

## Lecture 2 Syntax Analysis

Dr. Wilson Rivera

ICOM 4036: Programming Languages Electrical and Computer Engineering Department University of Puerto Rico

Some slides adapted from Sebesta's textbook

# Lecture Outline

- Introduction to syntax analysis
- Regular expressions
- Finite Automata
- Context free grammars
- Lexical Analysis
- Parsing

#### Introduction





## Introduction

- Syntax: the form or structure of the expressions, statements, and program units
- Semantics: the meaning of the expressions, statements, and program units
- Syntax and semantics provide a language's definition
  - Syntactic specification of a programming language
  - Construction of efficient parsers
  - Error detection
  - Extension language and new constructs

# Introduction



## **Regular Expressions**

Regular Expressions: A concise notation for regular sets

- (1)  $\Phi$  denotes the regular set  $\Phi$ .
- (2)  $\Lambda$  denotes the regular set { $\Lambda$ }.
- (3)  $\alpha$  denotes the regular set { $\alpha$ }.

(4) If p and q are regular expressions denoting the regular sets P and Q respectively, then

- (a) (p | q) denotes  $P \cup Q$  (Union)
- (b) (pq) denotes P Q (Concatenation)
- (c) (p)\* denotes P\* (Kleene Closure)

(5) Nothing else is a regular expression.

Notation:

p+ = p\* p (non reflexive closure)

#### **Regular Expressions**

a|b (a|b)(a|b) a\* a\*b (a|b)\*

RE

**Regular language Description** 

{a,b}
{aa, ab, ba, bb}
{aa, aa, aaa, aaa, ...}
{b,ab, aab, aaab, ...}
{e, a, aa, aaa, ..., b, bb, bbb...
ab, abb,.... aab, aabb, ....}

## **Regular Expressions**

Regular Expression

 (10\*1(0|1)\*) | (01\*0(0|1)\*)



#### Deterministic Finite Automaton (DFA)

$$\mathbf{M} = (\mathbf{Q}, \boldsymbol{\Sigma}, \boldsymbol{\delta}, \mathbf{q}_0, \mathbf{F})$$

where

(1) Q is a finite non-empty set of states

- (2)  $\Sigma$  is a finite set of input symbols
- (3)  $q_0 \in Q$  (initial state)
- (4)  $F \in Q$  (final states)

(5)  $\delta$  is a partial mapping from Q x  $\Sigma$  to Q (transition function: move function)

#### Non-deterministic Finite Automaton (NFA)

- May have a choice of moves, i.e.  $\delta$  is a mapping from Q x  $\Sigma$  to  $2^Q$
- Also allows  $\in$ -transitions, i.e.,  $\delta (q, \in) \subseteq Q$

#### (a|b)\*abb = { abb, aabb, babb, aaabb, bbabb,...}



### **DFA Example**



#### DFA:

- 1. No state has an e-transition
- 2. For each state *S* and input symbol *a*, there is at most one edge labeled *a* leaving *S*.

### Finite Automata Construction

- For every NFA  $M_1$  there is a DFA  $M_2$  for which  $L(M_2) = L(M_1)$
- Thompson's Construction:
  - Systematically generate an NFA for a given Regular Expression.
- Subset construction algorithm:
  - Converts an NFA to an equivalent DFA
  - *Key:* identify sets of states of NFA that have similar behavior, and make each set a single state of the DFA.

