# Lexical Analyzer for the C Language



National Institute of Technology Karnataka, Surathkal

Date: 22 August 2020

Submitted To:

**Prof. P. Santhi Thilagam**

**CSE Dept, NITK**

Group Members:

**Bhaskar Kataria, 181CO213**

**Ketan Bhujange, 181CO227**

**Omanshu Mahawar, 181CO237**

**Shrvan Warke, 181CO151**

# Abstract

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. Compiler operates in 4 phases, namely **Lexical Analyzer, Parser**, **Semantic Analysis** and **Intermediate Code generator.** This module focuses on Lexical Analyzer.

**LEXICAL ANALYSIS** is the very first phase in the compiler designing. Lexical analyzer reads the characters from source code and converts it into tokens.

Analysis of the following class of tokens and statements are made by the Lexical Analyzer:

**Keywords**
```
auto, const, default, enum, extern, register, return, sizeof, static,
struct, typedef, union, volatile, break, continue, goto, else, switch, if,
case, for, do, while, char, double, float, int, long, short, signed,
unsigned, void
```

**Comments**
Single line and multiline comments,

**Identifiers**
Identification of valid identifiers used in the language,

It supports nested for and while loops, nested if...else-if...else statements, and nested conditional statement,

**Operators**
ADD(+), MULTIPLY(*), DIVIDE(/), MODULO(%), AND(&), OR(|)

Function construct of the language, Syntax: int func(int x)

# Contents

- Introduction
    - Lexical Analyzer
    - Flex Script
    - C Program
- Design of Programs
    - Code
    - Explanation
- Test Cases
    - Without Errors
    - With Errors
- Implementation
- Results / Future work
- References

# List of Figures and Tables

# Introduction

## Lexical Analyzer

Lexical Analysis is the first phase of the compiler, also known as a scanner. It converts the High-level input program into a sequence of **Tokens**. The main task of the lexical analyzer is to read the input characters of the source program, group them into **lexemes**, and produce as output a sequence of tokens for each lexeme in the source program. The stream of tokens is sent to the parser for syntax analysis.

**Tokens:** A lexical token is a sequence of characters that can be treated as a unit in the grammar of the programming languages.
Example
- Type token (id, number, real, . . . )
- Punctuation tokens (IF, void, return, . . . )
- Alphabetic tokens (keywords)

**Lexeme:** The sequence of characters matched by a pattern to form the corresponding token or a sequence of input characters that comprises a single token is called a lexeme. E.g., "float", "abs_zero_Kelvin", "=", "-", "273", ";".
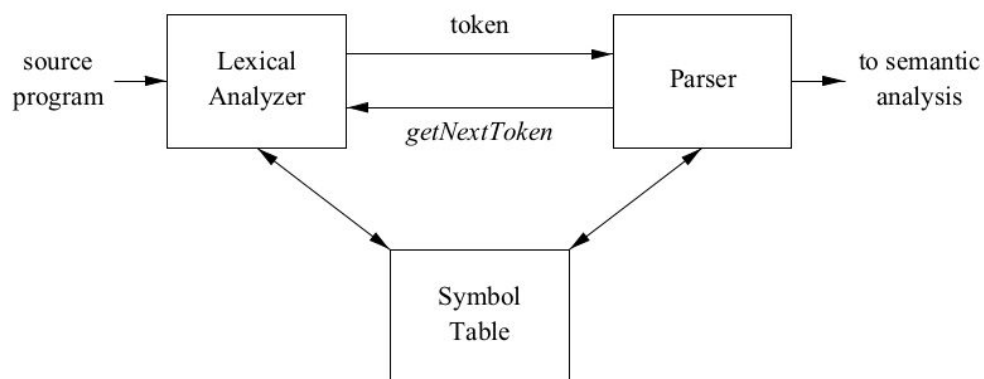
If



Figure 3.1: Interactions between the lexical analyzer and the parser

the lexical analyzer finds a token invalid, it generates an error. The lexical analyzer works closely with the syntax analyzer. It reads character streams from the source code, checks for legal tokens, and passes the data to the syntax analyzer when it demands.

