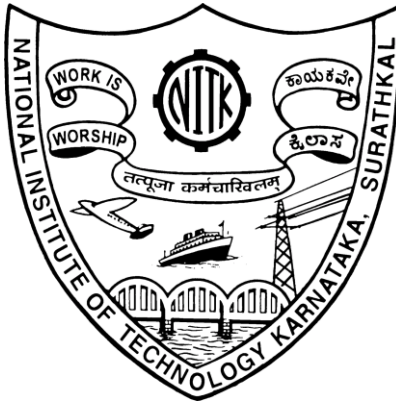


NATIONAL INSTITUTE OF TECHNOLOGY KARNATAKA,
SURATHKAL



COMPILER LAB PROJECT-1

LEXICAL ANALYZER USING FLEX

Submitted By:

Chirag Jamadagni (13CO117)
Abhijith Anilkumar (13CO102)
George CM (13CO119)

COMPILER LAB

Phase 1 - Lexical Analyzer

Abstract

A compiler, in general, is a computer program that transforms source code written in a programming language into another computer language. A compiler has many parts like the Lexical Analyzer, Parser, Semantic checker, Intermediate-code generator, Code optimizer and Code generator.

The Lexical analyzer of the compiler is the part which identifies each entry in the source program and differentiate them into lexemes of different tokens. It uses regular expression for achieving this task. Apart from identifying the token it also reports a few lexical errors.

This project makes a Lexical analyzer for a subset of C language using the flex tool. The flex tool allows users to create a scanner, with C language as the base and provides the provision for specifying the regular expression associated with each token. It then by default creates the DFA associated with it and identifies and returns the token as the output.

TABLE OF CONTENTS

Title	Page Number
Introduction	5
Design	9
Test Cases	13
Implementation	16
Results	17
Examples	18
Future Work	24
References	24

List of Tables

Title	Page No.
Test Cases	13

List of Figures

Title	Page No.
Scanner code - 1	9
Scanner code - 2	10
Scanner code - 3	10
Scanner code - 4	11
Scanner code - 5	11
Scanner code - 6	12
Run file (Shell Script)	12
Example 1 input	18
Example 1 output	18

Phase 1 – Lexical Analyzer

Example 2 input	19
Example 2 output	19
Example 3 input	20
Example 3 output	20
Example 4 input	21
Example 4 output	21
Example 5 input	22
Example 5 output	23

Introduction

This project involves building a compiler. There are four phases which leads to the development of the following parts: Lexical Analyzer, Parser, Semantic Checker and Intermediate Code Generator.

Lexical Analyzer

Lexical Analyzer is a component of a compiler that goes through the source program directly and identifies tokens.

The main task of Lexical Analyzer is to read a stream of characters as an input and produce a sequence of tokens such as names, keywords, punctuation marks etc. for syntax analyzer. It discards the white spaces and comments between the tokens and also keep track of line numbers. It reads a lexeme and identifies which token it belongs to. It also makes an entry into the symbol table

Lexical Analyzer also generates errors in the following cases:

- **Unidentified token:** When the lexeme does not match any of the specified regular expressions.
- **Unterminated String:** When the right number of inverted commas are not provided.
- **Nested Comments:** Nested comments are not supported.
- **Nested Strings:** Nestes strings are also not supported.
- **Unmatched Parenthesis:** If there are missing parenthesis, an error message is generated.

Flex Script

Flex - Fast Lexical Analyzer Generator

It is a tool for generating programs that perform pattern-matching on text.

Flex is basically a tool for generating **scanners**: programs which recognized lexical patterns in text. Flex reads the given input files, or its standard input if no file names are given, for a description of a scanner to generate. The description is in the form of pairs of regular expressions and C code, called **rules**. flex generates as output a C source file, 'lex.yy.c', which defines a routine 'yylex()'. This file is compiled to produce an executable. When the executable is run, it analyzes its input for occurrences of the regular expressions. Whenever it finds one, it executes the corresponding C code.

These programs perform character parsing and tokenizing via the use of a deterministic finite automaton (DFA). A DFA is a theoretical machine accepting regular languages. These machines are a subset of the collection of Turing machines. DFAs are equivalent to read-only right moving Turing machines. The syntax is based on the use of regular expressions. See also nondeterministic finite automaton.

