

Lexical Analyzer using Flex

COMPILER DESIGN

PROJECT - 1

19.01.2018

Submitted by

Saurabh Yadav, 15C0145

Kiran S, 15C0125

Phase 1 - Lexical Analyzer

Introduction

A compiler is a special program that processes statements written in a particular programming language and turns them into machine language or "code" that a computer's processor uses. Typically, a programmer writes language statements in a language such as C one line at a time using an *editor*. The file that is created contains what are called the *source statements*.

The programmer then runs the appropriate language compiler, specifying the name of the file that contains the source statements. When executing (running), the compiler first scans all of the language statements syntactically one after the other and then, in one or more successive stages or "passes", builds the output code, making sure that statements that refer to other statements are referred to correctly in the final code.

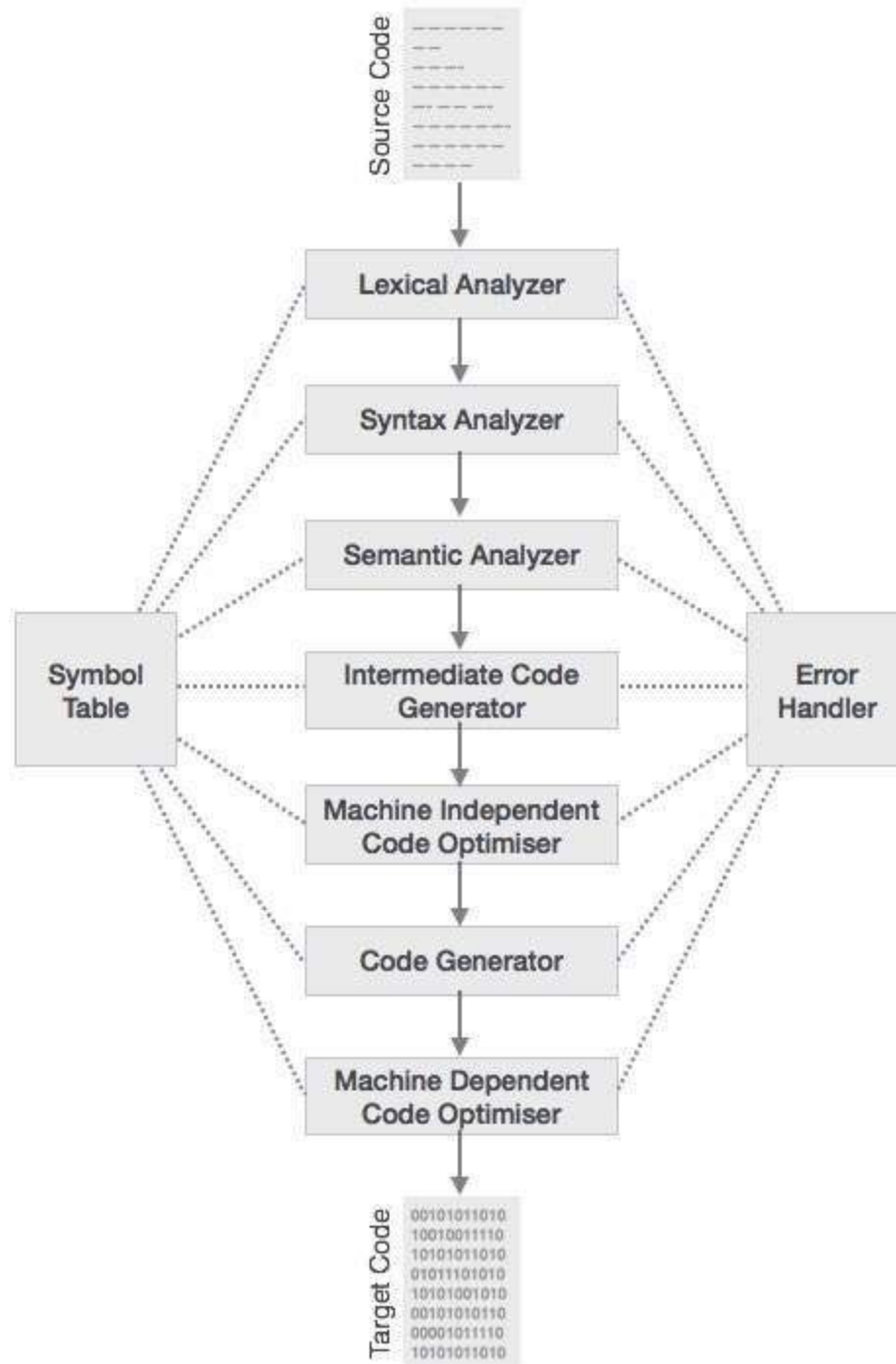
The compilation process is a sequence of various phases. Each phase takes input from its previous stage, has its own representation of source program, and feeds its output to the next phase of the compiler.

