# COMPILER DESIGN PROJECT REPORT - 1

# SCANNER FOR C

# PROGRAMMING LANGUAGE

Submitted By :

Arvind Ramachandran   - 15CO111

Aswanth P P            - 15CO112

DATE : 20/01/2018

# 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. A compiler translates (or compiles) a program written in a high-level programming language that is suitable for human programmers into the low-level machine language that is required by computers. During this process, the compiler will also attempt to spot and report obvious programmer mistakes. When executing (running), the compiler first parses (or analyzes) 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. We basically have two phases of compilers, namely **Analysis phase** and **Synthesis phase.**

Analysis phase creates an intermediate representation from the given source code. It is also termed as front end of the compiler. Analysis phase consists of :

1. **Lexical Analysis** :
   a. The lexical phase reads the characters in the source program and groups them into a stream of tokens in which each token represents a logically cohesive sequence of characters, such as, an identifier, a keyword, a punctuation character.
   b. The character sequence forming a token is called the lexeme for the token.
2. **Syntax Analysis** :
   a. Syntax analysis imposes a hierarchical structure on the token stream. This hierarchical structure is called syntax tree.
   b. A syntax tree has an interior node is a record with a field for the operator and two fields containing pointers to the records for the left and right children.
   c. A leaf is a record with two or more fields, one to identify the token at the leaf, and the other to record information about the token.
3. **Semantic Analysis** :
   a. This phase checks the source program for semantic errors and gathers type information for the subsequent code-generation phase.
   b. It uses the hierarchical structure determined by the syntax-analysis phase to identify the operators and operands of expressions and statements.
   c. An important component of semantic analysis is type checking.

4. **Intermediate Code Generation** :
    a. The syntax and semantic analysis generate a explicit intermediate representation of the source program.
    b. The intermediate representation should have two important properties:
        i. It should be easy to produce,
        ii. And easy to translate into target program.
    c. Intermediate representation can have a variety of forms. One of the forms is: three address code; which is like the assembly language for a machine in which every location can act like a register.
    d. Three address code consists of a sequence of instructions, each of which has at most three operands.

Synthesis phase creates an equivalent target program from intermediate representation. It is the back end of the compiler. Synthesis phase consists of :
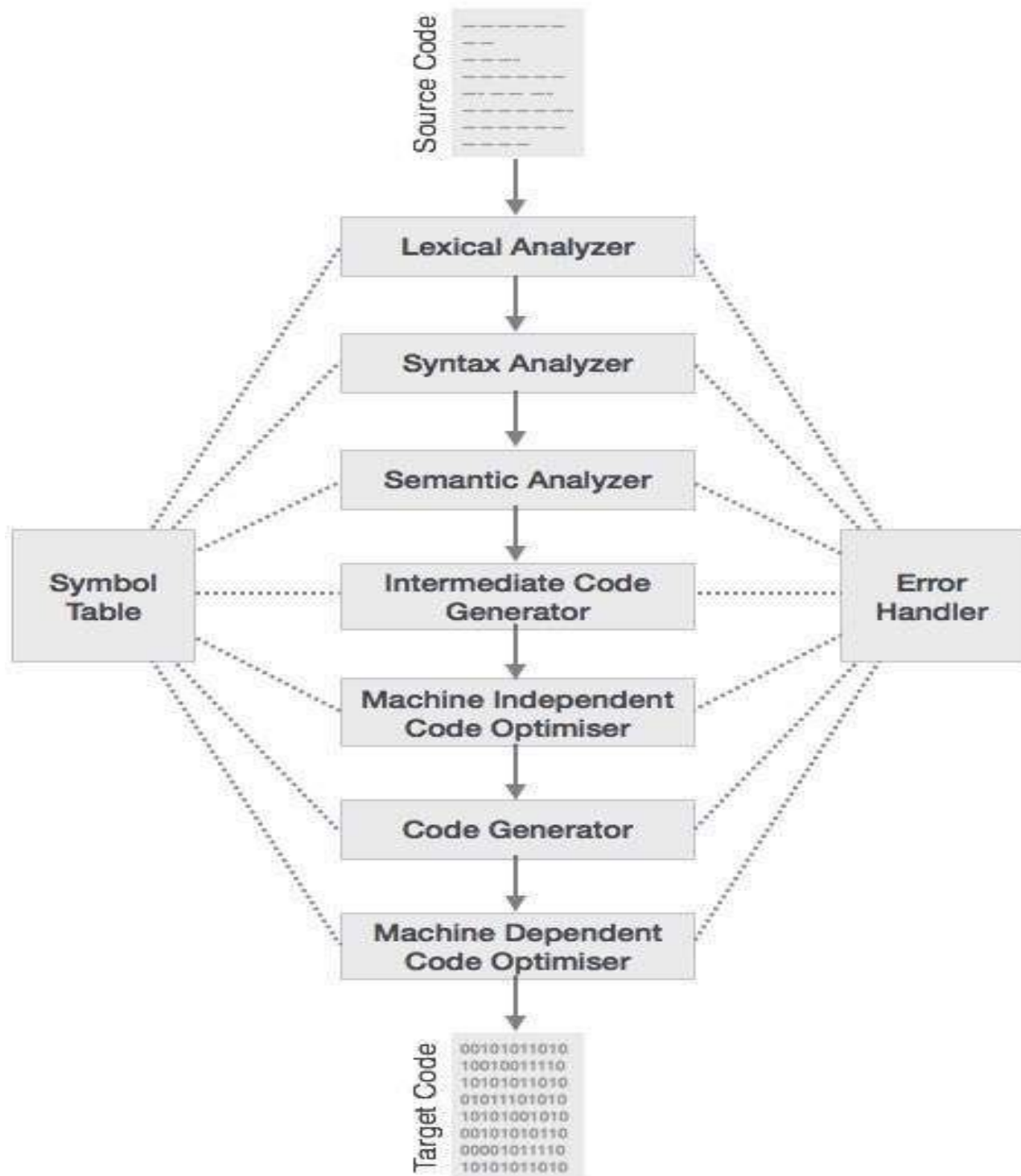
1. **Code Optimization** :
    a. Code optimization phase attempts to improve the intermediate code, so that faster-running machine code will result.
2. **Code Generation** :
    a. The final phase of the compiler is the generation of target code, consisting normally of relocatable machine code or assembly code.
    b. Memory locations are selected for each of the variables used by the program.
    c. Then, the each intermediate instruction is translated into a sequence of machine instructions that perform the same task.