#### **Thomson's Construction**





6 Combine single NFAs for complex structures

Example : (a|b)\*abb







#### Subset Construction

Operation

ε-closure(s)

ε-closure(T)

Move(T,a)

#### Description

Set of NFA states reachable from an NFA state s on e-transitions along

Set of NFA states reachable from some NFA state s in T on e-transitions along

Set of NFA states reachable from some NFA state set with a transition on input symbol a

#### **Subset Construction**

Initially, e-closure  $(s_0)$  is the only states in D and it is unmarked while there is an unmarked state T in D do mark T; for each input symbol a do U:=e-closure(move(T,a));if U is not in D then add U as an unmarked state to D Dtran[T,a] := U;end(for) end(while)

#### Example

- e-closure(0)={0,1,2,4,7}=A
- e-closure(move(A,a))=
   e-closure({3,8})={1,2,3,4,6,7,8}=B
   Thus, Dtran[A,a]=B
- e-closure(move(A,b))=
   e-closure({5})={1,2,4,5,6,7}=C
   Thus, Dtran[A,b]=C

- e-closure(move(B,a))=
   e-closure({3,8})=B
   Thus, Dtran[B,a]=B
- e-closure(move(B,b))=
   e-closure({5,9})={1,2,4,5,6,7,9}=D
   Thus, Dtran[B,b]=D

- e-closure(move(C,a))=
   e-closure({3,8})=B
   Thus, Dtran[C,a]=B
- e-closure(move(C,b))=
   e-closure({5})=C
   Thus, Dtran[C,b]=C

- e-closure(move(D,a))=
   e-closure({3,8})=B
   Thus, Dtran[D,a]=B
- e-closure(move(D,b))=
   e-closure({5,10})={1,2,4,5,6,7,10}=E
   Thus, Dtran[D,b]=E

- e-closure(move(E,a))=
   e-closure({3,8})=B
   Thus, Dtran[E,a]=B
- e-closure(move(E,b))=
   e-closure({5})=C
   Thus, Dtran[E,b]=C



# From a DFA to a minimal DFA

#### Hopcroft Algorithm

- 1. Construct an initial partition P of set of states with two groups: The accepting states F and the nonaccepting states S-F
- 2. Apply the following procedure to P to construct a new partition  $P_{new}$

#### for each group G of P do

divide G into subgroups such that two states

of G are in the same subgroup if and only if for all input symbol a, states s and t have transitions on a to states in the same group of P.

Replace G in  $P_{new}$  by the set of all subgroup formed.

- 3. If  $P_{new} = P$  let  $= P_{final} = P$  and go to 4. Otherwise repeat (2) with  $P = P_{new}$ .
- 4. Choose one state in each group of the partition  $P_{\text{final}}$  as the representative for that group
- 5. Contruct the new transition table by replacing the states in a group by the respresentative

## From a DFA to a minimal DFA

- 1.  $P = \{ (ABCD), (E) \}$
- 2.  $P_new = \{(ABC), (D), (E)\}$   $D \rightarrow E$
- 3.  $P_new = \{(AC), (B), (D), (E)\} \quad B \rightarrow D$
- 4. P\_final={(AC), (B), (D), (E)}
- 5. We choose A as the representative for the subgroup (AC)



#### **RE to NFA: Example**

 Desing a NFA for the following regular expression

(a|b)\*a(a|b)



### NFA to DFA: Example

#### •Reduce the NFA to a DFA

- E(0)={0 1 2 4 7}=A
- E(move A a) =  $E(3 \ 8) = \{1 \ 2 \ 3 \ 4 \ 6 \ 7 \ 8 \ 9 \ 11\} = B$
- E(move A b)=E(5)={1 2 4 5 6 7}=C
- E(move B a)=E(3 8 10)={1 2 3 4 6 7 8 9 10 13}=D
- E(move B b)=E(5 12)={1 2 4 5 6 7 12 13}=E
- E(move C a) = E(3 8) = B
- E(move C b)= E (5) = C
  - E(move D a) = E(3 8 10) = D
- E(move D b) = E(5) = C
- E(mov E a)=E(3 8)=B
- E(mov E b) = E(5) = C

### NFA to DFA



# Regular Expressions (summary)

- REs are commonly used to define search patterns and the lexical structure of programming Languages
- REs are a convenient mechanism for describing languages but they do not give a method for recognizing tokens
  - Construct a NFA for the regular expression.
  - Convert the NFA to an equivalent DFA.
  - Minimize the number of states in the DFA

### More on Regular Expressions

- {a<sup>n</sup>b<sup>n</sup> : n>0} is not regular since the number of a's controls the number of b's
  - This type of operation is not allowed according to the definition of regular expressions
- Many real world REs engines implement features that cannot be described by REs theory
  - The Python regular expression  $r"(a^*)b \setminus 1"$  recognizes  $a^n ba^n$ 
    - \1 matches the same string matched by the parenthesized subexpression.