# Flex Script

It is a tool/computer program for generating lexical analyzers (scanners or lexers).
Lex reads an input stream specifying the lexical analyzer and outputs source code implementing the lexer in the C programming language.

**Program Structure:**
The structure of our flex script is intentionally similar to that of a yacc file; files are divided
into three sections, separated by lines that contain only two percent signs, as follows:

```
%{
        // Definitions
%}

%%
        // Rules Section
%%
```

C code section

- The **Definition section** defines macros and imports header files written in C. It is also possible to write any C code here, which will be copied verbatim into the generated source file.
- The **Rules section** associates regular expression patterns with C statements. When the lexer sees text in the input matching a given pattern, it will execute the associated C code.
- The **C code section** contains C statements and functions that are copied verbatim to the generated source file. These statements presumably contain code called by the rules in the rules section.

## Symbol Table and Constant Table

The symbol table is an important data structure created and maintained by the compiler to keep track of semantics of variables, i.e. it stores information about the scope and binding information about names, information about instances of various entities such as variable and function names, classes, objects, etc.

# C Program

This section describes the input C program which is fed to the flex script in order to generate the lex file after taking all the rules mentioned in account. Finally, a file called lex.yy.c is generated, which when executed recognizes the tokens present in the C program which was given as an input. The script also has an option to take standard input instead of taking input from a file.

# Design of C Programs

## Code

### Symbol Table and Constant Table

```
Lexical-Analyzer > ≡ lexAnlysr.l
  1    %{
  2        #include <stdio.h>
  3        #include <string.h>
  4
  5        struct Constant_Table
  6        {
  7            char token_name[100];
  8            char token_type[100];
  9            int size;
 10        }CT[1001];
 11
 12        struct Symbol_Table
 13        {
 14            char token_name[100];
 15            char token_type[100];
 16            int size;
 17        }ST[1001];
```

### Hash Function

```
 18
 19        int find_hash(char *str)
 20        {
 21            int value = 0;
 22            for(int i = 0 ; i < strlen(str) ; i++)
 23            {
 24                value = 10*value + (str[i] - 'A');
 25                value = value % 1001;
 26                while(value < 0)
 27                    value = value + 1001;
 28            }
 29            return value;
 30        }
 31
```

## Look Up Function

```
31
32      int lookup(char *str , int lookup_Mode)
33      {
34          if(lookup_Mode == 0)
35          {
36              int value = find_hash(str);
37              if(ST[value].size == 0)
38              {
39                  return 0;
40              }
41              else if(strcmp(ST[value].token_name,str)==0)
42              {
43                  return 1;
44              }
45              else
46              {
47                  for(int i = value + 1 ; i!=value ; i = (i+1)%1001)
48                  {
49                      if(strcmp(ST[i].token_name,str)==0)
50                      {
51                          return 1;
52                      }
53                  }
54                  return 0;
55              }
56          }
57          else
58          {
59              int value = find_hash(str);
60              if(CT[value].size == 0)
61                  return 0;
62              else if(strcmp(CT[value].token_name,str)==0)
63                  return 1;
64              else
65              {
66                  for(int i = value + 1 ; i!=value ; i = (i+1)%1001)
67                  {
68                      if(strcmp(CT[i].token_name,str)==0)
69                      {
70                          return 1;
71                      }
72                  }
73                  return 0;
74              }
```

