Lecture 4: Syntactic Analysis

COMP 524 Programming Language Concepts Stephen Olivier January 27, 2009

Based on notes by N. Fisher, F. Hernandez-Campos, and D. Stotts



Goal of Lecture



• The main task of parsing is to **identify the syntax**.



Review: Regular Expression Rules

• A RE consist of:

- A character (e.g., "0", "1", ...)
- The empty string (i.e., "ε")
- Two REs next to each other (e.g., "non_negative_digit digit") to denote concatenation.
- Two REs separated by "|" next to each other (e.g., "non_negative_digit | digit") to denote one RE or the other.
- An RE followed by "*" (called the Kleene star) to denote zero or more iterations of the RE.
- Parentheses (in order to avoid ambiguity).

Review: Regular Expression Rules

• A RE consist of:

- A character (e.g., "0", "1", ...)
- The A RE is NEVER defined in terms of itself!
 Two Thus, REs cannot define recursive statements.
- Two REs separated by "|" next to each other (e.g., "non_negative_digit | digit") to denote one RE or the other.
- An RE followed by "*" (called the Kleene star) to denote zero or more iterations of the RE.
- Parentheses (in order to avoid ambiguity).

The University of North Carolina at Chapel Hill

igit") to

Review: Regular Expression Rules

• A RE consist of:

- A character (e.g., "0", "1", ...)
- The
 Two dend
 For example, REs cannot define arithmetic expressions with parentheses.
- Two REs separated by "|" next to each other (e.g., "non_negative_digit | digit") to denote one RE or the other.
- An RE followed by "*" (called the Kleene star) to denote zero or more iterations of the RE.
- Parentheses (in order to avoid ambiguity).

Context-Free Grammars

• Context-Free Grammars (CFGs) are similar to REs except that they can handle recursion.

Arithmetic expression with parentheses

 $expr \rightarrow id | number | - expr | (expr) | expr op expr$

















- Technically, the Kleene star (*) and parentheses are not allowed under the CFG rules, called Backus-Naur Form (BNF).
- However, for convenience, we will use Extended BNF that includes the Kleene star, parentheses, and the Kleene Plus (⁺), which stands for "one or more iterations."



EBNF Example.

• The Kleene star and parentheses can be we written as follows

$$id_list \rightarrow id (, id)^*$$

$$id_list \rightarrow id$$

$$id_list \rightarrow id_list, id$$



• A derivation is "a series of replacement operations that derive a string of terminals from the start symbol."







Parse Tree

• A parse tree is the graphical representation of the derivation.



Pa This parse tree constructs the formula (slope*x) + intercept A parse tree is the graphical representation of the derivation.



Parse Tree

• A parse tree is the graphical representation of the derivation.



Parse Tree

Lets try deriving "2*a*b+c"

• A parse tree is the graphical representation of the derivation.



Parse Tree (Ambiguous)

• This grammar is ambiguous and can construct the following parse tree.



Parse Tree (Ambiguous)

• This grammar is ambiguous and can construct the following parse tree.





Parse Tree (Ar Lets try deriving "2*a*b+c"

• This grammar, is ambiguous and can construct the following parse tree.



- The problem with our original grammar was that we did **not fully express the grammatical structure** (i.e., associativity and precedence).
- To create an unambiguous grammar, we need to fully specify the grammar.

$$expr \rightarrow term | expr add_op term$$

$$term \rightarrow factor | term mult_op factor$$

$$factor \rightarrow id | number | - factor | (expr)$$

$$add_op \rightarrow + | -$$

$$mult_op \rightarrow * | /$$



Ihe

- The problem with our original grammar was that we did **not fully express the grammatical structure (in a secondativity and** precedence Gives precedence to multiply
- To create an unampiguous grammar, we need to fully specify the grammar.

$$expr \rightarrow term | expr add_op term$$

```
term \rightarrow factor | term mult_op factor
```

$$add_op \rightarrow + | -$$

 $mult_op \rightarrow * | /$



The Unive



- The problem did not fully express the Lets try deriving "3*4+5*6+7" ty and precedence).
- To create an unambiguous grammar, we need to fully specify the grammar.

$$expr \rightarrow term | expr add_op term$$

$$term \rightarrow factor | term mult_op factor$$

$$factor \rightarrow id | number | - factor | (expr)$$

$$add_op \rightarrow + | -$$

$$mult_op \rightarrow * | /$$

I he



- The problem did not fully **express the** Lets try deriving "3*4" & "3+4" ty and precedence).
- To create an unambiguous grammar, we need to fully specify the grammar.

$$expr \rightarrow term | expr add_op term$$

$$term \rightarrow factor | term mult_op factor$$

$$factor \rightarrow id | number | - factor | (expr)$$

$$add_op \rightarrow + | -$$

$$mult_op \rightarrow * | /$$



I he







Java Spec

- Available on-line
 - <u>http://java.sun.com/docs/books/jls/second_edition/html/j.title.doc.html</u>
- Examples
 - Comments: <u>http://java.sun.com/docs/books/jls/</u> <u>second_edition/html/lexical.doc.html#48125</u>
 - Multiplicative Operators: <u>http://java.sun.com/docs/books/jls/</u> <u>second_edition/html/expressions.doc.html#239829</u>
 - Unary Operators: <u>http://java.sun.com/docs/books/jls/</u> <u>second_edition/html/expressions.doc.html#</u>4990


LL and LR Derivation

- CFGs can be parsed in O(n^3) time, where n is length of the program.
- This is too long for most code; however, there are two types of grammars that can be parsed in linear time, i.e., O(n).
 - LL: "Left-to-right, Left-most derivation"
 - LR "Left-to-right, Right-most derivation"





- LL parsers are top-down, i.e., they identify the nonterminals first and terminals second.
 - LL grammars are grammars that can be parsed by an LL parser.