# Chomsky Hierarchy

Туре-0	unrestricted	Turing machine (TM)
Type-1	Context sensitive	Linear bounded TM
Type-2	Context free	Non deterministic pushdown automaton
Type-3	Regular	Finite State Automaton



# Chomsky Hierarchy

Context-free Grammars	×.
Regular Grammars	
Regular Expressions	Pushdown Automata
Finite State Automata	

For Regular Grammars there is always a corresponding deterministic automaton. For Context-free grammars there is none unless the grammar is unambiguous.
#### Context-free Grammar

#### A formal grammar is define as

- G = (T, N, S, P)
- 1) <u>Tokens</u> (terminals)
- 2) <u>Constructs</u> (nonterminals)
- *3) <u>Starting nonterminal</u>* (main construct)
- 4) <u>Productions</u> (rules that define nonterminal in terms of sequences of terminals and nonterminals)
- A context-free grammar (CFG) is a formal grammar in which every production rule is of the form  $V \rightarrow w$  where
- V is a single nonterminal symbol
- w is a string of terminals and/or non-terminals (w can be empty).

The language of a grammar G=(T;N;R; S) usually denoted by L(G), is defined as L(G)={s | S-> \* s, s in T\*}

## **BNF and Context-Free Grammars**

- Context–Free Grammars
  - Developed by Noam Chomsky in the mid-1950s
  - Meant to describe the syntax of natural languages
  - Context free means replacing non-terminals in any order (i.e., regardless of context) produces same result (as long as you use same productions).
  - If we use the full power of the context-free languages we get compilers which in general are inefficient, and probably not so good in handling erroneous input
- Backus-Naur Form (1959)
  - Invented by John Backus to describe Algol 58
  - BNF is equivalent to context-free grammars

#### **BNF Fundamentals**

•A rule has a left-hand side (LHS), which is a nonterminal, and a right-hand side (RHS), which is a string of terminals and/or nonterminals

- Terminals are lexemes
- Nonterminals are often enclosed in angle brackets

- Examples of BNF rules: <ident\_list> → identifier | identifier, <ident\_list> <if stmt> → if <logic expr> then <stmt>

- •Grammar: a finite non-empty set of rules
- •A *start symbol* is a special element of the nonterminals of a grammar



 An abstraction (or nonterminal symbol) can have more than one RHS

> <stmt> → <single\_stmt> | begin <stmt\_list> end

## **Describing Lists**

Syntactic lists are described using recursion

<ident\_list> -> ident | ident, <ident list>

 A derivation is a repeated application of rules, starting with the start symbol and ending with a sentence (all terminal symbols)

#### Example: Grammar sentences

Consider the following grammar

 $\langle s \rangle \rightarrow \langle A \rangle a \langle B \rangle b$  $\langle A \rangle \rightarrow \langle A \rangle b \mid b$  $\langle B \rangle \rightarrow a \langle B \rangle \mid a$ 

Which of the following sentences are in the language generated by this grammar?

- 1. baab
- 2. bbbab
- 3. bbaaaa
- 4. bbaab

#### Example: Regular Grammar

Regular Grammar

 $\langle S \rangle \rightarrow a \langle S \rangle$  $\langle S \rangle \rightarrow b \langle A \rangle$  $\langle A \rangle \rightarrow e \mid c \langle A \rangle$ 

 Regular expression a\*bc\*

## Non Regular Grammar

- Context free grammar  $\langle S \rangle \rightarrow a \langle S \rangle b \mid ab$
- {a<sup>n</sup>b<sup>n</sup> : n>0} is not regular since the number of a's controls the number of b's
  - This type of operation is not allowed according to the definition of regular expressions