The following diagram shows the various phases of a compiler.

Source Code

Lexical Analyzer

Syntax Analyzer

Semantic Analyzer

Symbol Table

Intermediate Code Generator

Error Handler

Machine Independent Code Optimiser

Code Generator

Machine Dependent Code Optimiser

Target Code

```
00101011010
10010011110
10101011010
01011101010
10101001010
00101010110
00001011110
10101011010
```

3

# LEXICAL ANALYSIS

The word "lexical" in the traditional sense means "pertaining to words". In terms of programming languages, words are objects like variable names, numbers, keywords etc. Such words are traditionally called tokens.

Lexical analysis is the first phase of compiler which is also termed as scanning. A token is a sequence of characters that represent lexical unit, which matches with the pattern, such as keywords, operators, identifiers etc. Source program is scanned to read the stream of characters and those characters are grouped to form a sequence called lexemes which produces token as output.

A lexical analyser , or lexer for short, will as its input take a string of individual letters and divide this string into tokens. Additionally, it will filter out whatever separates the tokens (the so-called white-space) , i.e., lay-out characters (spaces,newlines,etc.) and comments.

## Tokens, Patterns and Lexemes

**Token** : It is a valid sequence of characters which are given by lexeme. In a programming language, keywords, constant, identifiers, numbers, operators and punctuations symbols are possible tokens to be identified. Example of tokens:

- Type token (id, number, real, . . . )
- Punctuation tokens (IF, void, return, . . . )
- Alphabetic tokens (keywords)

**Pattern** : A pattern describes a rule that must be matched by sequence of characters (lexemes) to form a token. It can be defined by regular expressions or grammar rules.

For example - [A-Za-z][A-Za-z_0-9]*

**Lexeme** : A lexeme is a sequence of characters that matches the pattern for a token i.e., instance of a token. Eg: c=a+b*5;

The sequence of tokens produced by lexical analyzer helps the parser in analyzing the syntax of programming languages.

## Role of a Lexical Analyzer

A lexical analyzer performs the following tasks :

- Reads the source program, scans the input characters, group them into lexemes and

produce the token as output.
- Enters the identified token into the symbol table.
- Strips out white spaces and comments from source program.
- Correlates error messages with the source program i.e., displays error message with its occurrence by specifying the line number.
- Expands the macros if it is found in the source program.

## Need of Lexical Analyzer

- **Simplicity of design of compiler -** The removal of white spaces and comments enables the syntax analyzer for efficient syntactic constructs.
- **Compiler efficiency is improved -** Specialized buffering techniques for reading characters speed up the compiler process.
- **Compiler portability is enhanced**

## Lexical Errors

A character sequence that cannot be scanned into any valid token is a lexical error. Lexical errors are uncommon, but they still must be handled by a scanner. Usually, a lexical error is caused by the appearance of some illegal character, mostly at the beginning of a token.
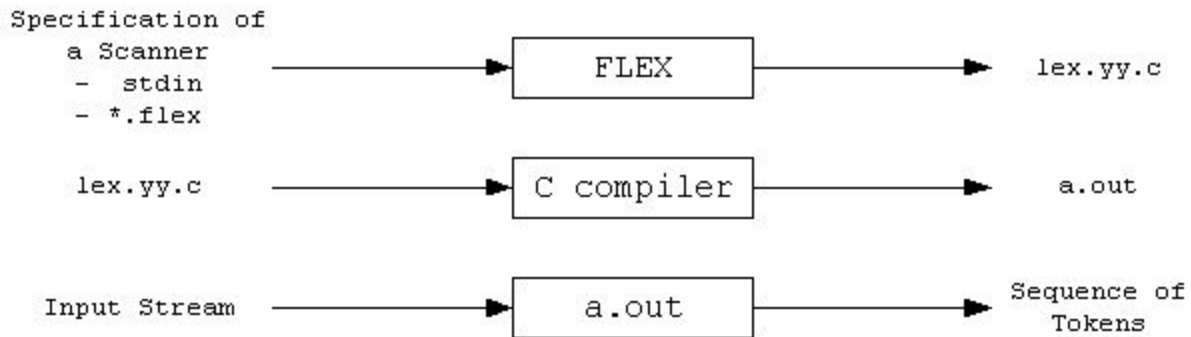
Lexical Analyzer also generates errors in the following cases:

1) **Unterminated String** : When the right number of inverted commas are not provided.
2) **Nested Comments** : Nested comments are not supported.
3) **Unmatched Parenthesis** : If there are missing parenthesis, an error message is generated.
4) **Invalid Identifier** : If the entered identifier does not match the identifier forming rules, an error message is displayed.

# FLEX

The scanner performs lexical analysis of a certain program. It reads the source program as a sequence of characters and recognizes "larger" textual units called tokens.

FLEX stands for Fast Lexical Analyzer Generator. It is a tool for generating scanners. Instead of writing a scanner from scratch, you only need to identify the vocabulary of a certain language, write a specification of patterns using regular expressions (e.g. DIGIT [0-9]), and FLEX will construct a scanner for you. FLEX is generally used in the manner depicted here:

Specification of
a Scanner
- stdin      →      FLEX      →      lex.yy.c
- *.flex

lex.yy.c      →      C compiler      →      a.out

Input Stream      →      a.out      →      Sequence of
Tokens

First, FLEX reads a specification of a scanner either from an input file *.lex, or from standard input, and it generates as output a C source file lex.yy.c. Then, lex.yy.c is compiled and linked with the "-lfl" library to produce an executable a.out. Finally, a.out analyzes its input stream and transforms it into a sequence of tokens.

- **\*.lex** is in the form of pairs of regular expressions and C code.
- **lex.yy.c** defines a routine yylex() that uses the specification to recognize tokens.
- **a.out** is actually the scanner.

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.

# FORMAT OF THE FLEX 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 conditions**, which are explained in a later section. Name definitions have the form:

*name definition*

The "name" is a word beginning with a letter or an underscore ('_') followed by zero or more letters, digits, '_', or '-' (dash). The definition is taken to begin at the first non-whitespace character following the name and continuing to the end of the line. The definition can subsequently be referred to using "{name}", which will expand to "(definition)". For example,

*DIGIT    [0-9]*
*ID      [a-z][a-z0-9]\**

defines "DIGIT" to be a regular expression which matches a single digit, and "ID" to be a regular expression which matches a letter followed by zero-or-more letters-or-digits. A subsequent reference to

*{DIGIT}+"."{DIGIT}\**

is identical to

*([0-9])+"."([0-9])\**

and matches one-or-more digits followed by a '.' followed by zero-or-more digits.


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.

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.

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 "*/".

# IMPLEMENTATION

## scanner.l

This is the lex program that contains various regular expressions for all the specific actions that are to be carried out by a lexical analyzer. This file is converted to **lex.yy.c** which is compiled to get the executable **a.out**.

CODE :



```
1   %{
2   #include <stdio.h>
3   #include <stdlib.h>
4   #include <malloc.h>
5   #include <string.h>
6   int var=0,i,nc=0,cLine=0,flag=0;
7   int lineNumber=1;
8   int cBrac=0;
9
10  FILE *symbol,*constants;
11  char *comment,*inputFile, s_comment[1000];
12
13  void insertToTable(char *yytext,char type);
14  void displayComment(char *yytext);
15  void storeSingleLineComment(char *yytext);
16
17  struct Node {
18      char *tname;
19      int av;
20      struct Node *next;
21  }*head=NULL;
22
23  %}
24
25  digit   [0-9]
26  letter  [a-zA-Z]
27  keyword "auto"|"break"|"case"|"char"|"const"|"continue"|"default"|"do"|"double"|"else"|"enum"|"extern"|"float"|"for"|"goto"|"if"|"int"|"long"|"register"|"return"
28  datatype "int"|"char"|"void"
29  sign "signed"|"unsigned"
30  modifiers "long"|"short"
31  relational >|<|<=|>=|!=|==
32  logical \&|\^|\~
33  arithmatic \+|\-|\*|\/|\%
34  puncuator \(|\)|\[|\]|\;|\,|\:|\.
35  assignment =
36  quote \'|\"|\\
37  whitespace [ \t]+
38  newline "\n"
39  singlelinecomment (\/\/.*)
40  multilinecommentstart (\/\*)
41  multilinecommentend (\*\/)
```