Lexical analysis : This is the initial part of reading and analysing the program text. The text is read and divided into tokens, each of which corresponds to a symbol in the programming language, e.g., a variable name, keyword or number.

Syntax analysis : This phase takes the list of tokens produced by the lexical analysis and arranges these in a tree-structure (called the syntax tree) that reflects the structure of the program. This phase is often called parsing.

Semantic Analysis : This phase analyses the syntax tree to determine if the program violates certain consistency requirements, e.g., if a variable is used but not declared or if it is used in a context that does not make sense given the type of the variable, such as trying to use a boolean value as a function pointer.

Intermediate code generation : The program is translated to a simple machine independent intermediate language.



Lexical Analyzer

Lexical analysis is the first phase of a compiler. It is the process of taking an input string of characters (such as the source code of a computer program) and producing a sequence of symbols called lexical tokens, or just tokens, which may be handled more easily by a parser. The lexical analyzer reads the source text and, thus, it may perform certain secondary tasks:

- Eliminate comments and white spaces in the form of blanks, tab and newline characters.
- Correlate errors messages from the compiler with the source program (eg, keep track of the number of lines).

The interaction with the parser is usually done by making the lexical analyzer be a sub-routine of the parser.

Token: A token is a group of characters having collective meaning: typically a word or punctuation mark, separated by a lexical analyzer and passed to a parser.

Pattern : A rule that describes the set of strings associated to a token. Expressed as a regular expression and describing how a particular token can be formed. For example, `[A-Za-z][A-Za-z_0-9]*` The pattern matches each string in the set.

Lexeme : It is an actual character sequence in the source text forming a specific instance of a token that is matched by the pattern for a token.

When more than one pattern matches a lexeme, the lexical analyzer must provide additional information about the particular lexeme. The lexical analyzer collects information about tokens into their associated attributes.

In practice, a token has usually only a single attribute - a pointer to the symbol-table entry in which the information about the token is kept such as: the lexeme, the line number on which it was first seen.

Source Code

```
%{
    #include <stdio.h>
    #include <string.h>
    #include <stdlib.h>

    int count=0,nest_cmt=0,line_count=1,comm_line;
    struct token_list
    {
        int attr_num;
        char *name,token_type[25];
        struct token_list *next;
    }*header;
    void insert(char *yytext,char type);
}%

/*=====*/

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"

letter      [a-zA-Z_]
loop        "for"|"while"
id          {letter}({letter}|{digit})*
digit       [0-9]
hexa        [a-fA-F0-9]
oct         [0-7]
exp         [Ee][+-]?{digit}+
```

```

float_suff      (f|F|l|L)
int_suff        (u|U|l|L|ul|UL|ll|LL|ull|ULL)
singlecomm      (\/\/.*)
comm_beg        (\/\/\*)
comm_end        (\/\*)
space           [ \t]+
punctuator      ";"|"{"|"}"|","|":"|"("|")"|"."
operator
">>="|"<<="|"+="|"-="|"*="|"/="|"%=|"&="|"^="|"|="|">>"|"<<"|"++"|"--"|">"|"<"|"&&"|"||"|"<="|">="|"=="|"!="|"&"|"!"|"~"|"-"|"+"|"*"|" "/"|"="|"%"|"<"|">"|"^"|"|"|"?"

%x COMM

/*=====*/

%%

/* Line Count */
<*>\n      {line_count++;}

/* Pre-processor directives */
^#([-a-zA-Z0-9.]|<|>|{space})*      {insert(yytext,'d');}

/* RegEx for Keywords */
{keyword}      { insert(yytext,'k'); }

/* RegEx for Constants Literals */
0[xX]{hexa}+{int_suff}?      {insert(yytext,'1');}
0{oct}+{int_suff}?      {insert(yytext,'2');}

