections 2.1 and 2.2

# Formal Languages

- Basis for the design and implementation of programming languages
- **Alphabet**: finite set  $\Sigma$  of symbols
- **String**: finite sequence of symbols
  - Empty string  $\varepsilon$ : sequence of length zero
  - $\Sigma^*$  - set of all strings over  $\Sigma$  (incl.  $\varepsilon$ )
  - $\Sigma^+$  - set of all non-empty strings over  $\Sigma$
- **Language**: set of strings  $L \subseteq \Sigma^*$ 
  - E.g., for Java,  $\Sigma$  is Unicode, a string is a program, and  $L$  is defined by a grammar in the language spec

# Formal Grammars

- $G = (N, T, S, P)$ 
  - Finite set of **non-terminal symbols**  $N$
  - Finite set of **terminal symbols**  $T$
  - Starting non-terminal symbol  $S \in N$
  - Finite set of **productions**  $P$
  - Describes a language  $L \subseteq T^*$
- Production:  $\mathbf{x} \rightarrow \mathbf{y}$ 
  - $\mathbf{x}$  is a non-empty sequence of terminals and non-terminals;  $\mathbf{y}$  is a seq. of terminals and non-terminals
- Applying a production:  $\mathbf{uxv} \Rightarrow \mathbf{uyw}$

# Example: Non-negative Integers

- $N = \{ I, D \}$
- $T = \{ 0, 1, 2, 3, 4, 5, 6, 7, 8, 9 \}$
- $S = I$
- $P = \{ \begin{array}{l} I \rightarrow D, \\ I \rightarrow DI, \\ D \rightarrow 0, \\ D \rightarrow 1, \\ \dots, \\ D \rightarrow 9 \end{array} \}$

# More Common Notation

$I \rightarrow D \mid DI$  - two production alternatives

$D \rightarrow 0 \mid 1 \mid \dots \mid 9$  - ten production alternatives

- Terminals: 0 ... 9
- Starting non-terminal: I
  - Shown first in the list of productions
- Examples of production applications:

$\underline{I} \Rightarrow \underline{DI}$

$D6\underline{I} \Rightarrow D6\underline{D}$

$D\underline{I} \Rightarrow D\underline{DI}$

$\underline{D}6D \Rightarrow \underline{3}6D$

$D\underline{DI} \Rightarrow D\underline{6I}$

$36\underline{D} \Rightarrow 36\underline{1}$

# Languages and Grammars

- String derivation
  - $w_1 \Rightarrow w_2 \Rightarrow \dots \Rightarrow w_n$ ; denoted  $w_1 \xRightarrow{*} w_n$
  - If  $n > 1$ , non-empty derivation sequence:  $w_1 \xRightarrow{+} w_n$
- Language generated by a grammar
  - $L(G) = \{ w \in T^* \mid S \xRightarrow{+} w \}$
- Fundamental theoretical characterization:  
Chomsky hierarchy (Noam Chomsky, MIT)
  - Regular languages  $\subset$  Context-free languages  $\subset$   
Context-sensitive languages  $\subset$  Unrestricted languages
  - Regular languages in PL: for **lexical analysis**
  - Context-free languages in PL: for **syntax analysis**

# Regular Languages (1/5)

- Operations on languages
  - **Union**:  $L \cup M$  = all strings in  $L$  or in  $M$
  - **Concatenation**:  $LM$  = all  $ab$  where  $a$  in  $L$  and  $b$  in  $M$
  - $L^0 = \{ \varepsilon \}$  and  $L^i = L^{i-1}L$
  - **Closure**:  $L^* = L^0 \cup L^1 \cup L^2 \cup \dots$
  - **Positive closure**:  $L^+ = L^1 \cup L^2 \cup \dots$
- Regular expressions: notation to express languages constructed with the help of such operations
  - Example:  $(0|1|2|3|4|5|6|7|8|9)^+$

## Regular Languages (2/5)

- Given some alphabet, a **regular expression** is
  - The empty string  $\epsilon$
  - Any symbol from the alphabet
  - If  $r$  and  $s$  are regular expressions, so are  $r|s$ ,  $rs$ ,  $r^*$ ,  $r^+$ ,  $r?$ , and  $(r)$
  - $^*/^+/?$  have higher precedence than concatenation, which has higher precedence than  $|$
  - All are left-associative



