**The University of North Carolina at Chapel Hill**

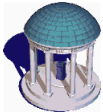**COMP 144 Programming Language Concepts**
**Spring 2002**

## Lecture 4:
## Syntax Specification
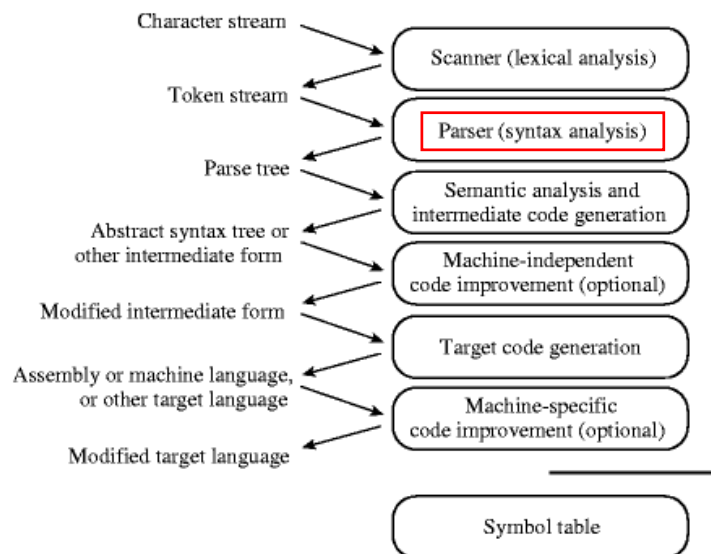
Felix Hernandez-Campos

Jan 16

COMP 144 Programming Language Concepts
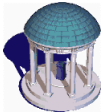Felix Hernandez-Campos

1

# Phases of Compilation



2

# Syntax Analysis

- Syntax:
  - Webster's definition: *1 a : the way in which linguistic elements (as words) are put together to form constituents (as phrases or clauses)*

- The syntax of a programming language
  - Describes its form
    - » *i.e.* **Organization of tokens** *(elements)*
  - Formal notation
    - » Context Free Grammars (CFGs)

COMP 144 Programming Language Concepts
Felix Hernandez-Campos

3

# Review: Formal definition of tokens

- A set of tokens is a set of strings over an alphabet
  - {read, write, +, -, *, /, :=, 1, 2, …, 10, …, 3.45e-3, …}

- A set of tokens is a *regular set* that can be defined by comprehension using a *regular expression*

- For every regular set, there is a *deterministic finite automaton* (DFA) that can recognize it
  - *i.e.* determine whether a string belongs to the set or not
  - Scanners extract tokens from source code in the same way DFAs determine membership
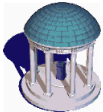
COMP 144 Programming Language Concepts
Felix Hernandez-Campos

4

# Review: Regular Expressions

- A regular expression (RE) is:
  - A single character
  - The empty string, $\varepsilon$
  - The <u>concatenation</u> of two regular expressions
    - » *Notation:* $RE_1$ $RE_2$ (*i.e.* $RE_1$ followed by $RE_2$)
  - The <u>union</u> of two regular expressions
    - » *Notation:* $RE_1$ | $RE_2$
  - The <u>closure</u> of a regular expression
    - » *Notation:* RE*
    - » * is known as the *Kleene star*
    - » * represents the concatenation of 0 or more strings

COMP 144 Programming Language Concepts
Felix Hernandez-Campos

5

# Review: Token Definition Example

- Numeric literals in Pascal
  - Definition of the token *unsigned_number*

$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$unsigned\_integer \rightarrow digit \; digit*$

$unsigned\_number \rightarrow unsigned\_integer \; ( \; ( \; . \; unsigned\_integer \; ) \mid \varepsilon \; )$
$( \; ( \; e \; ( \; + \mid - \mid \varepsilon ) \; unsigned\_integer \; ) \mid \varepsilon \; )$

- **Recursion is not allowed!**

COMP 144 Programming Language Concepts
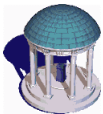Felix Hernandez-Campos

6

## **Exercise**

$digit \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$

$unsigned\_integer \rightarrow digit\ digit*$

$unsigned\_number \rightarrow unsigned\_integer\ (\ (\ .\ unsigned\_integer\ )\ |\ \varepsilon\ )$
$(\ (\ e\ (\ +\ |\ -\ |\ \varepsilon\ )\ unsigned\_integer\ )\ |\ \varepsilon\ )$

- Regular expression for
  – Decimal numbers

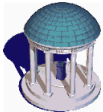$number \rightarrow \ldots$

## **Exercise**

$digit \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$

$unsigned\_integer \rightarrow digit\ digit*$

$unsigned\_number \rightarrow unsigned\_integer\ (\ (\ .\ unsigned\_integer\ )\ |\ \varepsilon\ )$
$(\ (\ e\ (\ +\ |\ -\ |\ \varepsilon\ )\ unsigned\_integer\ )\ |\ \varepsilon\ )$