## Insert Function

```
78      void insert(char *str1, char *str2, int insert_Mode)
79      {
80          if(insert_Mode == 0)
81          {
82              if(lookup(str1, 0))
83              {
84                  return;
85              }
86              else
87              {
88                  int value = find_hash(str1);
89                  if(ST[value].size == 0)
90                  {
91                      strcpy(ST[value].token_name,str1);
92                      strcpy(ST[value].token_type,str2);
93                      ST[value].size = strlen(str1);
94                      return;
95                  }
96
97                  int pos = 0;
98
99                  for (int i = value + 1 ; i!=value ; i = (i+1)%1001)
00                  {
01                      if(ST[i].size == 0)
02                      {
03                          pos = i;
04                          break;
05                      }
06                  }
07
08                  strcpy(ST[pos].token_name,str1);
09                  strcpy(ST[pos].token_type,str2);
10                  ST[pos].size = strlen(str1);
11              }
12          }
113         else
114         {
115             if(lookup(str1, 1))
116                 return;
117             else
118             {
119                 int value = find_hash(str1);
120                 if(CT[value].size == 0)
121                 {
122                     strcpy(CT[value].token_name,str1);
123                     strcpy(CT[value].token_type,str2);
124                     CT[value].size = strlen(str1);
125                     return;
126                 }
127
128                 int pos = 0;
```

```
127
128                    int pos = 0;
129
130                    for (int i = value + 1 ; i!=value ; i = (i+1)%1001)
131                    {
132                        if(CT[i].size == 0)
133                        {
134                            pos = i;
135                            break;
136                        }
137                    }
138
139                    strcpy(CT[pos].token_name,str1);
140                    strcpy(CT[pos].token_type,str2);
141                    CT[pos].size = strlen(str1);
142                }
143            }
144        }
```

## Printing Tables

```
146        void printTables()
147        {
148            printf("\n\n-------------------------------------------------");
149            printf("\nSYMBOL TABLE\n");
150            printf("-------------------------------------------------\n");
151            for(int i = 0 ; i < 1001 ; i++)
152            {
153                if(ST[i].size == 0)
154                    continue;
155
156                printf("\t%s\t\t\t%s\n",ST[i].token_name, ST[i].token_type);
157            }
158            printf("-------------------------------------------------\n");
159            printf("\n\n-------------------------------------------------");
160            printf("\nCONSTANT TABLE\n");
161            printf("-------------------------------------------------\n");
162            for(int i = 0 ; i < 1001 ; i++)
163            {
164                if(CT[i].size == 0)
165                    continue;
166
167                printf("\t%s\t\t\t%s\n",CT[i].token_name, CT[i].token_type);
168            }
169            printf("-------------------------------------------------\n");
170        }
171    %}
172
```

