## **System Programming and Compiler Construction**

VI Semester (Computer)

Academic Year: 22-23

# **Experiment No 5**

Aim: Study of Lexical analyzer tool -Flex/ Lex

**Leraning Objective:** Recognise lexical pattern from given input file

Theory:

Lex is a program generator designed for lexical processing of character input streams. It accepts a high-level, problem oriented specification for character string matching, and produces a program in a general purpose language which recognizes regular expressions. The regular expressions are specified by the user in the source specifications given to Lex. The Lex written code recognizes these expressions in an input stream and partitions the input stream into strings matching the expressions. At the boundaries between strings program sections provided by the user are executed. The Lex source file associates the regular expressions and the program fragments. As each expression appears in the input to the program written by Lex, the corresponding fragment is executed. The user supplies the additional code beyond expression matching needed to complete his tasks, possibly including code written by other generators. The program that recognizes the expressions is generated in the general purpose programming language employed for the user's program fragments. Thus, a high level expression language is provided to write the string expressions to be matched while the user's freedom to write actions is unimpaired. This avoids forcing the user who wishes to use a string manipulation language for input analysis to write processing programs in the same and often inappropriate string handling language.

Lex is not a complete language, but rather a generator representing a new language feature which can be added to different programming languages, called `host languages." Just as general purpose languages can produce code to run on different computer hardware, Lex can write code in different host languages. The host language is used for the output code generated by Lex and also for the program fragments added by the user. Compatible runtime libraries for the different host languages are also provided. This makes Lex adaptable to different environments and different users. Lex itself exists on UNIX, GCOS, and OS/370; but the code generated by Lex may be taken anywhere where appropriate compilers exist.

Lex turns the user's expressions and actions (called source in this pic) into the host general-purpose language; the generated program is named yylex. The yylex program will recognize expressions in a stream (called input in this pic) and perform the specified actions for each expression as it is detected.

## **System Programming and Compiler Construction**

Academic Year: 22-23

VI Semester (Computer)

+-----+
Source -> | Lex | -> yylex
+----+
+----+
Input -> | yylex | -> Output
+-----+
An overview of Lex

For a trivial example, consider a program to delete from the input all blanks or tabs at the ends of lines.

%%

$$[\t]+$;$$

is all that is required. The program contains a %% delimiter to mark the beginning of the rules, and one rule. This rule contains a regular expression which matches one or more instances of the characters blank or tab (written \t for visibility, in accordance with the C language convention) just prior to the end of a line. The brackets indicate the character class made of blank and tab; the + indicates ``one or more ...''; and the \$ indicates ``end of line," as in QED. No action is specified, so the program generated by Lex (yylex) will ignore these characters. Everything else will be copied. To change any remaining string of blanks or tabs to a single blank, add another rule:

%%

 $[\t]+$;$ 

The finite automaton generated for this source will scan for both rules at once, observing at the termination of the string of blanks or tabs whether or not there is a newline character, and executing the desired rule action. The first rule matches all strings of blanks or tabs at the end of lines, and the second rule all remaining strings of blanks or tabs.

# **System Programming and Compiler Construction**

Academic Year: 22-23

## VI Semester (Computer)

Lex can be used alone for simple transformations, or for analysis and statistics gathering on a lexical level. Lex can also be used with a parser generator to perform the lexical analysis phase.

The general format of Lex source is:

{definitions}

%%

{rules}

%%

{user subroutines}

where the definitions and the user subroutines are often omitted. The second %% is optional, but the first is required to mark the beginning of the rules. The absolute minimum Lex program is thus

%%

(no definitions, no rules) which translates into a program which copies the input to the output unchanged. In the outline of Lex programs shown above, the rules represent the user's control decisions; they are a table, in which the left column contains regular expressions and the right column contains actions, program fragments to be executed when the expressions are recognized. Thus an individual rule might appear integer printf("found keyword INT"); to look for the string integer in the input stream and print the message ``found keyword INT" whenever it appears. In this example the host procedural language is C and the C library function printf is used to print the string. The end of the expression is indicated by the first blank or tab character. If the action is merely a single C expression, it can just be given on the right side of the line; if it is compound, or takes more than a line, it should be enclosed in braces.