- Regular expression for
  – Decimal numbers

$number \rightarrow (\ +\ |\ -\ |\ \varepsilon)\ unsigned\_integer\ (\ (\ .\ unsigned\_integer\ )\ |\ \varepsilon\ )$

## **Exercise**

$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$unsigned\_integer \rightarrow digit\ digit*$

$unsigned\_number \rightarrow unsigned\_integer\ ((\ .\ unsigned\_integer\ ) \mid \varepsilon\ )$
$((\ e\ (\ +\mid-\mid\varepsilon)\ unsigned\_integer\ ) \mid \varepsilon\ )$

- Regular expression for
  – Identifiers

$identifier \rightarrow ...$

COMP 144 Programming Language Concepts
Felix Hernandez-Campos

9

## **Exercise**

$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$unsigned\_integer \rightarrow digit\ digit*$

$unsigned\_number \rightarrow unsigned\_integer\ ((\ .\ unsigned\_integer\ ) \mid \varepsilon\ )$
$((\ e\ (\ +\mid-\mid\varepsilon)\ unsigned\_integer\ ) \mid \varepsilon\ )$

- Regular expression for
  – Identifiers

$identifier \rightarrow letter\ (\ letter \mid digit \mid \_\ )*$

$letter \rightarrow a \mid b \mid c \mid \ldots \mid z$

COMP 144 Programming Language Concepts
Felix Hernandez-Campos

10

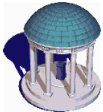# Context Free Grammars

- CFGs
  - Add recursion to regular expressions
    - » Nested constructions
  - Notation

    *expression → identifier | number | – expression*
    *| ( expression )*
    *| expression operator expression*

    *operator → + | – | * | /*

    - » **Terminal symbols**
    - » *Non-terminal symbols*
    - » Production rule (i.e. substitution rule)
      terminal symbol → terminal and non-terminal symbols

COMP 144 Programming Language Concepts
Felix Hernandez-Campos

11

# Backus-Naur Form

- Backus-Naur Form (BNF)
  - Equivalent to CFGs in power
  - CFG

    *expression → identifier | number | – expression*
    *| ( expression )*
    *| expression operator expression*

    *operator → + | – | * | /*

  - BNF

    ⟨expression⟩ → ⟨identifier⟩ | ⟨number⟩ | - ⟨expression⟩
    | ( ⟨expression⟩ )
    | ⟨expression⟩ ⟨operator⟩ ⟨expression⟩
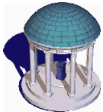
    ⟨operator⟩ → + | - | * | /

COMP 144 Programming Language Concepts
Felix Hernandez-Campos

12

# Extended Backus-Naur Form

- Extended Backus-Naur Form (EBNF)
  - Adds some convenient symbols
    » Union                                    |
    » Kleene star                              *
    » Meta-level parentheses          ( )
  - It has the same expressive power

COMP 144 Programming Language Concepts
Felix Hernandez-Campos

13

# Extended Backus-Naur Form

- Extended Backus-Naur Form (EBNF)
  - It has the same expressive power

**BNF**
  ⟨digit⟩ → 0
  ⟨digit⟩ → 1
  …
  ⟨digit⟩ → 9
  ⟨unsigned_integer⟩ → ⟨digit⟩
  ⟨unsigned_integer⟩ → ⟨digit⟩ ⟨unsigned_integer⟩

**EBNF**
  ⟨digit⟩ → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
  ⟨unsigned_integer⟩ → ⟨digit⟩ ⟨digit⟩*

COMP 144 Programming Language Concepts
Felix Hernandez-Campos

14

## Derivations

- A derivation shows how to generate a syntactically valid string
  - Given a CFG
  - Example:
    - » CFG

    $expression \rightarrow identifier \mid number \mid - expression$
    $\qquad\qquad \mid ( expression )$
    $\qquad\qquad \mid expression\ operator\ expression$
    $operator \rightarrow + \mid - \mid * \mid /$

    » Derivation of

    ```
    slope * x + intercept
    ```

COMP 144 Programming Language Concepts
Felix Hernandez-Campos

15

## Derivation Example

- Derivation of `slope * x + intercept`

  $expression \Rightarrow expression\ operator\ \underline{expression}$
  $\qquad\qquad \Rightarrow expression\ \underline{operator}\ intercept$
  $\qquad\qquad \Rightarrow \underline{expression} + intercept$
  $\qquad\qquad \Rightarrow expression\ operator\ \underline{expression} + intercept$
  $\qquad\qquad \Rightarrow expression\ \underline{operator}\ x\ + intercept$
  $\qquad\qquad \Rightarrow \underline{expression} * x\ + intercept$
  $\qquad\qquad \Rightarrow slope * x\ + intercept$