## Rules

```
174    operator [[<][=]|[>][=]|[=][=]|[!][=]|[>]|[<]|[\|][\|]|[&][&]|[\!]|[=]|[\^]|[\
       +][=]|[\-][=]|[\*][=]|[\/][=]|[\%][=]|[\+][\+]|[\-][\-]|[\+]|[\-]|[\*]|[\/]|
       [\%]|[&]|[\|]|[~]|[<][<]|[>][>]]
175    floatNumber ([0-9]*)\.([0-9]+)
176    intNumber [1-9][0-9]*|0
177    postNumber [;|,|" "|\)|<|>|=|\!|\||&|\+|\-|\*|\/|\%|~|\]|\}|:|\n|\t|\^]
178    postIdentifier [" "|;|,|\(|\)|<|>|=|\!|\||&|\+|\-|\*|\/|\%|~|\n|\.|\{|\^|\t|\]]
179
180    %%
181
182    \n  {yylineno++;}
183    ([#]|" "]*(include)[ ]*([<]?)([A-Za-z]+)[.]?([A-Za-z]*)([>]?))/["\n"|\/|" "|
       "\t"] {printf("line no: %d\t%s \tPre Processor directive\n",yylineno, yytext);}
184    ([#]|" "]*(define)[" "]*([A-Za-z]+)(" ")*({floatNumber}|{intNumber}))/["\n"|\/|
       " "|"\t"] {printf("line no: %d\t%s \tMacro\n",yylineno, yytext);}
185    \/\/(.*) {printf("line no: %d\t%s \tSINGLE LINE COMMENT\n", yylineno, yytext);}

186    \/\*([^*]|[\r\n]|(\*+([^*/]|[\r\n])))*\*+\/ {printf("line no: %d\t%s \tMULTI
       LINE COMMENT\n", yylineno, yytext);}
187    [ \n\t] ;
188    ; {printf("line no: %d\t%s \tSEMICOLON\n", yylineno, yytext);}
189    , {printf("line no: %d\t%s \tCOMMA\n", yylineno, yytext);}
190    \{ {printf("line no: %d\t%s \tOPENING BRACES\n", yylineno, yytext);}
191    \} {printf("line no: %d\t%s \tCLOSING BRACES\n", yylineno, yytext);}
192    \( {printf("line no: %d\t%s \tOPENING BRACKETS\n", yylineno, yytext);}
193    \) {printf("line no: %d\t%s \tCLOSING BRACKETS\n", yylineno, yytext);}
194    \[ {printf("line no: %d\t%s \tSQUARE OPENING BRACKETS\n", yylineno, yytext);}
195    \] {printf("line no: %d\t%s \tSQUARE CLOSING BRACKETS\n", yylineno, yytext);}
196    \: {printf("line no: %d\t%s \tCOLON\n", yylineno, yytext);}
197    \\ {printf("line no: %d\t%s \tFSLASH\n", yylineno, yytext);}
198    \. {printf("line no: %d\t%s \tDOT\n", yylineno, yytext);}
199
200    auto|double|int|struct|break|else|long|switch|case|enum|register|typedef|char|
       extern|return|union|continue|for|signed|void|do|if|static|while|default|goto|
       sizeof|volatile|const|float|short|unsigned|main/[\(|" "|\{|;|:|"\n"|"\t"]
       {printf("line no: %d\t%s \tKEYWORD\n", yylineno, yytext); insert(yytext,
       "KEYWORD", 0);}
201    \"[^\n]*\"/[;|,|\)|\}] {printf("line no: %d\t%s \tSTRING CONSTANT\n",
       yylineno, yytext); insert(yytext,"STRING CONSTANT", 1);}
202    \'[A-Z|a-z]\'/[;|,|\)|:|\}] {printf("line no: %d\t%s \tCharacter CONSTANT\n",
       yylineno, yytext); insert(yytext,"Character CONSTANT", 1);}
203    [a-z|A-Z]([a-z|A-Z]|[0-9])*/\[ {printf("line no: %d\t%s \tARRAY IDENTIFIER\n",
       yylineno, yytext); insert(yytext, "IDENTIFIER", 0);}
```

```
205    {operator}/[a-z]|[0-9]|;|" "|[A-Z]|\(|\"|\'|\)|\n|\t {printf("line no: %d\t%s
       \tOPERATOR\n", yylineno, yytext);}
206
207    {intNumber}/{postNumber} {printf("line no: %d\t%s \tNUMBER CONSTANT\n",
       yylineno, yytext); insert(yytext, "NUMBER CONSTANT", 1);}
208    {floatNumber}/{postNumber} {printf("line no: %d\t%s \tFloating CONSTANT\n",
       yylineno, yytext); insert(yytext, "Floating CONSTANT", 1);}
209    [A-Za-z_][A-Za-z_0-9]*/{postIdentifier} {printf("line no: %d\t%s
       \tIDENTIFIER\n", yylineno, yytext); insert(yytext, "IDENTIFIER", 0);}
210
211
```

## Identifying Lexical Errors

```
212   (.?) {
213       printf("-------------------------------------------------\n");
214       if(yytext[0]=='#')
215       {
216           printf("ERROR at line no. %d : Error in Pre-Processor directive \n",
              yylineno);
217       }
218       else if(yytext[0]=='/')
219       {
220           printf("ERROR at line no. %d: UNMATCHED_COMMENT \n",yylineno);
221       }
222       else if(yytext[0]=='"')
223       {
224           printf("ERROR at line no. %d: INCOMPLETE_STRING \n",yylineno);
225       }
226       else
227       {
228           printf("ERROR at line no. %d\n",yylineno);
229       }
230       printf("\t%s\n", yytext);
231       printf("\n-------------------------------------------------\n");
232       return 0;
233   }
234
235   %%
```