```

```

{digit}+{int_suff}?                {insert(yytext,'3');}
{digit}+{float_suff}?              {insert(yytext,'4');}
{digit}+{exp}{float_suff}?         {insert(yytext,'4');}
{digit}*"."{digit}+({exp})?{float_suff}? {insert(yytext,'4');}
{digit}+"."{digit}*({exp})?{float_suff}? {insert(yytext,'4');}

/*RegEx for Identifier */
{id}                                { insert(yytext,'i'); }
{digit}+{letter}+                   { printf("Line %d: error: Invalid
Identifier\n",line_count);}

/*RegEx to identify functions and loops*/
{loop}{space}?"("                   {insert(yytext,'l');}
{keyword}{space}?"("                {insert(yytext,'k');}
{id}{space}?"("                     {insert(yytext,'f');}

/*
{id}{space}?"("(({keyword}{space}""*"?{id}?",")*{keyword}{space}""*"?{i
d}))?"")" {insert(yytext,'f');} */

/* {id}{space}?"("(("*"?{id}?",")*""*"?{id}))?"")"
{insert(yytext,'f');}*/

/*RegEx for Comments*/
{singlecomm} {}

{comm_beg} {
    BEGIN(COMM);
    nest_cmt++;
    comm_line=line_count;

```

```

    }

{comm_end}      { printf("Line %d: error: Invalid Comment
Terminator\n",line_count);}

<COMM>{comm_beg}      {
                        nest_cmt++;

                        if(nest_cmt>1)      printf("Line %d: error: Nested
Comment\n",line_count);

                        }

<COMM>{comm_end}      {
                        if(nest_cmt>0)      nest_cmt--;
                        if(nest_cmt==0)      BEGIN(INITIAL);
                        }

<COMM>. {}

/* RegEx for String Literals */
\"(\\.|[^\"])*\"      {insert(yytext,'s');}
\"(\\.|[^\"])*      {printf("Line %d: error: Incomplete String
Literal\n",line_count);}

/*RegEx for Operators*/
{operator}      {insert(yytext,'o');}

/*RegEx for Punctuators*/
{punctuator}      {}

/* RegEx to ignore unnecessary spaces */

```



```

[ \t\v\n\f]    {}
"["|"]"|"\'\\0\'" {}

.    {printf("Line %d: error: Bad Character\n",line_count);}

%%

/*=====*/

void insert(char *yytext,char type)
{
    int len1 = strlen(yytext);
    char token_type[25];
    struct token_list *lexeme,*temp,*ptr;
    ptr = header;
    switch(type)
    {
        case 'k':
            strcpy(token_type,"Keyword");
            if(yytext[strlen(yytext)-1]=='(')
yytext[strlen(yytext)-1]='\0';
            break;
        case 'l':
            strcpy(token_type,"Keyword (Loop)");
            yytext[strlen(yytext)-1]='\0';
            break;
        case 'i':
            strcpy(token_type,"Identifier (Variable)");
            break;
        case 'f':

```

```
        strcpy(token_type, "Identifier (Function)");
        yytext[strlen(yytext)-1]='\0';
        break;
    case 's':
        strcpy(token_type, "String Literal");
        break;
    case 'o':
        strcpy(token_type, "Operator");
        break;
    case 'd':
        strcpy(token_type, "Preprocessor Directive");
        break;
    case '1':
        strcpy(token_type, "Constant (Hexadecimal)");
        break;
    case '2':
        strcpy(token_type, "Constant (Octal)");
        break;
    case '3':
        strcpy(token_type, "Constant (Integer)");
        break;
    case '4':
        strcpy(token_type, "Constant (Float)");
        break;
}
if(nest_cmt==0)
{
    int i;
    for(i=0;i<count;i++,ptr=ptr->next)
```

```
{
    lexeme = ptr;
    if(strcmp(lexeme->name,yytext)==0) break;
}
if(i==count)
{
    temp = (struct token_list*)malloc(sizeof(struct
token_list));
    temp->attr_num=line_count;
    temp->name = (char*)malloc((len1+1)*sizeof(char));
    strcpy(temp->name,yytext);
    strcpy(temp->token_type,token_type);
    temp->next = NULL;
    if(count==0) header=temp;
    else lexeme->next = temp;
    count++;
    fprintf(yyout,"\n%35s %30s
%25d",temp->name,temp->token_type,temp->attr_num);
}
}
}

/* main() function */

int main()
{
    yyin=fopen("testcase.txt","r");
    yyout=fopen("output.txt","w");
```

```
fprintf(yyout, "\t\t\t\t\t-----\n");

    fprintf(yyout, "\t\t\t\t\tSYMBOL TABLE\n");

fprintf(yyout, "\t\t\t\t\t-----\n");

    fprintf(yyout, "\t\t\tLexeme\t\t\tToken\t\tAttribute Value\n");

    yylex();

    if(nest_cmt!=0)      printf("Line %d: error: Unterminated Comment\n", comm_line);

    fprintf(yyout, "\n");

    fclose(yyout);

}


int yywrap()
{

    return 1;

}
```