# Regular Languages (3/5)

- Each regular expression  $r$  defines a language  $L(r)$ 
  - $L(\varepsilon) = \{ \varepsilon \}$
  - $L(a) = \{ a \}$  for alphabet symbol  $a$
  - $L(r|s) = L(r) \cup L(s)$
  - $L(rs) = L(r)L(s)$
  - $L(r^*) = (L(r))^*$
  - $L(r^+) = (L(r))^+$
  - $L(r?) = \{ \varepsilon \} \cup L(r)$
  - $L((r)) = L(r)$
- Example: what is the language defined by  
 $0(x|X)(0|1|\dots|9|a|b|\dots|f|A|B|\dots|F)^+$

# Regular Languages (4/5)

- **Regular grammars**

- All productions are  $A \rightarrow wB$  and  $A \rightarrow w$

- $A$  and  $B$  are non-terminals;  $w$  is a sequence of terminals
- This is a right-regular grammar

- Or all productions are  $A \rightarrow Bw$  and  $A \rightarrow w$

- Left-regular grammar

- Example:  $L = \{ a^n b \mid n > 0 \}$  is a regular language

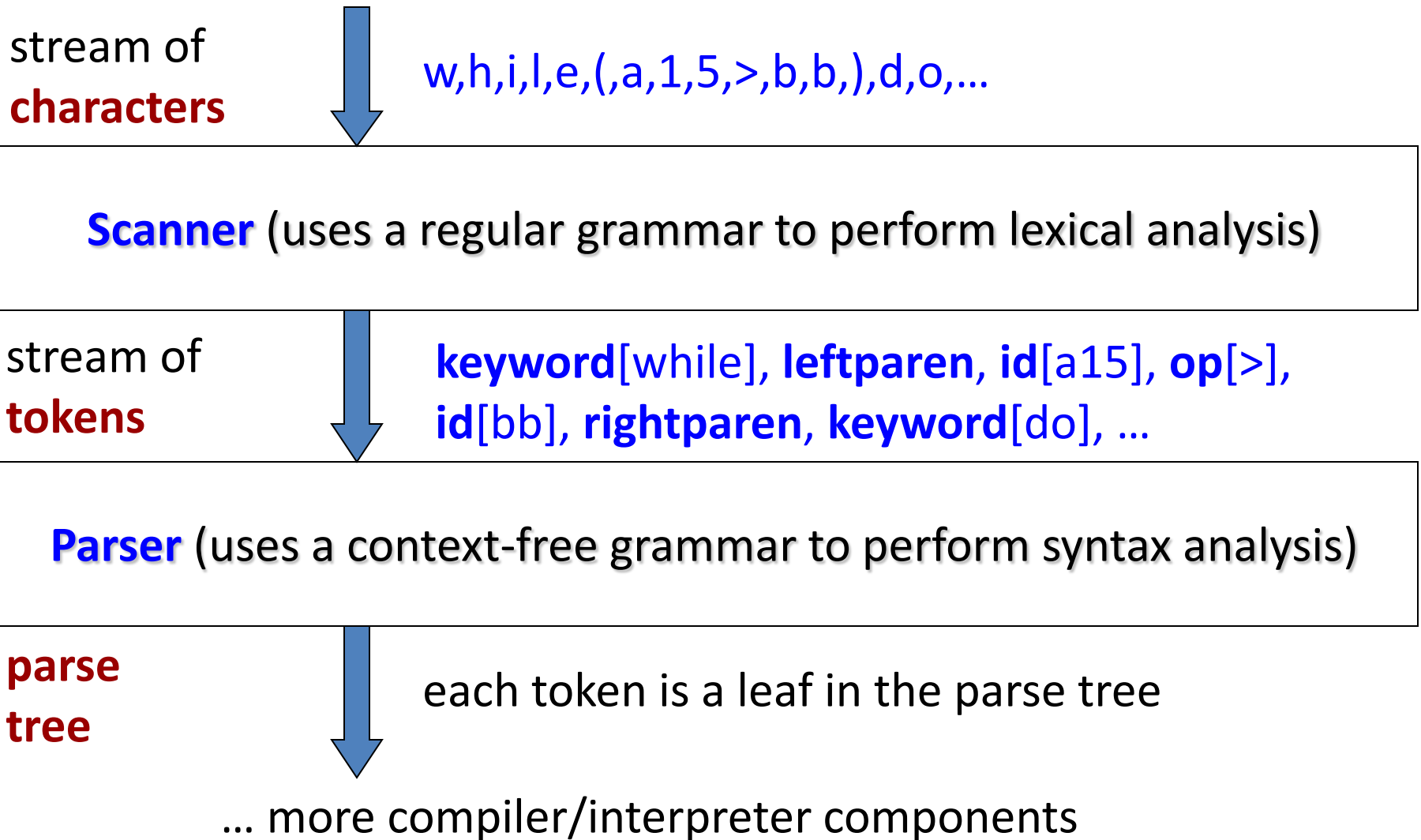
- $S \rightarrow Ab$  and  $A \rightarrow a \mid Aa$