# Context free grammars

- If we remove the regular-grammar restrictions, we get a class of grammars known as the Type 2 or context-free grammars. Such grammars can "count" arbitrarily high, at least under the right circumstances.
  - For example, here is a language of correctly nested parentheses
    - S  $\rightarrow$  S '(' S ')' | e
    - It recognizes the empty string, '()', '()()', '(())', '(()(()))'

#### **Example: Grammar Derivations**

```
<program> -> <stmts>
```

```
<stmts> \rightarrow <stmt> | <stmt> ; <stmts>
```

<stmt $> \rightarrow <$ var> = <expr>

 $\langle var \rangle \rightarrow a \mid b \mid c \mid d$ 

 $\langle expr \rangle \rightarrow \langle term \rangle + \langle term \rangle | \langle term \rangle - \langle term \rangle$ 

<term> → <var> | const

#### **Example: Grammar Derivation**



- => <var> = <expr>
- => a = <expr>
- => a = <term> + <term>
- => a = <var> + <term>
- => a = b + <term>
- => a = b + const

#### Derivations

- Every string of symbols in a derivation is a sentential form
- A *sentence* is a sentential form that has only terminal symbols
- A *leftmost derivation* always choose leftmost non-terminals to be expanded.

#### Parse Tree

A hierarchical representation of a derivation



## Ambiguity in Grammars

- A grammar is *ambiguous* if and only if it generates a sentential form that has two or more distinct parse trees
- Checking a grammar is ambiguous is undecidable
  - Reduction from Post Correspondence problem
- Ambiguous grammars lead to ambiguous semantics!

9-7

#### An Ambiguous Expression Grammar



#### An Unambiguous Expression Grammar

 If we use the parse tree to indicate precedence levels of the operators, we cannot have ambiguity

<expr> → <expr> + <term> | <term> <term> → <term> \* const| const



#### Associativity of Operators

- Operator associativity can also be indicated by a grammar
- <expr> -> <expr> + <expr> | const (ambiguous)
  <expr> -> <expr> + const | const (unambiguous)



#### **Extended BNF**

- Optional parts are placed in brackets []
  <proc\_call> -> ident [(<expr\_list>)]
- Alternative parts of RHSs are placed inside parentheses and separated via vertical bars

 $< term > \rightarrow < term > (+ | -) const$ 

 Repetitions (0 or more) are placed inside braces { }

<ident> -> letter {letter|digit}

#### **BNF and EBNF**

#### • BNF

• EBNF

<expr> → {<expr>(+ | -)} <term> <term> → {<term>(\* | /)}<factor>

#### Example: Grammar Definition

#### Function: type IDENTIFIER LEFT\_PARENT optional\_arguments RIGHT\_PAREN LEFT\_PAREN LEFT\_CURLY statements RIGHT\_CURLY

- Type: TYPE\_KEYWORD | user\_defined\_type
- Optional-arguments: null | argument | argument COMMA arguments
- Statements: statement | statement statements
- Statement: declaration | equals\_expr | .....
- Declaration: type IDENTIFIER optional\_declarations SEMICOLON
- equals\_expr: IDENTIFIER EQUAL\_SIGN expression
- Expression: plus\_expr | equals\_expr | ...
- plus\_expression: expression PLUS\_SIGN expression

#### Example: Java language specification

PrimitiveType: NumericType boolean

*NumericType: IntegralType FloatingPointType* 

IntegralType: one of byte short int long char

FloatingPointType: one of float double

- The syntax analysis portion of a language processor nearly always consists of two parts:
  - A low-level part called a *lexical analyzer* (mathematically, a finite automaton based on a regular grammar)
  - A high-level part called a *parser* (mathematically, a push-down automaton based on a context-free grammar)

# Lexical Analysis

- Breaks the text down into the smallest useful atomic units, known as tokens
  - while throwing away (or at least, putting to one side) information, such as white space and comments
- Groups characters into tokens to facilitate the parsing process
  - X+Y \*
  - X+Y++