Format of the Input File

The flex input file consists of three sections, separated by a line with just `%%' in it:

definitions

%%

rules

%%

user code

The **definitions** section contains declarations of simple **name** definitions to simplify the scanner specification, and declarations of **start condition**.

The **rules** section of the flex input contains a series of rules of the form:

Pattern action: Where the pattern must be unindented and the action must begin on the same line. The patterns in the input are written using an extended set of regular expressions.

Finally, the **user code** section is simply copied to `lex.yy.c' verbatim. It is used for companion routines which call or are called by the scanner. The presence of this section is optional; if it is missing, the second `%%' in the input file may be skipped, too.

Phase 1 – Lexical Analyzer

In the definitions and rules sections, any indented text or text enclosed in ``%{'` and ``%}'` is copied verbatim to the output (with the ``%{'`'s removed). The ``%{'`'s must appear unindented on lines by themselves.

In the rules section, any indented or `%{'` text appearing before the first rule may be used to declare variables which are local to the scanning routine and (after the declarations) code which is to be executed whenever the scanning routine is entered. Other indented or `%{'` text in the rule section is still copied to the output, but its meaning is not well-defined and it may well cause compile-time errors.

In the definitions section (but not in the rules section), an unindented comment (i.e., a line beginning with `/*`) is also copied verbatim to the output up to the next `*/`.

C Program

This project aims to create a compiler for subset of C language. C is a general-purpose, high-level language that was originally developed by Dennis M. Ritchie to develop the UNIX operating system at Bell Labs. C was originally first implemented on the DEC PDP-11 computer in 1972. In 1978, Brian Kernighan and Dennis Ritchie produced the first publicly available description of C, now known as the K&R standard.

C is a general-purpose, imperative computer programming language, supporting structured programming, lexical variable scope and recursion, while a static type system prevents many unintended operations. By design, C provides constructs that map efficiently to typical machine instructions, and therefore it has found lasting use in applications that had formerly been coded in assembly language, including operating systems, as well as various application software for computers ranging from supercomputers to embedded systems.

Since this compiler is built for C language, all testing has been done using C programs. Special C programs are written as test cases for testing the working of the compiler.

Design

Scanner.l

Explanation:

This code is the scanner which will be converted to lex.yy.cc. It contains regular expressions for classifying the various lexemes to tokens and also defines the actions that should be taken for the various cases.

Code:

```

1  /* Scanner Program */
2
3  /* DEFINITIONS */
4  D [0-9]
5  L [a-zA-Z_]
6  H [a-zA-Z0-9]
7  E [Ee][+-]?([0]+
8  FS (f|F|l|L)
9  IS (u|U|l|L)*
10
11 %{
12 #include <stdio.h>
13 #include <stdlib.h>
14 #include <malloc.h>
15 #include <string.h>
16 int found=0, var_cnt=0, l, nestedComment=0, bracCount=0, llineCount=0, flag=0;
17 char *comment;
18 void add_to_comment(char*);
19 void insert(char *yytext, char type);
20 struct holder
21 {
22     char *name;
23     int val;
24     struct holder *next;
25 }*st,*head;
26
27 %}
28 string L?\"([\\.|[^\"]])*\"
29 printf (printf((string)\\,([^\"])+\"))
30 prnterr (printf(.*))
31 scanf (scanf((string)\\,([^\"])+\"))
32 scanerr (scanf(.*))
33 comment (\\/\\.**)
34 constr (\\/\\.**)
35 comend (\\/\\)
36 keyword \"auto\"|\"break\"|\"case\"|\"char\"|\"const\"|\"continue\"|\"default\"|\"do\"|\"double\"|\"else\"|\"enum\"|\"extern\"|\"float\"|\"for\"|\"goto\"|\"if\"|\"int\"|\"long\"|\"register\"|\"return\"|\"short\"|
    \"signed\"|\"sizeof\"|\"static\"|\"struct\"|\"switch\"|\"typedef\"|\"union\"|\"unsigned\"|\"void\"|\"volatile\"|\"while\"
37 relop >|<|<=|>=|!=
38 ws [ \\t]+
39 %x C_COMMENT
40
41 %
42
43 /* RULES */
44 %M([a-zA-Z0-9_][relop][ws])* {insert(yytext, 'd');}
45 (printf)
46 (prnterr) {printf(\"ERROR: printferror\\n\");}
47 (scanf)