- $I \rightarrow D \mid DI$  and  $D \rightarrow 0 \mid 1 \mid \dots \mid 9$  : is this a regular grammar?

# Regular Languages (5/5)

- Equivalent formalisms for regular languages
  - Regular grammars
  - Regular expressions
  - Nondeterministic finite automata (NFA)
  - Deterministic finite automata (DFA)
  - Additional details: Sections 2.2 and 2.4
- What does this have to do with PLs?
  - Foundation for **lexical analysis** done by a **scanner**
  - You will have to implement a scanner for your interpreter project; Section 2.2 provides useful guidelines

# Regular Languages in Compilers & Interpreters



# Uses of Regular Languages

- Lexical analysis in compilers
  - E.g., an identifier token is a string from the regular language **letter (letter|digit)\***
  - Each token is a **terminal symbol** for the context-free grammar of the parser
- Pattern matching
  - **stdlinux> grep "a\+b" foo.txt**
  - Find every line from foo.txt that contains a string from the language  $L = \{ a^n b \mid n > 0 \}$ 
    - i.e., the language for reg. expr.  $a^+b$

# Context-Free Languages

- They subsume regular languages
  - Every regular language is a c.f. language
  - $L = \{ a^n b^n \mid n > 0 \}$  is c.f. but not regular
- Generated by a **context-free grammar**
  - Each production:  $A \rightarrow w$
  - $A$  is a non-terminal,  $w$  is a sequences of terminals and non-terminals
- BNF (Backus-Naur Form): traditional alternative notation for context-free grammars
  - John Backus and Peter Naur, for Algol-58 and Algol-60
    - Backus was also one of the creators of Fortran
  - Both are recipients of the ACM Turing Award

# Example: Non-negative Integers

- $I \rightarrow D \mid DI$  and  $D \rightarrow 0 \mid 1 \mid \dots \mid 9$
- BNF
  - $\langle \text{integer} \rangle ::= \langle \text{digit} \rangle \mid \langle \text{digit} \rangle \langle \text{integer} \rangle$
  - $\langle \text{digit} \rangle ::= 0 \mid 1 \mid \dots \mid 9$
- What if we wanted to disallow zeroes at the beginning?
  - e.g. 509 is OK, but 059 is not
    - Possible motivation: in C, leading 0 means an octal constant
  - Propose a context-free grammar that achieves this
    - Is this grammar regular? If not, can you change it to make it regular?