# Lexical Analysis

- Approaches to building a lexical analyzer:
  - Write a formal description of the tokens and use a software tool that constructs table-driven lexical analyzers given such a description
  - Design a state diagram that describes the tokens and write a program that implements the state diagram

#### Example: Lexical Analysis

#### int main ()

```
{
  float foo, bar, baz;
  foo = bar + baz;
}
```

Lexeme 'int'	Token Category TYPE_KEYWORD
'main'	IDENTIFIER
"("	LEFT_PAREN
')'	RIGHT_PAREN
<b>'</b> { <b>'</b>	LEFT_CURLY
'float'	TYPE_KEYWORD
'foo'	IDENTIFIER
6 J J	СОММА
'bar'	IDENTIFIER
4 J J	СОММА
'baz'	IDENTIFIER
6 _ 9 9	SEMICOLON
'foo'	IDENTIFIER
<b>'</b> = <b>'</b>	OPERATOR
'bar'	IDENTIFIER
<b>'+'</b>	OPERATOR
'baz'	IDENTIFIER
4 <u>-</u> 9 9	SEMICOLON
<b>'}'</b>	RIGHT_CURLY

#### Example: Lexical Analysis

Token Lexeme 'int' TYPE\_KEYWORD 'main' **IDENTIFIER** '(' LEFT\_PAREN Nothing nothing ')' **RIGHT\_PAREN** "{ ' LEFT\_CURLY 'float' **TYPE KEYWORD** 'foo' **IDENTIFIER** 6 7 COMMA . 'bar' **IDENTIFIER** 6.7 **COMMA** . . 'baz' **IDENTIFIER** 6.7 **SEMICOLON** 'foo' **IDENTIFIER** 6\_6 **OPERATOR** 'bar' **IDENTIFIER '+' OPERATOR** 'baz' **IDENTIFIER** 6.7 **SEMICOLON** Ŷ **RIGHT\_CURLY** 

CFG Definition type IDENTIFIER LEFT\_PAREN optional\_arguements RIGH\_PAREN LEFT\_CURLY Statements

**RIGHT\_CURLY** 

#### Example: Lexical Analysis

If (i== j)
z = 0; /\* comment \*/
else
z= 1;

tif(i== j) nttz = 0; /\* comment\*/ntelse nttz = 1; IF, LPAR, ID("i"), EQUALS, ID("j"), RPAR, ID("z"), ASSIGN,INTLIT("0"), SEMI, ELSE, ID("z"), ASSIGN, INTLIT("1"), SEMI

#### State Diagram



### Implementing the State Diagram

```
int CharClass;
char lexeme [100];
char nextChar;
int lexLen;
int LETTER=0;
int DIGIT=1;
int UNKNOWN=-1;
/*addChar - a function to add nextChar to lexeme */
void addChar() {
   if (lexLen <=99)
      lexeme[lexLen++]=nextChar;
   else printf("Error: Lexeme is too long n'');
}
/* getNonBlank - call getChar until it returns a non-whitespace character */
void getNonBlank() {
   while(isspace(nextChar))
      getChar();
}
```

#### Implementing the State Diagram

```
/*getChar - a function to get the next character of input and determine its
   character class*/
void getChar() {
   if(isalpha(nextChar))
      charClass = LETTER;
     else if (isdigit(nextChar))
         charClass = DIGIT;
       else
          charClass = UNKNOWN
}
/*lex - a simple lexical analyzer*/
int lex() {
     lexlen =0;
     static int first=1;
     if (first) {
         getChar();
         first=0;
     }
     getNonBlank();
```

```
switch (CharClass) {
```

## Implementing the State Diagram

```
switch (charClass) {
       case LETTER:
         addChar();
         getChar();
         while (charClass == LETTER || charClass == DIGIT) {
               addChar();
               getChar();
         }
        return lookup(lexeme);
        break;
       case DIGIT:
         addChar();
         getChar();
         while (charClass == DIGIT) {
               addChar();
               getChar();
         3
        return INT LIT;
        break;
   } /*End of switch*/
} /*End of lex*/
```