Line 1, Column 1    Tab Size: 4    Lex/Flex

**FIG. 1**

8

```
42  identifier ({letter}({letter}|{digit})*)|"_"({letter}|{digit})+
43
44  %x DETECT_COMMENT
45
46  %%
47
48  ^#([-a-zA-Z0-9.]|{relational}|{whitespace})* insertToTable(yytext,'d'); //preprocessor directive rule
49
50  {keyword} insertToTable(yytext,'k');
51  {sign}?{whitespace}{modifiers}?{whitespace}{datatype} insertToTable(yytext,'k'); //keyword rule
52
53  ^{sign}?{whitespace}*{modifiers}?{whitespace}*{datatype}{whitespace}*{identifier}\(.*\)  insertToTable(yytext,'j'); //procedure rule
54  {identifier}\[{digit}*\] insertToTable(yytext,'a'); // array rule
55  \*{identifier} insertToTable(yytext,'q'); // pointer rule
56  {identifier} insertToTable(yytext,'i'); // variable rule
57  {digit}+({letter}|{digit})+|"_" { printf("%s : %d : Invalid Identifier\n",inputFile,lineNumber); } // invalid identifier
58
59
60  {relational} insertToTable(yytext,'r'); //operator rules
61  {logical} insertToTable(yytext,'l');
62  {arithmatic} insertToTable(yytext,'o');
63  {assignment} insertToTable(yytext,'e');
64  {puncuator} insertToTable(yytext,'p');
65
66  {digit}+ insertToTable(yytext,'c'); //integer constants rule
67  \"(.)*\" insertToTable(yytext,'s'); //string constants rule
68  L?\"(\\.|[^\\"])*  {
69                      if(nc<=0) //invalid String
70                      printf("%s : %d : String does not End\n",inputFile,lineNumber);
71                  }
72  [-+]?{digit}*\.?{digit}+([eE][-+]?{digit}+)? insertToTable(yytext,'f');  // float constant rule
73  \'{letter}\' insertToTable(yytext,'z');  // character constant rule
74
75  {quote} ;
76  {whitespace} ;
77  {newline} lineNumber++;
78
79  "{"  {  cBrac++;
80      insertToTable(yytext,'p');
81      }
82  _
```

**FIG. 2**

```
81      }
82
83  "}"  {  cBrac--;
84      insertToTable(yytext,'p');
85      }
86
87  {singlelinecomment} {storeSingleLineComment(yytext);}
88
89  {multilinecommentstart}  {
90              BEGIN(DETECT_COMMENT);
91              nc++;
92              cLine++;
93              displayComment("\n\t");
94          }
95
96  <DETECT_COMMENT>{multilinecommentstart} {
97                  nc++;
98                  if(nc>1)
99                  {
100                     printf("%s : %d : Nested Comment\n",inputFile,lineNumber);
101                     flag = 1;
102                 }
103             }
104
105 <DETECT_COMMENT>{multilinecommentend} {
106                 if(nc>0)
107                     nc--;
108                 else
109                     printf("%s : %d : */ found before /*\n",inputFile,lineNumber);
110
111                 if(nc==0)
112                     BEGIN(INITIAL);
113             }
114
115 <DETECT_COMMENT>\n   {
116             cLine++;
117             lineNumber++;
118             displayComment("\n");
119             }
120
121 <DETECT_COMMENT>. {displayComment(yytext);}
```

**FIG. 3**

9

```
120
121   <DETECT_COMMENT>. {displayComment(yytext);}
122
123
124   %%
125
126   int main(int argc,char **argv)
127   {
128       comment = (char*)malloc(100*sizeof(char));
129       yyin=fopen(argv[1],"r");
130       inputFile=argv[1];
131
132       symbol=fopen("symbolTable.txt","w"); //File to write symbol table
133       fprintf(symbol,"\n Symbol Table:\n \t\tLexeme\t\t\tType\t\t\tAttribute Value\t\t\tLine Number\n");
134
135       constants=fopen("constantTable.txt","w"); // File to write constant table
136       fprintf(constants,"\n Constants Table:\n \t\tLexeme\t\t\tType\t\t\tAttribute Value\t\t\tLine Number\n");
137
138       yyout=fopen("parsedTable.txt","w"); // File to write all token in source program
139       fprintf(yyout,"\n Table:\n \t\tLexeme\t\t\tToken\t\t\t\tAttribute Value\t\t\t\tLine Number\n");
140
141       yylex();
142
143       if(nc!=0)
144           printf("%s : %d : Comment Does Not End\n",inputFile,lineNumber);
145
146       if(cBrac!=0)
147           printf("%s : %d : Unbalanced Parenthesis\n",inputFile,lineNumber);
148
149       fprintf(yyout,"\n");
150       if(flag==1)
151       {
152           cLine = 0;
153           fprintf(yyout,"\n\nComment (%d lines):\n",cLine);
154           printf("%s : %d : Nested Comment\n",inputFile,lineNumber);
155       }
156       else
157       {
158           int i;
159           fprintf(yyout,"\n\nMultiLineComment (%d lines):",cLine);
160           fputs(comment,yyout);
```