# Derivation Tree for a String

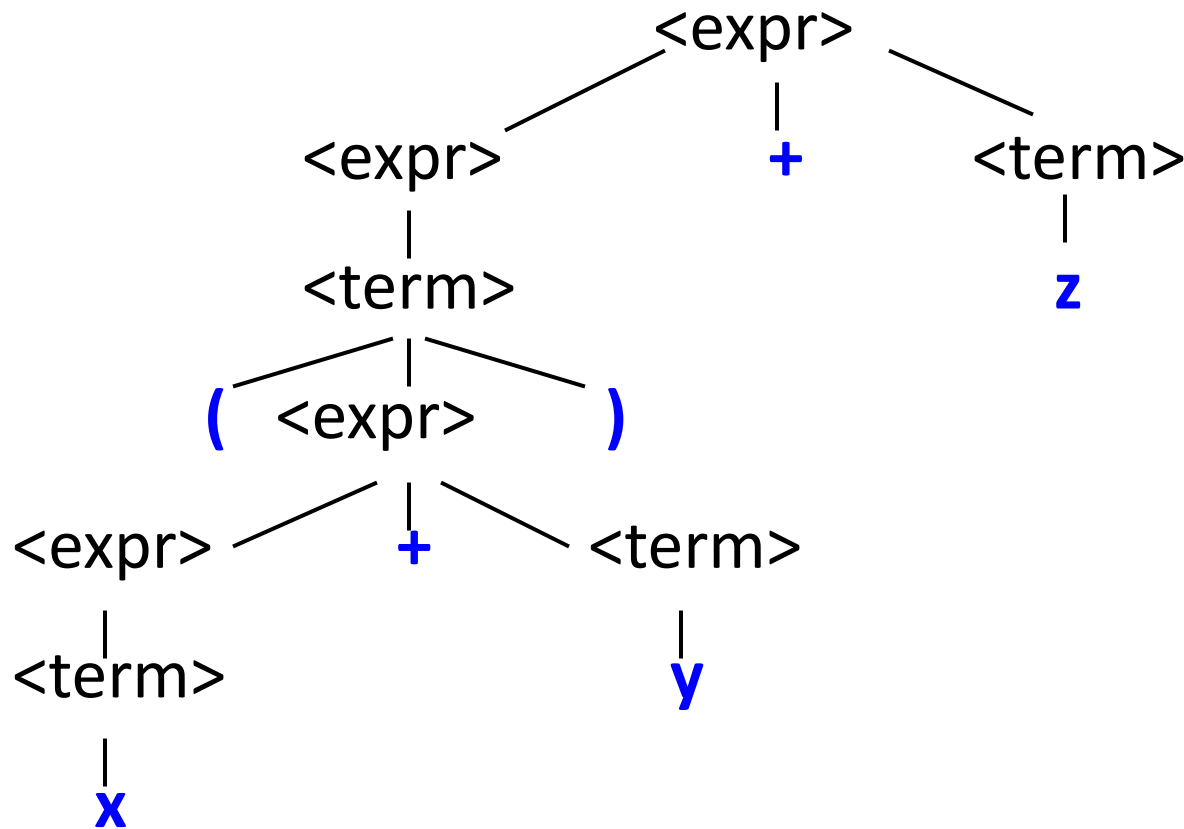
- Also called **parse tree** or **concrete syntax tree**
  - Leaf nodes: terminals
  - Inner nodes: non-terminals
  - Root: starting non-terminal of the grammar
- Describes a particular way to derive a string based on a context-free grammar
  - Leaf nodes from left to right are the string
  - To get this string: depth-first traversal of the tree, always visiting the leftmost unexplored branch



# Example of a Derivation Tree

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

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



Parse tree for  
 $(x+y)+z$

# Equivalent Derivation Sequences

The set of string derivations that are represented by the same parse tree

One derivation:

$$\begin{aligned} \langle \text{expr} \rangle &\Rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle \Rightarrow \langle \text{expr} \rangle + z \Rightarrow \\ \langle \text{term} \rangle + z &\Rightarrow (\langle \text{expr} \rangle) + z \Rightarrow \\ (\langle \text{expr} \rangle + \langle \text{term} \rangle) + z &\Rightarrow (\langle \text{expr} \rangle + y) + z \Rightarrow \\ (\langle \text{term} \rangle + y) + z &\Rightarrow (x + y) + z \end{aligned}$$

Another derivation:

$$\begin{aligned} \langle \text{expr} \rangle &\Rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle \Rightarrow \langle \text{term} \rangle + \langle \text{term} \rangle \Rightarrow \\ (\langle \text{expr} \rangle) + \langle \text{term} \rangle &\Rightarrow (\langle \text{expr} \rangle + \langle \text{term} \rangle) + \langle \text{term} \rangle \Rightarrow \\ (\langle \text{term} \rangle + \langle \text{term} \rangle) + \langle \text{term} \rangle &\Rightarrow (x + \langle \text{term} \rangle) + \langle \text{term} \rangle \Rightarrow \\ (x + y) + \langle \text{term} \rangle &\Rightarrow (x + y) + z \end{aligned}$$

Many more ...