#### State Diagram

 Problem: Design a state diagram to recognize one form of the comments of Cbased programming languages, those that begin with /\* and end with \*/



#### Example: state diagram

For-Keyword = for Identifier = [a-z][a-z0-9]\*



# Lexical Analyzer with Lex

- Lex is a domain-specific programming language for creating programs to process streams of input characters.
- A Lex program has the following form:

declarations %% translation rules %%

- auxiliary functions
- The declarations section can contain declarations of variables, manifest constants, and regular definitions. The declarations section can be empty.
- The translation rules are each of the form pattern {action}
  - Each pattern is a regular expression which may use regular definitions defined in the declarations section.
  - Each action is a fragment of C-code. The braces around the action may be omitted if the action is a single statement.

# Lexical Analyzer with Lex

a + b matches the string a + b. . matches any character except a newline. ^ matches the empty string at the beginning of a line. \$ matches the empty string at the end of a line. [a-z] matches any lowercase letter between a and z. [A-Za-z0-9] matches any alphanumeric character. [^abc] matches any character except an a, or a b, or a c. [^0–9] matches any nonnumeric character. a\* matches a string of zero or more a's. a+ matches a string of one or more a's. .a? matches a string of zero or one a's. a{2,5} matches any string consisting of two to five a's. a/b matches an a when followed by a b. \n matches a newline. \t matches a tab

## Example with Lex

```
int num_words = 0,num_numbers = 0, num_lines = 0;
word [A-Za-z]+
number [0–9]+
%%
{word} {++num_words;}
{number} {++num_numbers;}
n ++num_{lines;}
. { }
%%
int main() {
    yylex();
    printf("# of words = %d, # of numbers = %d, # of lines = %dn",
             num_words, num_numbers, num_lines );
}
```
# Example with Lex

- Put lex program into a file file.
- Compile the lex program with the command
  - lex file.l
    - This command produces an output file lex.yy.c.
- Compile this output file with the C compiler and the lex library -II
  - gcc lex.yy.c -ll
    - The resulting a.out is the lexical processor.

## The Parsing Problem

- Given an input program:
  - Operates on tokens and groups them into useful grammatical structures.
  - Find all syntax errors; for each, produce an appropriate diagnostic message and recover quickly
  - Produce the parse tree, or at least a trace of the parse tree, for the program

## The Parsing Problem (cont.)

- Two categories of parsers
  - Top down produce the parse tree, beginning at the root
    - Order is that of a leftmost derivation
    - Traces or builds the parse tree in preorder
  - Bottom up produce the parse tree, beginning at the leaves
    - Order is that of the reverse of a rightmost derivation

# The Parsing Problem (cont.)

#### Top-down Parsers

- Given a sentential form,  $xA\alpha$ , the parser must choose the correct A-rule to get the next sentential form in the leftmost derivation, using only the first token produced by A
- The most common top-down parsing algorithms:
  - Recursive descent a coded implementation
  - LL parsers table driven implementation

# The Parsing Problem (cont.)

#### Bottom-up parsers

- Given a right sentential form,  $\alpha$ , determine what substring of  $\alpha$  is the right-hand side of the rule in the grammar that must be reduced to produce the previous sentential form in the right derivation
- The most common bottom-up parsing algorithms are in the LR (Left-to-right, Rightmost derivation) type.

