

# CSE 3302 Programming Languages

## Syntax

Chengkai Li  
Fall 2007

Lecture 3 - Syntax, Fall 2007      CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007      1

## Phases of Compilation

[Programming Language Pragmatics, by Michael Scott]

Lecture 3 - Syntax, Fall 2007      CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007      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, Fall 2007      CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007      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, Fall 2007      CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007      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 syntactic structure

Lecture 3 - Syntax, Fall 2007      CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007      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, Fall 2007      CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007      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, Fall 2007

CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007

7

## Principle of Longest Substring



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

Lecture 3 - Syntax, Fall 2007

CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007

8

## White Space



- Principle of longest substring requires that tokens are separated by white space.
- White space (token delimiters):
  - Blanks, newlines, tabs, comments
  - 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, Fall 2007

CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007

9

## Indentation in Python



```
def perm(l):                                #error: first line indented
for i in range(len(l)):                    #error: not indented
    s = l[:i] + l[i+1:]
    p = perm(l[:i] + l[i+1:]) #error: unexpected indent
    for x in p:
        r.append(l[:i+1] + x)
    return r                                #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+1] + x)
    return r
```

Lecture 3 - Syntax, Fall 2007

CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007

10

## Regular Expression



- Tokens: defined using grammar or regular expressions.
- Regular expressions: Useful for describing text patterns:
  - e.g., `grep`
- Example:
 

```
[0-9]+(\.[0-9]+)?
```

Lecture 3 - Syntax, Fall 2007

CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007

11

## Scanner



regular expression description of the tokens  
→ (lex/flex)  
scanner of a language

- Example: Figure 4.1 (page 82)

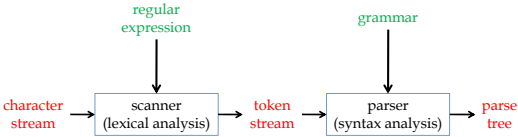
Lecture 3 - Syntax, Fall 2007

CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007

12

## Scanning and Parsing

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



```

graph LR
    CS[character stream] --> S[scanner  
(lexical analysis)]
    RE[regular expression] --> S
    S --> TS[token stream]
    TS --> P[parser  
(syntax analysis)]
    G[grammar] --> P
    P --> PT[parse tree]
    
```

Lecture 3 - Syntax, Fall 2007      CSE3302 Programming Languages, UT-Arlington ©Chengkai Li, 2007      13

## Grammar

Example:

- (1) *sentence* → *noun-phrase verb-phrase .*
- (2) *noun-phrase* → *article noun*
- (3) *article* → *a | the*
- (4) *noun* → *girl | dog*
- (5) *verb-phrase* → *verb noun-phrase*
- (6) *verb* → *sees | pets*

Figure 4.2 (page 83)

Lecture 3 - Syntax, Fall 2007      CSE3302 Programming Languages, UT-Arlington ©Chengkai Li, 2007      14

## Grammar

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

Lecture 3 - Syntax, Fall 2007      CSE3302 Programming Languages, UT-Arlington ©Chengkai Li, 2007      15

## 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 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)
- Define derivation? Figure 4.3 (page 85)

Lecture 3 - Syntax, Fall 2007      CSE3302 Programming Languages, UT-Arlington ©Chengkai Li, 2007      16

## CFG and BNF

- 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)
- Backus-Naur Forms (BNF)
  - John Backus/Peter Naur: for describing the syntax of Algol60.
  - BNF: equivalent to CFG
    - Use only metasymbols → | ( )


Lecture 3 - Syntax, Fall 2007      CSE3302 Programming Languages, UT-Arlington ©Chengkai Li, 2007      17

## 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, Fall 2007      CSE3302 Programming Languages, UT-Arlington ©Chengkai Li, 2007      18

## Recursion



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

Figure 4.4 (page 87)


$number \Rightarrow number digit \Rightarrow number digit digit \Rightarrow digit digit digit$   
 $\Rightarrow 2 digit digit \Rightarrow 23 digit \Rightarrow 234$

$expr \Rightarrow expr + expr \Rightarrow expr + expr * expr$   
 $\Rightarrow \dots \Rightarrow 3 + 4 * 5$

- How many programs(sentences) are in this language?
- How to derive  $(1+2) * 3 * (4*5)$ ? How many different derivations? Is the process lengthy and tedious?
- $((1+3))$  legal?  $+(4*5)$  legal?

Lecture 3 - Syntax, Fall 2007
CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007
19

## Parse Trees




- There can be many different derivations for capturing the same syntactic structure, and it can become tedious.

Question:  
Find another derivation for  $number \Rightarrow \dots \Rightarrow 234$

- Parse trees capture the intrinsic structure in a more intuitive way.

Lecture 3 - Syntax, Fall 2007
CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007
20

## Examples of Parse Trees



*number*

```

graph TD
    n1[number] --- n2[number]
    n1 --- d1[digit]
    n2 --- d2[digit]
    n2 --- d3[digit]
    d2 --- 2[2]
    d3 --- 3[3]
    d1 --- 4[4]
            
```

page 90

*expr*


```

graph TD
    e1[expr] --- e2[expr]
    e1 --- plus[+]
    e1 --- e3[expr]
    e2 --- n1[number]
    n1 --- d1[digit]
    d1 --- 3[3]
    e3 --- star[*]
    e3 --- e4[expr]
    e3 --- e5[expr]
    e4 --- n2[number]
    n2 --- d2[digit]
    d2 --- 4[4]
    e5 --- n3[number]
    n3 --- d3[digit]
    d3 --- 5[5]
            
```

page 91

Lecture 3 - Syntax, Fall 2007
CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007
21

## Parse Tree Nodes




- Root: start symbol
- Leaf nodes: terminals
- Internal nodes: nonterminals
- Children nodes of a node : A replacement (derivation) by the corresponding production rule  
e.g.,  $expr \rightarrow expr + expr$

```

graph TD
    e[expr] --- e1[expr]
    e --- plus[+]
    e --- e2[expr]
            
```

Lecture 3 - Syntax, Fall 2007
CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007
22

## Abstract Syntax Trees



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


```

graph TD
    n1[number] --- n2[number]
    n1 --- d1[digit]
    n2 --- d2[digit]
    n2 --- d3[digit]
    d2 --- 2[2]
    d3 --- 3[3]
    d1 --- 4[4]
    
    a1[2] --- a2[3]
    a1 --- a3[4]
            
```

page 91

Lecture 3 - Syntax, Fall 2007
CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007
23

## Another Example of Abstract Syntax Tree



```

graph TD
    e1[expr] --- e2[expr]
    e1 --- plus1[+]
    e1 --- e3[expr]
    e2 --- n1[number]
    n1 --- d1[digit]
    d1 --- 3[3]
    e3 --- star[*]
    e3 --- e4[expr]
    e3 --- e5[expr]
    e4 --- n2[number]
    n2 --- d2[digit]
    d2 --- 4[4]
    e5 --- n3[number]
    n3 --- d3[digit]
    d3 --- 5[5]
    
    plus2[+] --- 3[3]
    plus2 --- star[*]
    star --- 4[4]
    star --- 5[5]
            
```

Lecture 3 - Syntax, Fall 2007
CSE3302 Programming Languages, UT-Arlington  
©Chengkai Li, 2007
24

## Syntax-Directed Semantics (or Semantics-Directed Syntax)

- $sentence \rightarrow noun\text{-}phrase\ verb\text{-}phrase \ .$

*What is meant? (The subject performs some action.)*

- $expr \rightarrow expr + expr$

*What is meant? What is the computation executed?*