## Main Function

```
238   int main(int argc , char **argv){
239
240       int i;
241       for (i=0;i<1001;i++)
242       {
243           ST[i].size=0;
244           CT[i].size=0;
245       }
246
247       printf("\n\n");
248
249       yyin = fopen(argv[1],"r");
250       yylex();
251
252       printTables();
253   }
254
255   int yywrap(){
256       return 1;
257   }
```

# Explanation

## Definition Section

- All header files are included in the definition section.
- It also contains the structure of the `Symbol table` and the `Constant table` along with the various functions supporting both the tables lookup, insert, and hash function.
- The `find_hash` function takes a string input and returns an integer hash of that string.
- The `lookup` function checks if entry is already present in the table or not.
- The `insert` function inserts an entry into the corresponding table, linear probing hashing technique is used to handle collisions.
- The `print_table` function is used to neatly print both the tables.

## Rules Section

- The rules section includes all the necessary rules in the form of regular expressions.
- Regular expressions are included to identify the following:
  - Preprocessor Directive
  - Keywords
  - Identifiers
  - Operators
  - Single and Multiline comments
  - Parenthesis
  - ERRORS in
    - Preprocessor directives
    - Unmatched comments
    - Incomplete string

## C Section

- In the C section, both the symbol and constant table are initialized and the `print_table` function is called to show the results
- yylex() function is called to run the program on the given input file.

# Test Cases

## Without Errors

Table 1

| Serial No. | Test Case | Expected Output | Result |
|---|---|---|---|
| 1 | #include <stdio.h> | Preprocessor directive | PASS |
| 2 | /* Comment */ | Comment | PASS |
| 3 | int n = 10; | int:      Keyword<br>n         Identifier<br>=          Operator<br>10        Number Constant | PASS |
| 4 | while() | While    Keyword<br>(          Opening Bracket<br>)          Closing Bracket | PASS |

## With Errors

Table 2

| Serial No. | Test Case | Expected Output | Result |
|---|---|---|---|
| 1 | if ( | If       Keyword<br>(         Opening Bracket<br>Error  Missing Bracket | PASS |
| 2 | /* Comment | Error  Unmatched Comment | PASS |
| 3 | string a = "ac | string:    Keyword<br>a          Identifier<br>=           Operator<br>Error      Unmatched Comment | PASS |
| 4 | # inc <stdio.> | Error in preprocessor directive | PASS |

## Identifying Comments

Figure 1



Figure 2

## Nested Loop

Figure 3

```
1    #include <stdio.h>
2    /*
3     Nested Loop no errors
4     Prints a pyramid
5    */
6
7    void main(){
8      int i = 1, n = 15;
9
10     while(i < n){
11       for (int j = 0; j < i; j++) {
12         printf("*", i);
13       }
14
15       printf("\n");
16       i++;
17     }
18   }
19
```

Figure 4

## Missing parenthesis

Figure 5



Figure 6

## Incomplete String

Figure 7



Figure 8

## Error in Preprocessor directive

Figure 9

```
 1    #incl < stdio.h>
 2
 3    //Testing String Error
 4
 5    void main()
 6    {
 7      int age = 20;
 8      char firstname = "Omanshu";
 9      char lastname = "Mahawar;
10
11          printf("%s %s\n", firstname, lastname);
12    }
13
```

Figure 10

```
omanshu>Lexical-Analyzer $./a.out < Test/9_preprocessor_err.c


------------------------------------------------------
ERROR at line no. 1 : Error in Pre-Processor directive
        #

----------------------------------------------------


----------------------------------------------------
SYMBOL TABLE
----------------------------------------------------
----------------------------------------------------


----------------------------------------------------
CONSTANT TABLE
----------------------------------------------------
----------------------------------------------------
```

# Implementation

The Regular Expressions used for each different segment of the C programming language are listed below :

- **Preprocessor directives**

  Statements processed :  #include<stdio.h>, #define
  Identified using:

  `[#][" "]*(include)[ ]*([<]?)([A-Za-z]+)[.]?([A-Za-z]*)([>]?))/["\n"|\/|" "|"\t"]`

- **Keywords:**

  Statements processed:  auto, const, default, enum, extern, register, return, sizeof, static, struct, typedef, union, volatile, break, continue, goto, else, switch, if, case, for, do, while, char, double, float, int, long, short, signed, unsigned, void and so on.