## **Recursive-Descent Parsing**

- There is a subprogram for each nonterminal in the grammar, which can parse sentences that can be generated by that nonterminal
- EBNF is ideally suited for being the basis for a recursive-descent parser, because EBNF minimizes the number of nonterminals

- A grammar for simple expressions:
- <expr> → <term> { (+ | -) <expr>}
  <term> → <factor> { (\* | /) <term>}
  <factor> → id | int\_constant | ( <expr> )

- Assume we have a lexical analyzer named lex, which puts the next token code in nextToken
- The coding process when there is only one RHS (Right Hand Side):
  - For each terminal symbol in the RHS, compare it with the next input token; if they match, continue, else there is an error
  - For each nonterminal symbol in the RHS, call its associated parsing subprogram

```
/* Function expr
   Parses strings in the language
   generated by the rule:
   \langle expr \rangle \rightarrow \langle term \rangle \{ (+ | -) \langle term \rangle \}
 */
void expr() {
/* Parse the first term */
  term();
/* As long as the next token is + or -, call
   lex to get the next token and parse the
   next term */
  while (nextToken == ADD OP ||
          nextToken == SUB OP) {
    lex();
    term();
  }
}
```

- A nonterminal that has more than one RHS requires an initial process to determine which RHS it is to parse
  - The correct RHS is chosen on the basis of the next token of input (the lookahead)
  - The next token is compared with the first token that can be generated by each RHS until a match is found
  - If no match is found, it is a syntax error

```
/* term
Parses strings in the language generated by the rule:
<term> -> <factor> { (* | /) <factor>)
*/
void term() {
 printf("Enter <term>\n");
/* Parse the first factor */
  factor();
/* As long as the next token is * or /,
   next token and parse the next factor */
  while (nextToken == MULT OP || nextToken == DIV OP) {
    lex();
   factor();
  }
  printf("Exit <term>\n");
} /* End of function term */
```

```
/* Function factor
  Parses strings in the language
  generated by the rule:
  <factor> -> id | (<expr>) */
void factor() {
/* Determine which RHS */
  if (nextToken) == ID CODE || nextToken == INT CODE)
/* For the RHS id, just call lex */
    lex();
/* If the RHS is (<expr>) - call lex to pass over the left parenthesis,
  call expr, and check for the right parenthesis */
  else if (nextToken == LP CODE) {
   lex();
     expr();
    if (nextToken == RP CODE)
        lex();
      else
      error();
    } /* End of else if (nextToken == ... */
  else error(); /* Neither RHS matches */
}
```

Trace of the recursive descent parser for the string a + b \* c

```
Call lex /* returns a */
Enter <expr>
Enter <term>
Enter <factor>
Call lex /* returns + */
Exit <factor>
Exit <term>
Call lex /* returns b */
Enter <term>
Enter <factor>
Call lex /* returns * */
Exit <factor>
Call lex /* returns c */
Enter <factor>
Call lex /* returns end-of-input */
Exit <factor>
Exit <term>
Exit <expr>
```

#### Bottom-up Parsing

- The Bottom-up Parsing problem is finding the correct RHS in a right-sentential form to reduce to get the previous rightsentential form in the derivation
- LR Parsers
  - Canonical LR (Knuth, 1965): original LR algorithm

- Advantages of LR parsers:
  - They will work for nearly all grammars that describe programming languages.
  - They work on a larger class of grammars than other bottom-up algorithms, but are as efficient as any other bottom-up parser.
  - The LR class of grammars is a superset of the class parsable by LL parsers.

#### Structure of An LR Parser

An LR configuration stores the state of an LR parser

 $(S_0X_1S_1X_2S_2...X_mS_m, a_ia_i+1...a_n$ 



- LR parsers are table driven, where the table has two components, an ACTION table and a GOTO table
  - The ACTION table specifies the action of the parser, given the parser state and the next token
    - Rows are state names; columns are terminals
  - The GOTO table specifies which state to put on top of the parse stack after a reduction action is done
    - Rows are state names; columns are nonterminals

- Initial configuration:  $(S_0, a_1...a_n)$
- Parser actions:
  - If ACTION[S<sub>m</sub>, a<sub>i</sub>] = Shift S, the next configuration is:

 $(S_0X_1S_1X_2S_2...X_mS_ma_iS, a_{i+1}...a_n$ 

- If ACTION[S<sub>m</sub>,  $a_i$ ] = Reduce A  $\rightarrow \beta$  and S = GOTO[S<sub>m-r</sub>, A], where r = the length of  $\beta$ , the next configuration is

 $(S_0X_1S_1X_2S_2...X_{m-r}S_{m-r}AS, a_ia_{i+1}...a_n$ 

- Parser actions (continued):
  - If ACTION[S<sub>m</sub>, a<sub>i</sub>] = Accept, the parse is complete and no errors were found.
  - If ACTION[ $S_m$ ,  $a_i$ ] = Error, the parser calls an error-handling routine.