```

Figure 1

Figure 2

143,148-154	45%
-------------	-----

Figure 3

Phase 1 – Lexical Analyzer

```
144 > (if(nestedComment==0) insert(yytext,'o'));
145 ^ (if(nestedComment==0) insert(yytext,'o'));
146 | (if(nestedComment==0) insert(yytext,'o'));
147 ? (if(nestedComment==0) insert(yytext,'o'));
148
149 [ \t\v\n\F] ({}
150 . { /* ignore bad characters */ }
151
152
153
154 /* USER CODE */
155
156 int main()
157 {
158     comment = (char*)malloc(100*sizeof(char));
159     yyin=fopen("input.txt","r");
160     yyout=fopen("out.txt","w");
161     fprintf(yyout,"Symbol Table Format is:\n Lexeme\t\t\t\t\tToken\t\t\t\t\tAttribute Value\n");
162     yylex();
163     if(nestedComment!=0)
164         printf("ERROR: Comment does not end\n");
165     if(bracCount!=0)
166         printf("ERROR: Bracket mismatch\n");
167     fprintf(yyout,"\n");
168     if(flag==1)
169     {
170         lineCount=0;
171         fprintf(yyout,"\n\nComment (%d lines):\n",lineCount);
172         fprintf(yyout,"ERROR: Nested Comment");
173     }
174     else
175     {
176         fprintf(yyout,"\n\nComment (%d lines):",lineCount);
177         fputs(comment,yyout);
178     }
179     fclose(yyout);
180 }
181 int yywrap()
182 {
183     return(1);
184 }
185
186 void add_to_comment(char *yytext)
187 {
188     /* Function to display comments separately! */
189
190     int len1,len2;
191     char *temp;
```

Figure 4

```
192     len1 = strlen(comment);
193     len2 = strlen(yytext);
194     temp = (char*)malloc((len1+1)*sizeof(char));
195     strcpy(temp,comment);
196     comment = (char*)malloc((len1+len2+1)*sizeof(char));
197     strcat(temp,yytext);
198     strcpy(comment,temp);
199 }
200
201 void insert(char *yytext,char type)
202 {
203     /* Function to insert symbols to the Symbol Table */
204
205     int len1 = strlen(yytext);
206     char token[20];
207     struct holder *symbol,*temp,*nextptr;
208     nextptr = head;
209     switch(type)
210     {
211         case 'c':
212             strcpy(token,"Constant");
213             break;
214         case 'v':
215             strcpy(token,"Variable");
216             break;
217         case 'p':
218             strcpy(token,"Punctuator");
219             break;
220         case 'o':
221             strcpy(token,"Operator");
222             break;
223         case 'k':
224             strcpy(token,"Keyword");
225             break;
226         case 's':
227             strcpy(token,"String Literal");
228             break;
229         case 'd':
230             strcpy(token,"Preprocessor Statement");
231             break;
232     }
233     if(nestedComment==0)
234     {
235         for(i=0;i<var_cnt;i++,nextptr=nextptr->next)
236         {
237             symbol = nextptr;
238             if(strcmp(symbol->name,yytext)==0)
239                 break;
```

Figure 5

256,108	Bot
---------	-----

Figure 6

Run

Explanation:

This is a shell script which automates the compiling and execution of the lex code.