FIG. 4

```
159           fprintf(yyout,"\n\nMultiLineComment (%d lines):",cLine);
160           fputs(comment,yyout);
161           fprintf(yyout,"\n\nSingleLineComment :\n");
162           fputs(s_comment,yyout);
163       }
164
165       fclose(yyout);
166       fclose(symbol);
167       fclose(constants);
168   }
169
170   void storeSingleLineComment(char *yytext)
171   {
172       int len = strlen(yytext);
173       int i, j=0;
174       char *temp;
175       temp = (char*)malloc((len+1)*sizeof(char));
176       for(i=2;yytext[i]!='\0';i++)
177       {
178           temp[j++] = yytext[i];
179       }
180       strcat(temp,"\n");
181       strcat(s_comment,temp);
182   }
183   void displayComment(char *yytext)
184   {
185       int l1, l2;
186       char *temp;
187
188       l1 = strlen(comment);
189       l2 = strlen(yytext);
190       temp = (char*)malloc((l1+1)*sizeof(char));
191       strcpy(temp,comment);
192       comment = (char*)malloc((l1+l2+1)*sizeof(char));
193       strcat(temp,yytext);
194       strcpy(comment,temp);
195   }
196   void insertToTable(char *yytext,char type)
197   {
198       int l1 = strlen(yytext), i;
199
```

FIG. 5

```
198        int l1 = strlen(yytext), i;
199
200        char token[30];
201        struct Node *current = NULL, *temp = NULL;
202
203        switch(type)
204        {
205            case 'd': strcpy(token,"Preprocessor Statement");break;
206
207            case 'k': strcpy(token,"Keyword");break;
208
209            case 'j': strcpy(token,"Procedure");break;
210
211            case 'a': strcpy(token,"Array");break;
212
213            case 'q' : strcpy(token,"Pointer");break;
214
215            case 'i': strcpy(token,"Identifier");break;
216
217            case 'r': strcpy(token,"Relational Op");break;
218
219            case 'p': strcpy(token,"Punctuator");break;
220
221            case 'o': strcpy(token,"Arithmetic Op");break;
222
223            case 'c': strcpy(token,"Integer Constant");break;
224
225            case 'f': strcpy(token,"Float Constant");break;
226
227            case 'z': strcpy(token,"Character Constant");break;
228
229            case 'e': strcpy(token,"Assignment Op");break;
230
231            case 'l': strcpy(token,"Logical Op");break;
232
233            case 's': strcpy(token,"String Literal");break;
234        }
235
236        if(nc<=0)
237        {
238            current = head;
```

Line 81, Column 7                                    Tab Size: 4        Lex/Flex

FIG. 6

```
237        {
238            current = head;
239            for(i=0;i<var;i++)
240            {
241                if(strcmp(current->tname,yytext)==0)
242                {
243                    break;
244                }
245                current = current->next;
246            }
247
248            if(i==var)
249            {
250                temp = (struct Node *)malloc(sizeof(struct Node));
251                temp->av = i;
252                temp->tname = (char *)malloc(sizeof(char)*(l1+1));
253                strcpy(temp->tname,yytext);
254                temp->next = NULL;
255
256                if(head==NULL)
257                {
258                    head = temp;
259                }
260                else
261                {
262                    current = head;
263                    while(current->next!=NULL)
264                    {
265                        current = current->next;
266                    }
267                    current->next = temp;
268                }
269
270                var++;
271            }
272        }
273
274        if(type =='i' || type == 'a' || type == 'q' || type=='j')
275        {
276            fprintf(symbol,"\n%20s%30s%30d%35d",yytext,token,i,lineNumber);
277        }
```

Line 81, Column 7                                    Tab Size: 4        Lex/Flex

FIG. 7

11

```
267                current->next = temp;
268            }
269
270            var++;
271        }
272    }
273
274    if(type =='i' || type == 'a' || type == 'q' || type=='j')
275    {
276        fprintf(symbol,"\n%20s%30s%30d%35d",yytext,token,i,lineNumber);
277    }
278    switch(type)
279    {
280        case 'c' : fprintf(constants,"\n%20s%20s%30d%35d",yytext,"int",i,lineNumber);
281                   break;
282
283        case 'f' : fprintf(constants,"\n%20s%20s%30d%35d",yytext,"float",i,lineNumber);
284                   break;
285
286        case 'z' : fprintf(constants,"\n%20s%20s%30d%35d",yytext,"char",i,lineNumber);
287                   break;
288    }
289
290    fprintf(yyout,"\n%20s%30s%30d%35d",yytext,token,i,lineNumber);
291 }
292
293 int yywrap()
294 {
295    return(1);
296 }
```

Line 296, Column 2                                    Tab Size: 4        Lex/Flex

## FIG. 8

# Input.c

This is the C source program for which the lexical analysis is done. Based on this, we generate the Symbol and Constants Table.

CODE :

```
1   #include <stdio.h>
2   int main()
3   {
4       /* My name is Arvind.
5       This is
6       a sample
7       multi-line
8       comment*/
9       struct node{
10      int a;
11      char name;
12      };
13
14      float fl = 5.01;
15      char letter = 'a';
16      int _a=0,b=0,c=5;
17      unsigned int arr[50];
18      char *ptr;
19      scanf("%d %d",&a,&b);
20      int sum=0;
21      sum=a+b;
22      printf("\n Sum : %d \n",sum);
23
24      //Enter Name
25      char *name = "Arvind";
26
27      //End the main function
28      return 1;
29  }
30  void abc()
31  {
32      printf("\nTest function ");
33  }
34
```

# EXECUTION OF CODE :

The lex code can be executed by the following commands :

- lex <filename1>
- cc lex.yy.c
- ./a.out <filename2>

Where <filename1> is the lex program (here **scanner.l**)