  $expression \Rightarrow^* slope * x\ + intercept$

  » Identifiers were not derived for simplicity
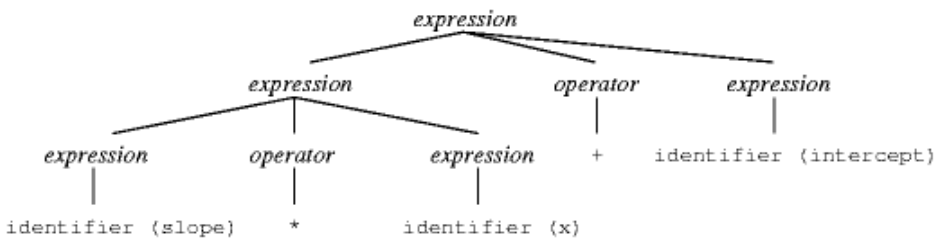
COMP 144 Programming Language Concepts
Felix Hernandez-Campos

16

# Parse Trees

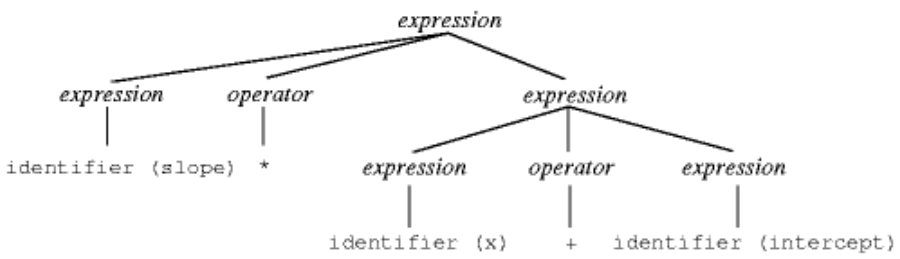- A parse is graphical representation of a derivation
- Example

17

# Ambiguous Grammars

- Alternative parse tree
  - same expression
  - same grammar



- This grammar is ambiguous

18

## Designing unambiguous grammars

- Specify more grammatical structure
  - In our example, left associativity and operator precedence
    - » 10 - 4 - 3 means $(10 - 4) - 3$
    - » 3 + 4 * 5 means $3 + (4 * 5)$

(1) *expression* $\longrightarrow$ *term* | *expression add_op term*
(2) *term* $\longrightarrow$ *factor* | *term mult_op factor*
(3) *factor* $\longrightarrow$ identifier | number | - *factor* | ( *expression* )
(4) *add_op* $\longrightarrow$ + | -
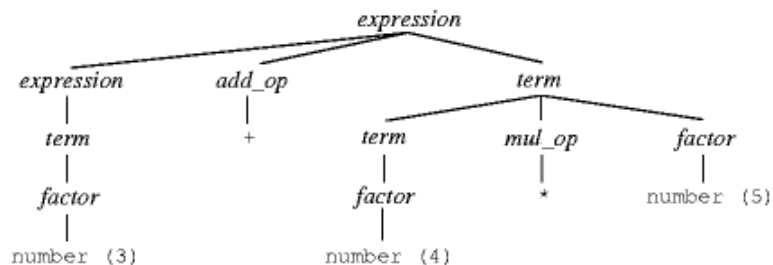(5) *mult_op* $\longrightarrow$ * | /

COMP 144 Programming Language Concepts
Felix Hernandez-Campos

19

## Example

- Parse tree for 3 + 4 * 5



- Exercise: parse tree for
  - 10 / 5 * 8 - 4 - 5

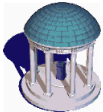COMP 144 Programming Language Concepts
Felix Hernandez-Campos

20

# Java Language Specification

- Available on-line
  - http://java.sun.com/docs/books/jls/second_edition/html/j.title.doc.html

- Examples
  - *Comments:* http://java.sun.com/docs/books/jls/second_edition/html/lexical.doc.html#48125
  - *Multiplicative Operators:* http://java.sun.com/docs/books/jls/second_edition/html/expressions.doc.html#239829
  - *Unary Operators:* http://java.sun.com/docs/books/jls/second_edition/html/expressions.doc.html#4990

COMP 144 Programming Language Concepts
Felix Hernandez-Campos

21

# Reading Assignment

- Scott's Chapter 2
  - Section 2.1.2
  - Section 2.1.3

- Java language specification
  - Chapter 2 (Grammars)
  - Glance at chapter 3
  - Glance at sections 15.17, 15.18 and 15.15

COMP 144 Programming Language Concepts
Felix Hernandez-Campos

22