Code:

5 lines (4 sloc) 49 Bytes	
1	#!/bin/sh
2	lex scanner.l
3	gcc lex.yy.c -ll
4	./a.out

Figure 7

Test Cases

Test Case No.	Input Type	Input	Status
1	Preprocessor statement	#include<stdio.h> #define x 5	Passed Passed
2	Constants	float x = 5.0; double y = 10.0; int a = 2; char c = 'c';	Passed Constant = 5.0 Passed Constant = 10.0 Passed Constant = 2.0 Passed Constant = c
3	Keywords	int a = 2; float x = 5.0; return (0);	Passed Keyword = int Passed Keyword = float Passed Keyword = return
4	Punctuators	int main() { int a,b int c; } int a = 5;	Passed Punctuator = (,) Passed Punctuator = { Passed Punctuator = , Passed Punctuator = ; Passed Punctuator = } Passed Punctuator = '='

Phase 1 – Lexical Analyzer

5	Variables	int amount; float rate; char a[50] = "Chirag"; double principle;	Passed Variable = amount Passed Variable = rate Passed Variable = a Passed Variable = principle
6	String Literal	char a[50] = "Chirag"; char z[50] = "Chirag	Passed String = Chirag Passed ERROR: String doesn't end
7	Operators	a = b * c; c = a + b; d = a/c; a++; --a; d = a>b?1:2; c = !a; d = a b; d = a!:b	Passed Operator = * Passed Operator = + Passed Operator = / Passed Operator = ++ Passed Operator = -- Passed Operator = >, ?: Passed Operator = ! Passed Operator = Passed No operator
8	Comments	//Single line comment /* a	Passed Passed

Phase 1 – Lexical Analyzer

		b */ /* This is a multi /*this */ line comment*/ /* This	Passed ERROR: Nested comment Passed ERROR: Comment does not end
9	scanf() / printf() error	scanf(""); printf("");	Passed ERROR: scanferror Passed ERROR: printferror

Implementation

We have implemented a Lexical Analyzer for a subset of C language using flex tool. This tool takes regular expressions as an input and gives the corresponding token as an output. This phase involves building the Lexical Analyzer which goes through the source program, identifies the tokens and inserts them into the symbol table and constant table.

Identification of tokens:

We have written regular expressions for all tokens. The flex tool performs the longest match and classifies the lexeme into the corresponding token. We have separate regular expressions for keywords, identifiers, strings, comments and so on. For each token, the action taken is also defined separately according to its requirements. For example, when it identifies a comment or a delimiter or a preprocessor directive, it simply strips it out. However, when it is not able to classify a lexeme as a token it throws an error saying “unidentified token”.

Symbol and Constant Table:

We identify keywords, constants, variables, punctuators, operators, preprocessor statements and string literals. The symbol and constant table is stored in a linked list and displayed in the *out.txt* file. The linked list contains the token name and a value. Attribute Value is a sequential identification number given to a lexeme, i.e. that node in the linked list will contain details of that lexeme. At the bottom of the symbol table we are printing the number of comment lines and the comment itself.

Error Generation:

This lexical analyzer also generates errors if found. We have initialized variables which is incremented and decremented to maintain the count of the various brackets and comment symbols. In case of any discrepancy, an appropriate error is thrown. For unidentified token, we have defined a regular

| Phase 1 – Lexical Analyzer

expression that matches to anything that doesn't match any of the tokens defined earlier. We also handle invalid string literals, bad comments, printf() and scanf() errors.

Result

We have made a scanner which scans for a subset of C languages. The example inputs in the next section demonstrate error generation and token identification.

Examples

Example 1

This example demonstrates bracket mismatches, nested and unterminated comments.

Input

```
1 #include<stdio.h>
2
3 int main()
4 {
5     /* This has
6     /* a nested comment */
7     */
8
9     /*this is an unterminated comment
10 }
11
```

Figure 8

Output

```
abhiyith :~/Code/CCompiler/Scanner$ sh run
ERROR: Nested Comment
ERROR: Comment does not end
ERROR: Bracket mismatch
abhiyith :~/Code/CCompiler/Scanner$
```

Figure 9

Example 2

This example demonstrates printf() errors, string errors and bracket mismatches.

Input

```
1 #include<stdio.h>
2
3 int main()
4 {
5     printf("Hi
6 }
7
```

Figure 10

Output

```
abhiijith :~/Code/CCompiler/Scanners$ sh run
ERROR: String does not end
ERROR: Bracket mismatch
abhiijith :~/Code/CCompiler/Scanners$
```

Figure 11

Example 3

This example demonstrates unidentified tokens.