# Ambiguous Grammars

- For some string, there are several different parse trees
- An ambiguous grammar gives more freedom to the compiler writer
  - e.g. for code optimizations, to choose the shape of the parse tree that leads to better performance
- For real-world programming languages, we typically have non-ambiguous grammars
  - We need a deterministic specification of the parser
  - To remove the ambiguity: add non-terminals

# Elimination of Ambiguity (1/2)

- $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid \langle \text{expr} \rangle * \langle \text{expr} \rangle$   
 $\mid \text{id} \mid ( \langle \text{expr} \rangle )$
- Two possible parse trees for  $a + b * c$ 
  - Conceptually equivalent to  $(a + b) * c$  and  $a + (b * c)$
- Operator **precedence**: when several operators are without parentheses, what is an operand of what?
  - Is  $a+b$  an operand of  $*$ , or is  $b*c$  an operand of  $+$ ?
- Operator **associativity**: for several operators with the same precedence, left-to-right or right-to-left?
  - Is  $a - b - c$  equivalent to  $(a - b) - c$  or  $a - (b - c)$ ?

## Elimination of Ambiguity (2/2)

- In most languages,  $*$  has higher precedence than  $+$ , and both are left-associative
- Problem: change  $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid \langle \text{expr} \rangle * \langle \text{expr} \rangle \mid \text{id} \mid ( \langle \text{expr} \rangle )$ 
  - Eliminate the ambiguity
  - Get the correct precedence and associativity
- Solution: add new non-terminals
  - $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{term} \rangle \mid \langle \text{term} \rangle$
  - $\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{factor} \rangle \mid \langle \text{factor} \rangle$
  - $\langle \text{factor} \rangle ::= \text{id} \mid ( \langle \text{expr} \rangle )$

# The “dangling-else” Problem

- Ambiguity for “else”

$\langle \text{stmt} \rangle ::= \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle$

|  $\text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle$

- **if** a **then if** b **then** c:=1 **else** c:=2

– Two possible parse trees

- Traditional solution: match the **else** with the last unmatched **then**

# Non-Ambiguous Grammar

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

$\langle \text{matched} \rangle ::= \langle \text{non-if-stmt} \rangle \mid$

$\text{if } \langle \text{expr} \rangle \text{ then } \langle \text{matched} \rangle \text{ else } \langle \text{matched} \rangle$

$\langle \text{unmatched} \rangle ::= \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \mid$

$\text{if } \langle \text{expr} \rangle \text{ then } \langle \text{matched} \rangle \text{ else } \langle \text{unmatched} \rangle$

# Extended BNF (EBNF)

- [ ... ] optional element
  - **if** <expr> **then** <stmt> [ **else** <stmt> ]
- { ... } repetition (0 or more times)
  - <IdList> ::= <id> { , <id> }
  - Sometimes shown as { ... }\*
- { ... }<sup>+</sup> repetition (1 or more times)
  - <block> ::= **begin** <stmt> { <stmt> } **end**
  - <block> ::= **begin** { <stmt> }<sup>+</sup> **end**
- Does not change the expressive power of the notation (we can always rewrite in plain BNF)



# Core: A Toy Imperative Language (1/2)

$\langle \text{prog} \rangle ::= \text{program } \langle \text{decl-seq} \rangle \text{ begin } \langle \text{stmt-seq} \rangle \text{ end}$

$\langle \text{decl-seq} \rangle ::= \langle \text{decl} \rangle \mid \langle \text{decl} \rangle \langle \text{decl-seq} \rangle$

$\langle \text{stmt-seq} \rangle ::= \langle \text{stmt} \rangle \mid \langle \text{stmt} \rangle \langle \text{stmt-seq} \rangle$

$\langle \text{decl} \rangle ::= \text{int } \langle \text{id-list} \rangle ; \quad \langle \text{id-list} \rangle ::= \text{id} \mid \text{id} , \langle \text{id-list} \rangle$

$\langle \text{stmt} \rangle ::= \langle \text{assign} \rangle \mid \langle \text{if} \rangle \mid \langle \text{loop} \rangle \mid \langle \text{in} \rangle \mid \langle \text{out} \rangle$