<filename2> is the C source program for which Lexical Analysis is done. (here **input.c**)

# OUTPUT :

symbolTable.txt :

```
Symbol Table:
    Lexeme        Type        Attribute Value    Line Number

        int main()          Procedure              1                    2
            node            Identifier             4                    9
               a            Identifier             6                   10
            name            Identifier             9                   11
              fl            Identifier            12                   14
          letter            Identifier            15                   15
              _a            Identifier            17                   16
               b            Identifier            20                   16
               c            Identifier            21                   16
         arr[50]               Array              24                   17
            *ptr             Pointer              25                   18
           scanf            Identifier            26                   19
               a            Identifier             6                   19
               b            Identifier            20                   19
             sum            Identifier            31                   20
             sum            Identifier            31                   21
               a            Identifier             6                   21
               b            Identifier            20                   21
          printf            Identifier            33                   22
             sum            Identifier            31                   22
           *name             Pointer              35                   25
        void abc()          Procedure             39                   30
          printf            Identifier            33                   32
```

13

## constantTable.txt

| Lexeme | Type | Attribute Value | Line Number |
|---|---|---|---|
| 5.01 | float | 14 | 14 |
| 'a' | char | 16 | 15 |
| 0 | int | 18 | 16 |
| 0 | int | 18 | 16 |
| 5 | int | 22 | 16 |
| 0 | int | 18 | 20 |
| 1 | int | 38 | 28 |

Constants Table:

## parsedTable.txt

Table:

| Lexeme | Token | Attribute Value | Line Number |
|---|---|---|---|
| #include <stdio.h> | Preprocessor Statement | 0 | 1 |
| int main() | Procedure | 1 | 2 |
| { | Punctuator | 2 | 3 |
| struct | Keyword | 3 | 9 |
| node | Identifier | 4 | 9 |
| { | Punctuator | 2 | 9 |
| int | Keyword | 5 | 10 |
| a | Identifier | 6 | 10 |
| ; | Punctuator | 7 | 10 |
| char | Keyword | 8 | 11 |
| name | Identifier | 9 | 11 |
| ; | Punctuator | 7 | 11 |
| } | Punctuator | 10 | 12 |
| ; | Punctuator | 7 | 12 |
| float | Keyword | 11 | 14 |
| fl | Identifier | 12 | 14 |
| = | Assignment Op | 13 | 14 |
| 5.01 | Float Constant | 14 | 14 |
| ; | Punctuator | 7 | 14 |
| char | Keyword | 8 | 15 |
| letter | Identifier | 15 | 15 |
| = | Assignment Op | 13 | 15 |
| 'a' | Character Constant | 16 | 15 |
| ; | Punctuator | 7 | 15 |
| int | Keyword | 5 | 16 |
| _a | Identifier | 17 | 16 |
| = | Assignment Op | 13 | 16 |
| 0 | Integer Constant | 18 | 16 |
| , | Punctuator | 19 | 16 |
| b | Identifier | 20 | 16 |
| = | Assignment Op | 13 | 16 |
| 0 | Integer Constant | 18 | 16 |
| , | Punctuator | 19 | 16 |
| c | Identifier | 21 | 16 |
| = | Assignment Op | 13 | 16 |
| 5 | Integer Constant | 22 | 16 |

Line 1, Column 1     Spaces: 2     Plain Text

**FIG. 1**

| 37 | , | Punctuator | 19 | 16 |
|---|---|---|---|---|
| 38 | c | Identifier | 21 | 16 |
| 39 | = | Assignment Op | 13 | 16 |
| 40 | 5 | Integer Constant | 22 | 16 |
| 41 | ; | Punctuator | 7 | 16 |
| 42 | unsigned | Keyword | 23 | 17 |
| 43 | int | Keyword | 5 | 17 |
| 44 | arr[50] | Array | 24 | 17 |
| 45 | ; | Punctuator | 7 | 17 |
| 46 | char | Keyword | 8 | 18 |
| 47 | *ptr | Pointer | 25 | 18 |
| 48 | ; | Punctuator | 7 | 18 |
| 49 | scanf | Identifier | 26 | 19 |
| 50 | ( | Punctuator | 27 | 19 |
| 51 | "%d %d" | String Literal | 28 | 19 |
| 52 | , | Punctuator | 19 | 19 |
| 53 | & | Logical Op | 29 | 19 |
| 54 | a | Identifier | 6 | 19 |
| 55 | , | Punctuator | 19 | 19 |
| 56 | & | Logical Op | 29 | 19 |
| 57 | b | Identifier | 20 | 19 |
| 58 | ) | Punctuator | 30 | 19 |
| 59 | ; | Punctuator | 7 | 19 |
| 60 | int | Keyword | 5 | 20 |
| 61 | sum | Identifier | 31 | 20 |
| 62 | = | Assignment Op | 13 | 20 |
| 63 | 0 | Integer Constant | 18 | 20 |
| 64 | ; | Punctuator | 7 | 20 |
| 65 | sum | Identifier | 31 | 21 |
| 66 | = | Assignment Op | 13 | 21 |
| 67 | a | Identifier | 6 | 21 |
| 68 | + | Arithmetic Op | 32 | 21 |
| 69 | b | Identifier | 20 | 21 |
| 70 | ; | Punctuator | 7 | 21 |
| 71 | printf | Identifier | 33 | 22 |
| 72 | ( | Punctuator | 27 | 22 |
| 73 | "\n Sum : %d \n" | String Literal | 34 | 22 |
| 74 | , | Punctuator | 19 | 22 |
| 75 | sum | Identifier | 31 | 22 |
| 76 | ) | Punctuator | 30 | 22 |

