


CSE 3302
Programming Languages




Syntax

Chengkai Li, Weimin He
Spring 2008

Lecture 3 - Syntax, Spring 2008 CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008 1


Phases of Compilation



[Programming Language Pragmatics, by Michael Scott]

Lecture 3 - Syntax, Spring 2008 CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008 2


Syntax and Semantics



- Defining a programming language:
 - Specifications of syntax
 - Syntax** – structure (form) of programs (the form a program in the language must take).
 - Specifications of semantics
 - Semantics** - the meaning of programs
- Precise definition, without ambiguity
 - Given a program, there is only one unique interpretation.

Lecture 3 - Syntax, Spring 2008 CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008 3


Purpose of Describing Syntax and Semantics




- Purpose
 - For **language designers**: Convey the design principles of the language
 - For **language implementers**: Define precisely what to be implemented
 - For **language programmers**: Describe the language that is to be used
- How to describe?
 - Natural language: ambiguous
 - Formal ways: especially for syntax

Lecture 3 - Syntax, Spring 2008 CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008 4

Scanning and Parsing




- **Lexical Structure**: The structure of tokens (words)
 - scanning phase (lexical analysis) : scanner/lexer
 - recognize tokens from characters
- **Syntactical Structure**: The structure of programs
 - parsing phase (syntax analysis) : parser
 - determines the syntactical structure



Lecture 3 - Syntax, Spring 2008 CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008 5

Tokens



Tokens (words): Building blocks of programs

- **Reserved words (keywords):** e.g., `if, while, int, return`
- **Literals/constants:**
 - **numeric literal:** `42`
 - **string literal:** `"hello"`
- **Special symbols:** e.g., `“,”, “<=”, “+”`
- **Identifiers:** e.g., `x24, monthly_balance, putchar`

Lecture 3 - Syntax, Spring 2008 CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008 6

Reserved words vs. Predefined identifiers



- **Reserved words:**
 - cannot be redefined.
 - e.g., `double if`; is illegal.
- **Predefined identifiers:**
 - have initial meaning
 - allow redefinition (not a good idea in practice)
 - e.g., `String, Object, System, Integer` in Java

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

7

Principle of Longest Substring



- `doif` vs. `do if`; `x12` vs. `x 12`
- The longest possible string of characters is collected into a single token.
- An exception: FORTRAN
 - `DO 99 I = 1.10` (the same as `DO99I=1.10`)
 - `DO 99 I = 1, 10`

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

8

White Space



- Principle of longest substring requires that tokens are separated by white space.
- **White space** (token delimiters):
 - Blanks, newlines, tabs
 - ignored except that they separate tokens
- **Free-format language:** format has no effect on the program structure
 - Most languages are free format
 - One exception: python

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

9

Indentation in Python



```
def perm(l):
    for i in range(len(l)):
        s = l[:i] + l[i+1:]
        p = perm(l[:i] + l[i+1:])
        for x in p:
            r.append(l[i:i+1] + x)
    return r
```

#error: first line indented
#error: not indented
#error: unexpected indent
#error: inconsistent dedent

```
def perm(l):
    for i in range(len(l)):
        s = l[:i] + l[i+1:]
        p = perm(l[:i] + l[i+1:])
        for x in p:
            r.append(l[i:i+1] + x)
    return r
```

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

10

Regular Expression



- A form for representing sets of strings
- Description of patterns of characters
- Basic operations:
 - Concatenation
 - Repetition
 - Selection

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

11

Regular Expression



Name	RE
epsilon	ϵ
symbol	<code>a</code>
concatenation	<code>AB</code>
selection	<code>A B</code>
repetition	<code>A*</code>

Shortcuts: `A+ = AA*` `A? = A| ϵ` `[a-z] = (a|b|...|z)`

`[a-z][a-z0-9]*` `(a|b)*aa(a|b)*` `[0-9]*(\.[0-9]*)?`

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

12

Tasks of a Scanner



- Recognizes keywords
- Recognizes special characters
- Recognizes identifiers, integers, reals, decimals, strings, etc.
- Ignores white spaces and comments

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

13

Scanner



regular expression description of the tokens
→ (Lex or JLex)
scanner of a language

- Example: Figure 4.1 (page 82)

Lecture 3 - Syntax, Spring 2008

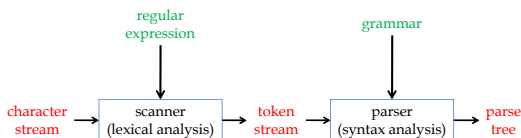
CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

14

Scanning and Parsing



- **Lexical Structure:** The structure of tokens (words)
- **Syntactical Structure:** The structure of programs



Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

15

Grammar



Example:

- (1) $sentence \rightarrow noun\text{-}phrase\ verb\text{-}phrase .$
- (2) $noun\text{-}phrase \rightarrow article\ noun$
- (3) $article \rightarrow a\ | \ the$
- (4) $noun \rightarrow girl\ | \ dog$
- (5) $verb\text{-}phrase \rightarrow verb\ noun\text{-}phrase$
- (6) $verb \rightarrow sees\ | \ pets$