Test Cases

Test Case (1)

```
test1.c
1  #include <stdio.h>
2
3  int main()
4  {
5      int year;
6
7      printf("Enter a year: ");
8      scanf("%d",&year);
9
10     if(year%4 == 0)
11     {
12         if( year%100 == 0)
13         {
14             // year is divisible by 400, hence the year is a leap year
15             if ( year%400 == 0)
16                 printf("%d is a leap year.", year);
17             else
18                 printf("%d is not a leap year.", year);
19         }
20         else
21             printf("%d is a leap year.", year );
22     }
23     else
24         printf("%d is not a leap year.", year);
25     return 0;
26 }
27
```

Output

```
saurabh@saurabh-Lenovo-Flex-2-14: ~
saurabh@saurabh-Lenovo-Flex-2-14:~$ lex lexer.l
saurabh@saurabh-Lenovo-Flex-2-14:~$ cc lex.yy.c
saurabh@saurabh-Lenovo-Flex-2-14:~$ ./a.out
saurabh@saurabh-Lenovo-Flex-2-14:~$
```

Symbol Table

output1.txt			
SYMBOL TABLE			
	Lexeme	Token	Attribute Value
1			
2			
3			
4			
5			
6	#include <stdio.h>	Preprocessor Directive	1
7	int	Keyword	3
8	main	Identifier (Function)	3
9	year	Identifier (Variable)	5
10	printf	Identifier (Function)	7
11	"Enter a year: "	String Literal	7
12	scanf	Identifier (Function)	8
13	"%d"	String Literal	8
14	&	Operator	8
15	if	Keyword	10
16	%	Operator	10
17	4	Constant (Integer)	10
18	==	Operator	10
19	0	Constant (Integer)	10
20	100	Constant (Integer)	12
21	if	Keyword	15
22	400	Constant (Integer)	15
23	"%d is a leap year."	String Literal	16
24	else	Keyword	17
25	"%d is not a leap year."	String Literal	18
26	return	Keyword	25
27			
28			
29			

Test Case (2)

```
test2.c x
1  #include<stdio.h>
2  #define X 5
3
4  main()
5  {
6      int a;
7      /* A multi-line comment
8         which is not closed
9         sad
10
11      scanf("%d",&a);
12      printf("%d",a*X);
13  }
14
```

Output

```
saurabh@saurabh-Lenovo-Flex-2-14: ~
saurabh@saurabh-Lenovo-Flex-2-14:~$ lex lexer.l
saurabh@saurabh-Lenovo-Flex-2-14:~$ cc lex.yy.c
saurabh@saurabh-Lenovo-Flex-2-14:~$ ./a.out
Line 7: error: Unterminated Comment
saurabh@saurabh-Lenovo-Flex-2-14:~$
```

Symbol Table

output2.txt			
1			
2			
3			
4			
5			
6			
7			
8			
9			
10			
11			
12			
13			
14			
15			
16			
17			

SYMBOL TABLE			
Lexeme	Token	Attribute	Value
#include<stdio.h>	Preprocessor Directive		1
#define X 5	Preprocessor Directive		2
main	Identifier (Function)		4
int	Keyword		6
a	Identifier (Variable)		6

Test Case (3)

```
test3.c
1  #include <stdio.h>
2  int main()
3  {
4      int n, 5flag = 0;          // 5flag is an invalid identifier
5
6      printf("Enter a positive integer: ");
7      scanf("%d",&n);
8
9      int var;
10     for(var=2; var<=n/2; ++var)
11     {
12         // condition for nonprime number
13         if(n @ var==0) // @ is a bad character
14         {
15             5flag=1;
16             break;
17         }
18     }
19
20     if (5flag==0)
21         printf("%d is a prime number." ,n);
22     else
23         printf("%d is not a prime number.,n);
24
25     return 0;
26 }
27
```