**FIG. 2**

| 70 | ; | Punctuator | 7 | 21 |
|---|---|---|---|---|
| 71 | printf | Identifier | 33 | 22 |
| 72 | ( | Punctuator | 27 | 22 |
| 73 | "\n Sum : %d \n" | String Literal | 34 | 22 |
| 74 | , | Punctuator | 19 | 22 |
| 75 | sum | Identifier | 31 | 22 |
| 76 | ) | Punctuator | 30 | 22 |
| 77 | ; | Punctuator | 7 | 22 |
| 78 | char | Keyword | 8 | 25 |
| 79 | *name | Pointer | 35 | 25 |
| 80 | = | Assignment Op | 13 | 25 |
| 81 | "Arvind" | String Literal | 36 | 25 |
| 82 | ; | Punctuator | 7 | 25 |
| 83 | return | Keyword | 37 | 28 |
| 84 | 1 | Integer Constant | 38 | 28 |
| 85 | ; | Punctuator | 7 | 28 |
| 86 | } | Punctuator | 10 | 29 |
| 87 | void abc() | Procedure | 39 | 30 |
| 88 | { | Punctuator | 2 | 31 |
| 89 | printf | Identifier | 33 | 32 |
| 90 | ( | Punctuator | 27 | 32 |
| 91 | "\nTest function " | String Literal | 40 | 32 |
| 92 | ) | Punctuator | 30 | 32 |
| 93 | ; | Punctuator | 7 | 32 |
| 94 | } | Punctuator | 10 | 33 |

```
 95
 96
 97  MultiLineComment (5 lines):
 98    My name is Arvind.
 99    This is
100    a sample
101    multi-line
102    comment
103
104  SingleLineComment :
105  Enter Name
106  End the main function
107
```

**FIG. 3**

15

# TEST CASES

| Test Case Filename | Test Case Type | Code | Status |
|---|---|---|---|
| Case1.c | PreProcessor statements | #include<stdio.h><br>#define count 10 | Passed |
| Case2.c | Constants<br>Keywords<br>Punctuators<br>Variables | int a=25;<br><br>char c='h';<br><br>char arr1[10]="hello";<br><br>char arr2[10]="hello; | Passed<br><br>Constants :<br><br>Integer : 25<br><br>Character :'h'<br><br>String :"Hello"<br><br>KeyWords :<br> int char<br><br>Identifiers:<br> a  c arr1 arr2<br><br><span style="color:red">Invalid String</span><br>"hello |
| Case3.c | Relational Operators<br><br>Logical Operators<br><br>Arithmetic Operators | int a=1,b=2,c=3,d;<br>a++;<br>--b;<br>c=a+b;<br>d+=a;<br>d=a\|\|b;<br>c=a%b;<br>d=a++c;<br>a=c+/c;<br>d=a>b; | Passed<br><br>Relational :<br>><br>Logical :<br>\|\|<br><br>Arithmetic :<br>+ - / % \ |
| Case4.c | Single Line Comments<br><br>Multi Line Comments | // single line comment<br><br>/// this is valid comment<br><br>/* this is<br>a multi line<br>comment */ | Passed<br><br>Passed<br><br>Passed |

| | | | |
|---|---|---|---|
| | | /* this is<br>a /* nested */<br>comment */<br><br>invalid comment */<br><br><br>/* invalid comment | Passed<br><span style="color:red">Error</span><br><br>Passed<br><span style="color:red">Error</span><br><br>Passed<br><span style="color:red">Error</span> |
| Case5.c | Control Statements<br>Looping Statements | int a=1,b=2,c=10;<br><br>if(a>b){<br>printf("\nInside if");<br>    }<br>else{<br>printf("\nInside else");<br>}<br><br>if(b>a){<br>    if(b<c){<br>printf("\nNested if ");<br>      }<br>    }<br><br>for(int i=1;i<=10;i++){<br><br>printf("\nIteration %d",i);<br>    }<br><br>for(int i=0;i<10;i++){<br>   for(int J=0;J<20;J++){<br>printf("\nNestedloop ");<br>    }<br>} | Passed<br><br>Keywords<br>If else for<br><br>Datatype<br>Int<br><br>Identifiers<br> J a b c i<br><br>Procedures :<br><br>printf () |
| Case6.c | Function Declaration | struct student{<br>    int a;<br>    char c;<br>};<br>union teacher{<br>    int q;<br>};<br>void abc(){<br>printf("\nHello World");<br>}<br>void main(){<br>    student *S;<br>    teacher *T; | Passed<br><br>Procedures :<br><br>main()<br>abc()<br>printf()<br><br>Identifiers :<br><br>Student<br>teacher |

| | | abc();<br>} | |
|---|---|---|---|

# RESULTS

## Case1.c

**Purpose :** Test case to check Preprocessor Directives declarations

**Source Code :**

```
1   // Test case to check Pre Processor Directives declarations
2
3   #include <stdio.h>
4   #include "userHeader.h"
5   #define size 10
6
7   #include<<math.h>
8
9   int main(){
10      int a=5;
11      char c='H';
12      char str[]="Hello World";
13      return 1;
14
15  }
```