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

$\langle \text{in} \rangle ::= \text{input } \langle \text{id-list} \rangle ; \quad \langle \text{out} \rangle ::= \text{output } \langle \text{id-list} \rangle ;$

$\langle \text{if} \rangle ::= \text{if } \langle \text{cond} \rangle \text{ then } \langle \text{stmt-seq} \rangle \text{ endif ;}$

$\mid \text{if } \langle \text{cond} \rangle \text{ then } \langle \text{stmt-seq} \rangle \text{ else } \langle \text{stmt-seq} \rangle \text{ endif ;}$

## Core: A Toy Imperative Language (2/2)

$\langle \text{loop} \rangle ::= \text{while } \langle \text{cond} \rangle \text{ begin } \langle \text{stmt-seq} \rangle \text{ endwhile ;}$

$\langle \text{cond} \rangle ::= \langle \text{cmpr} \rangle \mid ! \langle \text{cond} \rangle \mid ( \langle \text{cond} \rangle \text{ AND } \langle \text{cond} \rangle )$

$\mid ( \langle \text{cond} \rangle \text{ OR } \langle \text{cond} \rangle )$

$\langle \text{cmpr} \rangle ::= [ \langle \text{expr} \rangle \langle \text{cmpr-op} \rangle \langle \text{expr} \rangle ]$

$\langle \text{cmpr-op} \rangle ::= < \mid = \mid != \mid > \mid >= \mid <=$

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

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

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

# Parser vs. Scanner

- **id** and **const** are terminal symbols for the grammar of the language
  - **tokens** that are provided from the scanner to the parser
- But they are non-terminals in the regular grammar for the lexical analysis

– The terminals here are ASCII characters

$\langle \text{id} \rangle ::= \langle \text{letter} \rangle \mid \langle \text{id} \rangle \langle \text{letter} \rangle \mid \langle \text{id} \rangle \langle \text{digit} \rangle$

$\langle \text{letter} \rangle ::= \mathbf{A} \mid \mathbf{B} \mid \dots \mid \mathbf{Z} \mid \mathbf{a} \mid \mathbf{b} \mid \dots \mid \mathbf{z}$

$\langle \text{const} \rangle ::= \langle \text{digit} \rangle \mid \langle \text{const} \rangle \langle \text{digit} \rangle$

$\langle \text{digit} \rangle ::= \mathbf{0} \mid \mathbf{1} \mid \dots \mid \mathbf{9}$

Note: as written, this grammar is not regular, but can be easily changed to an equivalent regular grammar

# Notes for the **Core** Interpreter Project

- Consider  $9 - 5 + 4$ 
  - What is the parse tree? What is the problem?
  - For ease of implementation, we will not fix this
    - But if we wanted to fix it, how can we?
- Manually writing a scanner for this language
  - Ad hoc approach (next slide)
  - Systematic approach: write regular expressions for all tokens, convert to an NFA, convert that to a DFA, minimize it, write code that mimics the transitions of the DFA (Section 2.2)
    - There exist various tools to do this automatically, but you should **not** use them for the project (will use in CSE 5343)

# Outline of a Scanner for **Core** (1/2)

- The parser asks the scanner for the next token
- Skip white spaces and read next character  $x$
- If  $x$  is ; , ( ) [ ] = + - \* return the corresponding token
- If  $x$  is : , read the next character  $y$ 
  - If  $y$  is not = , report error, else return the token for :=
- If  $x$  is ! , peek at the next character  $y$ 
  - If  $y$  is not = , return the token for !
  - If  $y$  is = , read it and return the token for !=
  - Peeking can be done easily in C, C++, and Java file I/O

## Outline of a Scanner for **Core** (2/2)

- If  $x$  is  $<$  , peek at the next character  $y$ 
  - If  $y$  is not  $=$  , return the token for  $<$
  - If  $y$  is  $=$  , read it and return the token for  $<=$
- Similarly when  $x$  is  $>$
- If  $x$  is a letter, keep reading characters; stop before the first not-letter-or-digit character
  - If the string is a keyword, return the keyword token
  - Else return token **id** with the string name attached
- If  $x$  is a digit, keep reading characters; stop before the first not-digit character
  - Return token **const** with the integer value attached