# LR Parsing Table

1.  $E \rightarrow E+T$ 2.  $E \rightarrow T$ 3.  $T \rightarrow T*F$ 4.  $T \rightarrow F$ 5.  $F \rightarrow (E)$ 6.  $F \rightarrow id$ 

	Action						Goto		
State	id	+	*	(	)	\$	E	Т	F
0	S5			<b>S4</b>			1	2	3
1		S6				accept			
2		R2	S7		R2	R2			
3		R4	R4		R4	R4			
4	S5			S4			8	2	3
5		R6	R6		R6	R6			
6	\$5			S4				9	3
7	\$5			S4					10
8		S6			S11				
9		R1	S7		R1	R1			
10		R3	R3		R3	R3			
11		R5	R5		R5	R5			

## Trace of a parse

Stack	
0	
0id5	
0F3	
0T2	
0T2*7	
0T2*7(	4
0T2*7(	4id5
0T2*7(	4F3
0T2*7(	4T2
0T2*7(	4E8
0T2*7(	4E8+6
0T2*7(	4E8+6id5
0T2*7(	4E8+6F3
0T2*7(	4E8+6T9
0T2*7(	4E8
0T2*7(	4E8)11
0T2*7F	10
0T2	
0E1	

Input
id * (id + id)
* (id + id) \$
* (id + id) \$
* (id + id) \$
(id + id) \$
id + id) \$
+ id) \$
+ id)\$
+ id)\$
+ id)\$
id ) \$
)\$
)\$
)\$
)\$
\$
\$
\$
\$

Action Shift 5 Reduce 6 (Use GOTO[0, F]) Reduce 4 (Use GOTO[0, T]) Shift 7 Shift 4 Shift 5 Reduce 6 (Use GOTO[4, F]) Reduce 4 (Use GOTO[4, T]) Reduce 2 (Use GOTO[4, E]) Shift 6 Shift 5 Reduce 6 (Use GOTO[6, F]) Reduce 4 (Use GOTO[6, T]) Reduce 1 (Use GOTO[4, E]) Shift 11 Reduce 5 (Use GOTO[7, F]) Reduce 3 (Use GOTO[0, T]) Reduce 2 (Use GOTO[0, E]) ACCEPT

## Parsers Classification

- LL(k) Top down parser
  - LL(1) are simple parsers but cannot recognize all context free grammars
- LR(k) Bottom up parser
  - LR(1) are more powerful
  - LALR(1)
    - Not as powerful as full LR(1) but simpler to implement
    - Yacc uses LALR(1) parse tables to construct a restricted LR

## **Recommended Links**

- Parsing Simulator
- Compiler component generators
  - lexical analyzer generators: lex, flex
  - syntax analyzer generator: yacc, bison
  - <u>PLY (Python Lex-Yacc)</u>
  - ANTLR lexer-parser generator
  - Compact Guide to Lex & Yacc
  - <u>Flex</u>
  - GNU Bison
  - JavaCC

## Summary

- BNF and context-free grammars are equivalent metalanguages to describe syntax
- Regular expressions and Finite Automata
- Syntax analysis
  - Lexical Analysis (Produce tokens)
  - Parsing (Produces a parse tree)
- Lexical Analyzer
- Parsing
  - Top-down approach
    - Recursive-descent parser
  - Bottom-up approach
    - The LR family of shift-reduce parsers is the most common bottom-up parsing approach