## **Implementation Details:**

- 1. Open file in text editor
- 2. Enter keywords, rules for identifier and constant, operators and relational operators. In

the following format

a) % {

# **System Programming and Compiler Construction**

Academic Year: 22-23

## VI Semester (Computer)

Definition of constant /header files

%}

b) Regular Expressions

%%

Transition rules

%%

- c) Auxiliary Procedure (main() function)
- 3. Save file with .l extension e.g. Mylex.l
- 4. Call lex tool on the terminal e.g. [root@localhost]# lex Mylex.l This lex tool will convert
- ".l" file into ".c" language code file i.e. lex.yy.c
- 5. Compile the file lex.yy.c e.g. **gcc lex.yy.c** .After compiling the file lex.yy.c, this will create the output file **a.out**
- 6. Run the file a.out e.g. ./a.out
- 7. Give input on the terminal to the **a.out** file upon processing output will be displayed

## Sample Code

```
%{
#include<stdio.h>
int key_word=0;
%}
%%
"include"|"for"|"define" {key_word++;}
%%
int main()
{
printf("enter the sentence");
yylex();
printf("keyword are: %d\n ",key_word);
}
```

# **System Programming and Compiler Construction**

```
VI Semester (Computer)
                                                                    Academic Year: 22-23
int yywrap() { return 1; }
Example: Program for counting number of vowels and consonant
%{
#include <stdio.h>
int vowels = 0;
int consonants = 0;
%}
%%
[aeiouAEIOU] vowels++;
[a-zA-Z] consonants++;
[\n];
. ;
%%
int main()
{
printf ("This Lex program counts the number of vowels and ");
printf ("consonants in given text.");
printf ("\nEnter the text and terminate it with CTRL-d.\n");
yylex();
printf ("Vowels = \%d, consonants = \%d.\n", vowels, consonants);
```

## **System Programming and Compiler Construction**

# VI Semester ( Computer) return 0; } Output: #lex alphalex.l #gcc lex.yy.c #./a.out

This Lex program counts the number of vowels and consonants in given text.

Enter the text and terminate it with CTRL-d.

Iceream

Vowels =4, consonants =3.

#### **Test Cases:**

- 1. Input integer constant
- 2. Input special symbols

## **Output:**

```
universe@acer9: 🚛
                                                                                                                                                                                      universe@aser®: ~/Des
File Actions Edit View Help
                                                                                                       File Actions Edit View Help
               universe@acer9: ~/Desktop/9207
                                                                                                                     universe@acer9: ~/Desktop/9207
oash: cd: D: No such file or directory
                                                                                                       universe@acer9:~/Desktop/9207$ lex spcc_EXP_4_2.l && cc lex.yy.c
universe@acer9:~$ cd Desktop/
universe@acer9:~<mark>/Desktop</mark>$ ls
                                                                                                       universe@acer9:~/Desktop/9207$ ./a.out
9207 9550 9575 9626 9635

9218 9558 9579 9626 a.pdf 9637

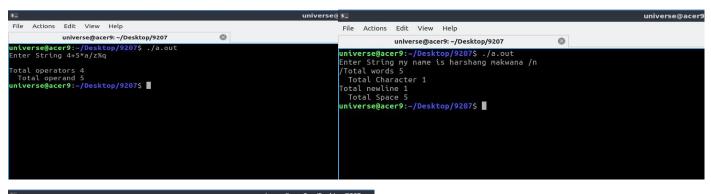
9525 9562 9613 9626.pdf ascc:

universe@acer9:-/besktop$ cd 9207
                                                    ascc2.jpg computer.desktop
block1.jpg exp5.jpg
                                                                                          mpexp2.asEnter String 345
                                     ascc1.jpg brita4.asm invert.png
                                                                                           Notepad+
                                                                                                       Total Positive 3
                                                                                                         Total Negative 0
universe@acer9:~/Desktop/9207$ ls
a.out lex.yy.c spcc_EXP_4.l
universe@acer9:~/Desktop/9207$ lex spcc_EXP_4.l && cc lex.yy.c
                                                                                                        Total floating Positive 0
                                                                                                       Total floating negative 0
universe@acer9:~/Desktop/9207$ ./a.out
                                                                                                      universe@acer9:~/Desktop/9207$
nter String vowelandConsonent
inal Answer
/owel are 6
onsonent are 11universe@acer9:~/Desktop/9207$
```

# **System Programming and Compiler Construction**

**Academic Year: 22-23** 

## VI Semester (Computer)



**Aim**: Study of Parser generator tool – Yacc

## **Leraning Objective:**

## Theory:

Parser for a grammar is a program which takes in the language string as its input and produces either a corresponding parse tree or a error. Syntax of a Language The rules which tells whether a string is a valid program or not are called the syntax Semantic's of Language The rules which give meaning to programs are called the semantic of a language Tokens When a string representing a program is broken into sequence of substrings, such that each substring represents a constant, identifier, operator, keyword etc of the language, these substrings are called the tokens of the language.

Lexical Analysis

# **System Programming and Compiler Construction**

Academic Year: 22-23

## VI Semester (Computer)

The function of lexical Analyzer is to read the input stream representing the source program, one character at a time and translate into valid tokens.