**Output:**

```
aswanth@hp-notebook: ~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ lex scanner.l
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ cc lex.yy.c
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ ./a.out Case1.c
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$
```

## Case2.c

**Purpose :** **Test case to check int and char constants,string constants,keywords, punctuators and variables**

**Source Code :**

```
1    // Test case to check int and char constants,string constants,
2    // keywords, punctuators and variables
3
4    #include<stdio.h>
5    void main(){
6        int a=25;
7        char c='h';
8        char c='abc
9        char arr1[10]="hello World";
10       char arr2[10]="hello;
11       return 1;
12
13
14   }
```

**Output :**

```
aswanth@hp-notebook: ~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases        ↑↓ En *
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ lex scanner.l
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ cc lex.yy.c
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ ./a.out Case2.c
Case2.c : 5 : Invalid Identifier
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ |
```

## Case3.c

**Purpose :** Test case to check relational, logical and arithmetic operators

**Source Code :**

```
1    // Test case to check relational, logical and arithmetic operators
2
3    #include<stdio.h>
4    void main(){
5
6        int a=1,b=2,c=3,d;
7        a++;
8        --b;
9        c=a+b;
10       d+=a;
11       d=a||b;
12       c=a%b;
13       d=a++c;
14       a=c+/c;
15
16       d=a>b;
17
18
19   }
```

19

**Output :**

```
aswanth@hp-notebook: ~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases     t↓ En *
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ lex scanner.l
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ cc lex.yy.c
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ ./a.out Case3.c
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ |
```

## Case4.c

**Purpose :** Test case to check single and multi line comments

**Source Code :**

```
1    // test case to check single and multi line comments
2
3    #include<stdio.h>
4    void main()
5    {
6        // this is a single line comment
7
8        /// this is also valid comment
9
10       /* this is
11       a multi line
12       comment */
13
14       /* this is
15       a /* nested */
16       comment */
17
18       invalid comment */
19
20       /* invalid comment
21
22
23   }
```

**Output :**

```
aswanth@hp-notebook: ~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases     ⌃ En *
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ lex scanner.l
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ cc lex.yy.c
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ ./a.out Case4.c
Case4.c : 15 : Nested Comment
Case4.c : 24 : Comment Does Not End
Case4.c : 24 : Nested Comment
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ |
```

## Case5.c

**Purpose :** Test case to check control and looping statements

**Source Code :**

```c
1   // test case to check control and looping statements
2
3   #include<stdio.h>
4   void main()
5   {
6       int a=1,b=2,c=10;
7
8       if(a>b){
9           printf("\nInside if");
10      }
11      else{
12          printf("\nInside else");
13      }
14
15      if(b>a){
16          if(b<c){
17              printf("\nNested if ");
18          }
19      }
20      for(int i=1;i<=10;i++){
21          printf("\nIteration %d",i);
22      }
23      for(int i=0;i<10;i++){
24          for(int j=0;j<20;j++){
25              printf("\nNested loop ");
26          }
27      }
28
29  }
30
```

**Output :**

```
aswanth@hp-notebook: ~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases          En
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ lex scanner.l
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ cc lex.yy.c
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ ./a.out Case5.c
Case5.c : 6 : Invalid Identifier
Case5.c : 20 : Invalid Identifier
Case5.c : 23 : Invalid Identifier
Case5.c : 24 : Invalid Identifier
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$
```

# Case6.c

**Purpose :** Test case to check function declaration ,struct and union declaration

**Source Code :**

```
1  // test case to check function decleration ,struct and union decleration
2
3  #include<stdio.h>
4  struct student{
5      int a;
6      char c;
7  };
8  union teacher{
9      int q;
10     int p;
11 }
12 void abc(){
13     printf("\nHello World");
14 }
15 void main(){
16     student *S;
17     teacher *T;
18     abc();
19 }
```

**Output :**

```
aswanth@hp-notebook: ~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases                En
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ lex scanner.l
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ cc lex.yy.c
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$ ./a.out Case7.c
aswanth@hp-notebook:~/Desktop/Sem6/CD/Compiler Design Project/Lexical Analyzer/test cases$
```

# CONCLUSION

We have made a scanner for C Programming Language. The following features have been taken care of :

1)  Identification of Keywords, Identifiers, Operators (Relational, Logical and Arithmetic), Punctuators, Constants (Integer, Character and Float) and String Literals.
2)  Arrays
3)  Pointers
4)  Single and Multi Line Comments
5)  Data Types (int, float and char) with modifiers (unsigned and signed) and types (short,long).
6)  Procedures

The following lexical errors have been handled :

1)  Nested Comments
2)  Invalid Identifiers
3)  Invalid String
4)  Balancing of Parentheses