Figure 4.2 (page 83)

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

16

Grammar



- Language: the programs (character streams) allowed
- **Grammar rules** (productions): "produce" the language
left-hand side, right-hand side
- **nonterminals** (structured names):
 $noun\text{-}phrase\ verb\text{-}phrase$
- **terminals** (tokens): $. \ dog$
- **metasymbols:** \rightarrow ("consists of") $|$ (choice)
- **start symbol:** the nonterminal that stands for the entire structure (sentence, program).
– $sentence$
 - E.g., $if\text{-}statement \rightarrow if\ (expression)\ statement\ else\ statement$

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

17

Grammars Produce Languages



- Language: the set of strings (of terminals) that can be generated from the start symbol by **derivation**:

```

sentence
⇒ noun-phrase verb-phrase . (rule 1)
⇒ article noun verb-phrase . (rule 2)
⇒ the noun verb-phrase . (rule 3)
⇒ the girl verb-phrase . (rule 4)
⇒ the girl verb noun-phrase . (rule 5)
⇒ the girl sees noun-phrase . (rule 6)
⇒ the girl sees article noun . (rule 2)
⇒ the girl sees a noun . (rule 3)
⇒ the girl sees a dog . (rule 4)
  
```

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

18

Context-Free Grammar

- Context-Free Grammars (CFG)
 - Noam Chomsky, 1950s.
 - Define *context-free languages*.
 - Four components:
 - terminals, nonterminals, one start symbol, productions (left-hand side: one single nonterminal)

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

19

What does "Context-Free" mean?

- Left-hand side of a production is always one single nonterminal:
 - The nonterminal is replaced by the corresponding right-hand side, no matter where the nonterminal appears. (i.e., there is no *context* in such replacement/derivation.)
- Context-sensitive grammar (context-sensitive languages)
- Why context-free?

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

20

Backus-Naur Form (BNF)

- A meta language used to describe CFG
- John Backus/Peter Naur: for describing the syntax of Algol60.
- BNF is formal and precise

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

21

BNF Example 1

- $$E \rightarrow \begin{array}{l} ID \\ NUM \\ E * E \\ E / E \\ E + E \\ E - E \\ (E) \end{array}$$

$$ID \rightarrow a \mid b \mid \dots \mid z$$

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

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

22

BNF Example 2

- $$S \rightarrow \begin{array}{l} \text{if } E \text{ then } S \text{ else } S \\ \text{begin } S \text{ L} \\ \text{print } E \end{array}$$

$$L \rightarrow \begin{array}{l} \text{end} \\ \text{; } S \text{ L} \end{array}$$

$$E \rightarrow \text{NUM} = \text{NUM}$$

$$\text{NUM} \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

23

Parse Tree

- Represents the derivation steps from start symbol to the string
- Given the derivations used in the parsing of an input sequence, a parse tree has
 - the start symbol as the root
 - the terminals of the input sequence as leaves
 - for each production $A \rightarrow X_1 X_2 \dots X_n$ used in a derivation, a node A with children $X_1 X_2 \dots X_n$

Lecture 3 - Syntax, Spring 2008

CSE3302 Programming Languages, UT-Arlington
©Chengkai Li, Weimin He, 2008

24

Parse Tree Example 1

CFG:
 $expr \rightarrow expr + expr \mid expr * expr \mid (expr) \mid number$
 $number \rightarrow number\ digit \mid digit$
 $digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

Input Sequence: **234**

Lecture 3 - Syntax, Spring 2008 CSE3302 Programming Languages, UT-Arlington ©Chengkai Li, Weimin He, 2008 25

Parse Tree Example 2

CFG:
 $expr \rightarrow expr + expr \mid expr * expr \mid (expr) \mid number$
 $number \rightarrow number\ digit \mid digit$
 $digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

Input Sequence: **3+4*5**

Lecture 3 - Syntax, Spring 2008 CSE3302 Programming Languages, UT-Arlington ©Chengkai Li, Weimin He, 2008 26

Abstract Syntax Tree

- Parse Tree:** still tedious, all terminals and nonterminals in a derivation are included in the tree.
- Abstract Syntax Tree:**
 - Remove "unnecessary" terminals and nonterminals
 - Still completely determine syntactic structure.

Lecture 3 - Syntax, Spring 2008 CSE3302 Programming Languages, UT-Arlington ©Chengkai Li, Weimin He, 2008 27

AST Example 1

Lecture 3 - Syntax, Spring 2008 CSE3302 Programming Languages, UT-Arlington ©Chengkai Li, Weimin He, 2008 28

AST Example 2

Lecture 2 - Syntax, Spring 2008 CSE3302 Programming Languages, UT-Arlington ©Chengkai Li, Weimin He, 2008 29