

# CS480

---

## Syntax Analysis

Ch 4 p 159-195

February 18, 2013

# Parser

- Receive stream of tokens from lexer
- Verify the stream is grammatically correct
- Perform semantic actions based on parse tree
  - Semantic checking

# Types of Parsing

- CYK (Cocke-Younger-Kasami) & Earley's algorithm
  - Can parse any grammar
    - What grammars are hard to parse?
  - JFLAP has a CYK parser
  - Too inefficient to actually use
  - Your book has references in the bibliography



# CS310 Problems

- Describe (in English) the language denoted by the regular expression  $((\epsilon|0)1^*)^*$
- Write regular definitions for:
  - all strings that begin with an aa
  - all strings that contain aa
  - all strings that do not contain aa
  - All are over the alphabet  $\{a,b\}$ .
- Construct an NFA for the regular expression  $((\epsilon| a)b^*)^*$

# CFGs

$\text{expr} \rightarrow \text{expr op expr} \mid (\text{expr}) \mid$   
 $\text{number} \mid \text{id}$

$\text{op} \rightarrow + \mid - \mid *$

- Backus-Naur Form

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle \langle \text{op} \rangle \langle \text{expr} \rangle \mid$   
 $(\langle \text{expr} \rangle) \mid \text{NUMBER}$

$\langle \text{op} \rangle ::= + \mid - \mid *$

# Notation from the Book

- Terminals

---

- Nonterminals
- String of terminals
- Greek Letters
- Alternate Forms
- Start production

# Derivations

$\Rightarrow$

- can derive with one application of a production

$\Rightarrow^*$

- can derive with zero or more applications of any productions

$\Rightarrow_{lm}$

Sentential form

$E \rightarrow ( E ) \mid a$   
 $E \rightarrow E + E$

Does  $E \Rightarrow^* ((a))$ ?

Does  $E \Rightarrow ((a))$ ?

Does  $E \Rightarrow^* (a)(a)$ ?



# Grammars

- $G1: A \rightarrow Aa \mid a$
  - $G2: B \rightarrow aB \mid a$
- 
- Do  $G1$  and  $G2$  describe the same language?
  - Are both  $G1$  and  $G2$  equivalent to  $a^*$ ?
  - Are they ambiguous?
    - How fix?
  - Right or Left recursive?
    - What problems could arise?
  - Does  $A \Rightarrow^* \epsilon$

# More...

- Give a CFG which generates sequences of one or more statements (s) separated by ;
  - (i.e.  $L(G) = \{s\ s;\ s\ s;\ s;\ s\ \dots\}$ )
- Give a CFG which generates sequences of one or more statements where the semicolon is a terminator and not a separator (i.e.  $L(G) = \{s;\ s;\ s;\ s;\ s;\ s;\ \dots\}$  )

# Parsing!

**expr -> expr op expr | ( expr ) |  
number**

**op -> + | - | \* | /**

- Problem?

1 + 3 \* 8

Left most? Right most?

- Ambiguity:

- Get rid of it      **OR**
- Use rules to limit its impact

# More..

**expr -> expr op expr | term**

**op -> + | - | \***

**term -> number**

- Ambiguous?
  - Why or why not?
  - Precedence?

# More still...

`stmt -> ifstmt | other`

`ifstmt -> if ( expr ) stmt |  
if ( expr ) stmt else stmt`

`expr -> T | F`

- Thoughts?
- Fixes?

# Immediate Left Recursion

- Immediate Left Recursion

$$E \rightarrow E + T \mid T$$
$$T \rightarrow T * F \mid F$$
$$F \rightarrow (E) \mid \text{id}$$

Differences?  
Why is this important?

- Nonimmediate Left Recursion:

$$S \rightarrow Aa \mid b$$
$$A \rightarrow Ac \mid Sd \mid e$$

How do you remove each type?

# Practice

**S** -> **Ba** | **b**

**B** -> **Sa** | **a**

- What is the language?
- Eliminate all the left recursion
  - Algorithm 4.1 on p 177