Output

```
saurabh@saurabh-Lenovo-Flex-2-14: ~
saurabh@saurabh-Lenovo-Flex-2-14:~$ lex lexer.l
saurabh@saurabh-Lenovo-Flex-2-14:~$ cc lex.yy.c
saurabh@saurabh-Lenovo-Flex-2-14:~$ ./a.out
Line 4: error: Invalid Identifier
Line 13: error: Bad Character
Line 15: error: Invalid Identifier
Line 20: error: Invalid Identifier
Line 23: error: Incomplete String Literal
saurabh@saurabh-Lenovo-Flex-2-14:~$
```

Symbol Table

output3.txt			
SYMBOL TABLE			
	Lexeme	Token	Attribute Value
1	#include <stdio.h>	Preprocessor Directive	1
2	int	Keyword	2
3	main	Identifier (Function)	2
4	n	Identifier (Variable)	4
5	=	Operator	4
6	0	Constant (Integer)	4
7	printf	Identifier (Function)	6
8	"Enter a positive integer: "	String Literal	6
9	scanf	Identifier (Function)	7
10	"%d"	String Literal	7
11	&	Operator	7
12	var	Identifier (Variable)	9
13	for	Keyword (Loop)	10
14	2	Constant (Integer)	10
15	<=	Operator	10
16	/	Operator	10
17	++	Operator	10
18	if	Keyword	13
19	==	Operator	13
20	1	Constant (Integer)	15
21	break	Keyword	16
22	if	Keyword	20
23	"%d is a prime number."	String Literal	21
24	else	Keyword	22
25			
26			
27			
28			
29			
30			
31			
32			

Test Case (4)

```

1  #include<stdio.h>
2
3  /* This is a
4     /* Nested
5         multi
6         line
7     */
8     Comment
9  */
10
11 int string_length(char *pointer)
12 {
13     int c = 0;
14     while( *(pointer + c) != '\0' )
15         c++;
16
17     return c;
18 }
19
20 void reverse(char *string)
21 {
22     int length, c;
23     char 7xyz=0, *begin, *end, temp;
24
25     length = string_length(string);
26     begin = string;
27     end = string;
28

```

```

28
29     for(c = 0; c < length - 1; c++)
30         end++;
31
32     for(c = 0; c < length/2; c++)
33     {
34         temp = *end;
35         *end = *begin;
36         *begin = temp;
37         7xyz++;
38         begin++;
39         end--;
40     }
41 }
42 main()
43 {
44     char str[100];
45     int 54ull, 54f, 54s; // 54s is not a valid integer constant
46     printf("Enter a string\n");
47     scanf("%d",str);
48     reverse(str);
49     printf("Reverse of entered string is %s", str);
50
51     return 0;
52 }
53

```

Output

```
saurabh@saurabh-Lenovo-Flex-2-14: ~  
saurabh@saurabh-Lenovo-Flex-2-14:~$ lex lexer.l  
saurabh@saurabh-Lenovo-Flex-2-14:~$ cc lex.yy.c  
saurabh@saurabh-Lenovo-Flex-2-14:~$ ./a.out  
Line 4: error: Nested Comment  
Line 14: error: Bad Character  
Line 23: error: Invalid Identifier  
Line 37: error: Invalid Identifier  
Line 45: error: Invalid Identifier  
saurabh@saurabh-Lenovo-Flex-2-14:~$
```

Symbol Table

SYMBOL TABLE			
Lexeme	Token	Attribute	Value
#include<stdio.h>	Preprocessor Directive		1
int	Keyword		11
string_length	Identifier (Function)		11
char	Keyword		11
*	Operator		11
pointer	Identifier (Variable)		11
c	Identifier (Variable)		13
=	Operator		13
0	Constant (Integer)		13
while	Keyword (Loop)		14
!=	Operator		14
++	Operator		15
return	Keyword		17
void	Keyword		20
reverse	Identifier (Function)		20
string	Identifier (Variable)		20
length	Identifier (Variable)		22
begin	Identifier (Variable)		23
end	Identifier (Variable)		23
temp	Identifier (Variable)		23
for	Keyword (Loop)		29
<	Operator		29
-	Operator		29
1	Constant (Integer)		29
/	Operator		32
2	Constant (Integer)		32
--	Operator		39
main	Identifier (Function)		42
str	Identifier (Variable)		44
100	Constant (Integer)		44
54ull	Constant (Integer)		45
54f	Constant (Float)		45
printf	Identifier (Function)		46
"Enter a string\n"	String Literal		46
scanf	Identifier (Function)		47
"%d"	String Literal		47
"Reverse of entered string is %s"	String Literal		49

*****THE END*****