Input

```
1 #include<stdio.h>
2
3 int main()
4 {
5     int 7abcd = 10;
6     printf("%d",7abcd);
7 }
8
```

Figure 12

Output

```
abhiyith :~/Code/CCompiler/Scanner$ sh run
ERROR: Bad Token
abhiyith :~/Code/CCompiler/Scanner$
```

Figure 13

Example 4

This example is an extensive test case with multiple errors.

Input

```
1 #include<stdio.h>
2 /* for for for This is a multi
3  /*I am not a nested comment*/
4  line comment*/
5 int main() {
6     int amount, rate, time, si;
7     printf("\nEnter Principal Amount : ");
8     scanf("%d", &amount);
9
10    printf("\nEnter Rate of Interest : ");
11    scanf("%d", &rate);
12
13    printf("\nEnter Period of Time : ");
14    scanf("%d", &time);
15    scanf(",");
16    si = (amount * rate * time) / 100;
17    printf("\nSimple Interest : %d", si);
18    printf(",");
19    return(0);
20 }
21 }
22 /* Heyy how are you? */
23 "I am a
```

Figure 14

Output

```
abhijith :~/Code/CCompiler/Scanner$ sh run
ERROR: Nested Comment
ERROR: scanferror
ERROR: printferror
ERROR: String does not end
ERROR: Bracket mismatch
abhijith :~/Code/CCompiler/Scanner$ █
```

Figure 15

Example 5

This example has no errors. The output shows the contents of the output file. The output file contains the symbol-constant table and information about comment lines.

Input

```
1 #include<stdio.h>
2 /*This is a multi
3 line comment*/
4 int main() {
5     int amount, rate, time, si;
6     printf("\nEnter Principal Amount : ");
7     scanf("%d", &amount);
8
9     printf("\nEnter Rate of Interest : ");
10    scanf("%d", &rate);
11
12    printf("\nEnter Period of Time : ");
13    scanf("%d", &time);
14
15    si = (amount * rate * time) / 100;
16    printf("\nSimple Interest : %d", si);
17    return(0);
18 }
19
20
~
~
~
```

Figure 16

Output

```

2 Symbol Table Format is:
3 Lexeme Token Attribute Value
4
5 #include<stdio.h> Preprocessor Statement 0
6 int Keyword 1
7 main Variable 2
8 ( Punctuator 3
9 ) Punctuator 4
10 { Punctuator 5
11 int Keyword 1
12 amount Variable 6
13 , Punctuator 7
14 rate Variable 8
15 , Punctuator 7
16 time Variable 9
17 , Punctuator 7
18 si Variable 10
19 ; Punctuator 11
20 ; Punctuator 11
21 ; Punctuator 11
22 ; Punctuator 11
23 ; Punctuator 11
24 ; Punctuator 11
25 ; Punctuator 11
26 si Variable 10
27 = Punctuator 12
28 ( Punctuator 3
29 amount Variable 6
30 * Operator 13
31 rate Variable 8
32 * Operator 13
33 time Variable 9
34 ) Punctuator 4
35 / Operator 14
36 100 Constant 15
37 ; Punctuator 11
38 ; Punctuator 11
39 return Keyword 16
40 ( Punctuator 3
41 0 Constant 17
42 ) Punctuator 4
43 ; Punctuator 11
44 } Punctuator 18
45
46
47 Comment (2 lines):
48 This is a multi
49 line comment

```

Figure 17

Future Work

Future work involves making the other parts of the compiler - parser, semantic checker and the intermediate code generator.

References

- [1] Flex, version 2.5, A fast scanner generator
<https://www.cs.princeton.edu/~appel/modern/c/software/flex/flex.html>, Edition 2.5, March 1995
- [2] Flex Tutorial, <http://alumni.cs.ucr.edu/~lgao/teaching/flex.html>
- [3] Flex: A Scanner Generator for C/C++, <http://web.cse.ohio-state.edu/~gurari/course/cse756/html/cse756se8.html>
- [4] LEX/FLEX Scanner Generator, <http://cse.iitkgp.ac.in/~goutam/compiler/lect/lect3.pdf>
- [5] Introduction to FLEX, <http://web.eecs.utk.edu/~bvz/cs461/notes/flex/>
- [6] An Overview of FLEX, with examples, http://web.mit.edu/gnu/doc/html/flex_1.html