Implementation Details

## 1: Create a lex file

The general format for lex file consists of three sections:

- 1. Definitions
- 2. Rules
- 3. User subroutine Section

Definitions consists of any external 'C' definitions used in the lex actions or subroutines. The other types of definitions are definitions are lex definitions which are essentially the lex substitution strings, lex start states and lex table size declarations. The rules is the basic part which specifies the regular expressions and their corresponding actions. The user Subroutines are the functions that are used in the Lex actions.

- 2 : Yacc is the Utility which generates the function 'yyparse' which is indeed the Parser. Yacc describes a context free, LALR(1) grammer and supports both bottom up and top-down parsing. The general format for the yacc file is very similar to that of the lex file.
  - 1. Declarations
  - 2. Grammar Rules
  - 3. Subroutines

In declarations apart from the legal 'C' declarations here are few Yacc specific declarations which begins with a % sign.

## **System Programming and Compiler Construction**

Academic Year: 22-23

## VI Semester (Computer)

1. % union It defines the Stack type for the Parser.

It is union of various datas/structures/objects.

- 2. % token These are the terminals returned by the yylex function to the yacc. A token cal also have type associated with it for good type checking and syntax directed translation. A type of a token can be specified as % token <stack member> tokenName.
- 3. %type The type of non-terminal symbol in the grammar rule can be specified with this. The format is %type <stack member> non termainal.
- 4. % noassoc Specifies that there is no associativity of a terminal symbol.
- 5. % left Specifies the left associativity of a terminal symbol.
- 6. % right Specifies the right associativity of a terminal symbol.
- 7. % start specifies the L.H.S. non-terminal symbol of a production rule which specifies starting point of grammar rules.
- 8. % prac changes the precedence level associated with a particular rule to that of the following token name or literal.

The Grammar rules are specified as follows:

Context free grammar production-

p->AbC

Yacc Rule-

P: A b C { /\* 'C' actions\*/}

## **System Programming and Compiler Construction**

#### VI Semester (Computer)

The general style of coding the rules is to have all Terminals in lower –case and all non-terminals in upper –case.

Academic Year: 22-23

To facilitate a proper syntax directed translation the Yacc has something calls pseudo-variables which forms a bridge between the values of terminals/non-terminals and the actions. These pseudo variables are \$\$, \$1, \$2, \$3,......The \$\$ is the L.H.S value of the rule whereas \$1 is the first R. H. S value of the rule, so is the \$2 etc. The default type for pseudo variables is integer unless they are specified by % type. %token <type> etc.

Perform the following steps, in order, to create the desk calculator example program:

1. Process the **yacc** grammar file using the **-d** optional flag (which tells the **yacc** command to create a file that defines the tokens used in addition to the C language source code):

```
yacc -d calc.yacc
```

2. Use the **li** command to verify that the following files were created:

**y.tab.c** The C language source file that the **yacc** command created for the parser.

y.tab.h A header file containing define statements for the tokens used by the parser.

3. Process the **lex** specification file:

```
lex calc.lex
```

4. Use the **li** command to verify that the following file was created:

**lex.yy.c** The C language source file that the **lex** command created for the lexical analyzer.

5. Compile and link the two C language source files:

```
cc y.tab.c lex.yy.c
```

6. Use the **li** command to verify that the following files were created:

y.tab.o The object file for the y.tab.c source file

# **System Programming and Compiler Construction**

Academic Year: 22-23

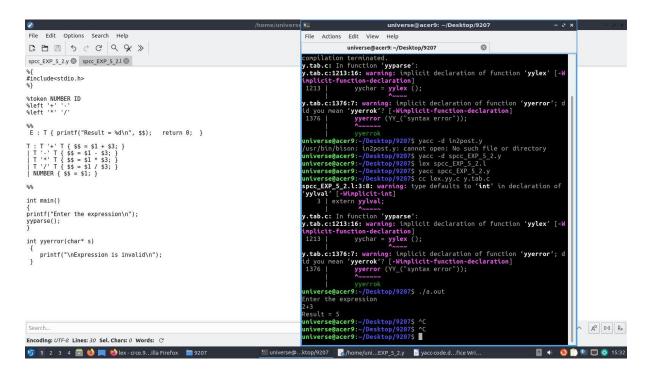
## VI Semester (Computer)

lex.yy.o The object file for the lex.yy.c source file

**a.out** The executable program file

- 7. To then run the program directly from the **a.out** file, enter:
- 8. \$ a.out

## **Output:**



#### Postlab:

- 1. Write the structure of Lex
- 2. Write the structure of Yacc