- LR parsers are bottom-up, i.e., they discover the terminals first and non-terminals second.
 - LR grammars are grammars that can be parsed by an LR parser.
 - All LL grammars are LR grammars but not vice versa.

















The University of North Carolina at Chapel Hill









A better bottom-up grammar

This grammar limits the number of "suspended" non-terminals.





A better bottom-up grammar





A better bottom-up grammar

However, it **cannot** be parsed by LL (top-down) parser. Since when the parser discovers an "id" it does not "know" the number of "id_list_prefixs"

 $id_list \rightarrow id_list_prefix;$ $id_list_prefix \rightarrow id_list_prefix, id$

id_list_prefix → id



A better bottom-up gramm

 $id_list \rightarrow id_list_prefix;$

id_list_prefix → *id_list_prefix*, id

id_list_prefix → id

Both of these are valid break downs, but we don't know which one. Therefore, is **NOT** a valid **LL grammar**, but it is a valid **LR** grammar.







This grammar (for the calculator) unlike the previous calculator grammar is an LL grammar, because when an "id" is encountered we know exactly where it belongs.



Let's try "c := 2^*A+B "

Let's try "2*A+B"



term → factor | term mult_op factor

factor → id | number | - factor | (expr)

add_op
$$\rightarrow + | -$$

$$mult_op \rightarrow * | /$$

The University of North Carolina at Chapel Hill







Recursive Descent & LL Parse Table

• There are two ways to code a parser for LL grammars:

- Recursive Descent, which is a recursive program with case statements that correspond to each one-to-one to nonterminals.
- LL Parse Table, which consist of an iterative driver program and a table that contains all of the nonterminals.





Recursive Descent



procedure stmt()

```
case in_tok of
```

```
id: match(id); match(:=); expr()
read: match(read); match(id)
write: match(write); expr()
```

else

return error

procedure match(expec)
if in_tok = expec
 consume in_tok
 else
 return error

Recursive Descent



Recursive Descent



procedure stmt()

```
case in_tok of
```

```
id: match(id); match(:=); expr()
read: match(read); match(id)
write: match(write); expr()
```

else

return error

procedure match(expec)
if in_tok = expec
 consume in_tok
 else
 return error



- Three functions allow us to label the branches
 - FIRST(a): The terminals (and ε) that can be the first tokens of the non-terminal symbol a.
 - FOLLOW(A): The terminals that can follow the terminal or nonterminal symbol A
 - PREDICT(A → a): The terminals that can be the first tokens as a result of the production A → a



- **FIRST**(*program*) = {id, read, write, \$\$}
- FOLLOW(*program*) = $\{\epsilon\}$
- PREDICT(program →
 stmt_list \$\$) = {id,
 read, write, \$\$}
- FOLLOW(id) = {+, -,
 *,/,), :=, id, read,
 write, \$\$}



- FIRST(factor_tail)
 ={ *, /, ε}
- FOLLOW(factor_tail) =
 {+,-,),id,read,
 write,\$\$}
- PREDICT(factor_tail → m_op factor factor_tail)
 = {*, /}
- PREDICT(factor_tail → ε)
 = {+, -,), id, read,
 write, \$\$}



These are all of the PREDICT() values from every production.



Recursive D

Constructing FIRST, FOLLOW, and PREDICT

- To construct the FIRST, FOLLOW, and PREDICT tables we iterate through the grammar building on knowledge.
 - First, we define all of the "obvious" FIRST and FOLLOW values
 - For example, \$\$ ∈ FOLLOW (stmt_list) and {id, read, write}∈
 FIRST(stmt)
 - Next, we build on this,
 - For example, {id, read, write}∈ FIRST(*stmt_list*) since *stmt_list* can begin with *stmt* and {id, read, write}∈ FIRST(*stmt*)
 - We then continue on until we get **no more knowledge**.



Let's try making tables for this grammar

expr → term | expr add_op term

term → factor | term mult_op factor

add_op
$$\rightarrow + | -$$

$$(mult_op \rightarrow *|/)$$

The University of North Carolina at Chapel Hill

Table-Driven





Parse Stack	Input Stream
program	read A read B
stmt_list \$\$	read A read B
stmt stmt_list \$\$	read A read B
read id <i>stmt_list</i> \$\$	read A read B
id stmt_list \$\$	A read B
Writing an LL(1) Grammar--Left Recursion

• Left recursion is where the leftmost symbol is a recursive non-terminal symbol.



• This can cause a grammar **not to be LL(1)**.

• It is desirable for LR grammars.

Writing an LL(1) Grammar-- Eliminating Left Recursion

To eliminate left recursion replace it with right recursion.





Writing an LL(1) Grammar--Common Prefix

• Common prefixes occur when there is more than one prefix for a given nonterminal.



• Again, this causes a grammar not to be LL(1).



Writing an LL(1) Grammar--Left factoring

• Common prefixes To get rid of common prefixes we use a technique called left factoring.

<pre>stmt → id := expr</pre>	<pre>stmt → id stmt_list_tail</pre>
<pre>stmt → id (arguments)</pre>	<pre>stmt_list_tail → := expr (arguments)</pre>



- Even if left recursion and common prefixes don't exist a language may not be LL(1).
- In Pascal, there is the problem that an else statement in if-then-else statements is optional. Because we don't know which if to match else to.





- In Pascal there is NO LL(1) parser that can handle this problem.
- Even though a proper LR(1) parser can handle this, it may not handle it in a method the programmer desires.





The University of North Carolina at Chapel Hill

• Thus, to write this code correctly (based on indention) "begin" and "end" statements must be added.