  Tokens generated: Keyword
  Identified using:

  `auto|double|int|struct|break|else|long|switch|case|enum|register|typedef|char|extern|return|union|continue|for|signed|void|do|if|static|while|default|goto|sizeof|volatile|const|float|short|unsigned|main/[\(|" "|\{|;|:|"\n"|"\t"]`

- **Identifiers:**

  Statements processed :  a, abc, a_b, a12b4
  Tokens generated: Identifier
  Identified using:

  `[a-z|A-Z]([a-z|A-Z]|[0-9])*/\`

- **Operators:**

  Statements processed :  +, -, *, /, %
  Tokens generated: Operators
  Identified using

  `operator=[[<][=]|[>][=]|[=][=]|[!][=]|[>]|[<]|[\]|[\]]|[&][&]|[\!]|[=]|[\^]|[\+][=]|[\-][=]|[\*][=]|[\/][=]|[\%][=]|[\+][\+]|[\-][\-]|[\+]|[\-]|[\*]|[\/]|[\%]|[&]|[\]]|[~]|[<][<]|[>][>]]`

  `{operator}/[a-z]|[0-9]|;|" "|[A-Z]|\(|\"|\'|\)|\n|\t {printf("\t\t\t%s\t\t\tOPERATOR\n", yytext);}`

- **Single-line comments:**

  Statements processed :  //...........
  Identified using: `\/\/(.*)`

- **Multi-line comments:**
  Statements processed :  /*............*/, /*.../*...*/ .... */
  Identified using: \/\*([^*]|[\r\n]|(\*+([^*/]|[\r\n])))*\*+\/

- **Parentheses (all types):**
  Statements processed :  (..), [..], [..]  (without errors) (..)..), [..]..], [..]..], (..., [... (with errors)

  Tokens generated: Parenthesis (without error) / Error with line number (with error)

There can be various Errors in the C program that should be displayed and some of them are handled here as follows:

- **Errors for Pre Processor directives:**
  Identified using: (yytext[o]=='#')

- **Errors for incomplete strings:**
  Statements processed:   char a[]= "abcd Error generated: Error Incomplete string and line number
  Identified using: (yytext[o]=='"')

- **Errors for nested comments:**
  Statements processed :  /*....../*....*/....
  Errors generated: Error with line number.

- **Errors for unmatched comments:**
  Statements processed :  /*.........
  Identified using: (yytext[o]=='/')

After recognizing all tokens, the lexer analyzes and prints the different identifiers, literal, and constants with the line number. The following is used for this :
- Two main structures are used to form the Symbol Table and the Constant Table which contains the identifiers and the constants.
- The lookup() function is used to check if the given identifier or constant is already present in the respective Symbol table or Constant table. The insert() function then adds the identifier or constant if it is not present.
- For every Identifier/Constant the lookup() and the insert() functions are called and then put them in their respective tables.
- At the end of the main() function, we call the printTables() function which prints the Symbol Table and the Constant Table.

# Result

The above lexical analyzer generates the following output:

```
Token Token_Type

Symbol Table
Token | Token_Type

Constant Table
Token | Token_Type
```

# Future Work

The Lexical Analyser helps to break down the source c program in tokens that are defined by the C Programming Language.
The flex script presented in this report takes care of all the rules of C language but is not fully exhaustive. Our future work would include making the script even more robust to handle all aspects of C language and making it more efficient
For the next phase, a parser will be designed, which will call the Flex program to give it tokens. The lexical analyzer cannot find the syntactical errors or find unmatched parenthesis, to do this the parser to be designed in the next phase is used.

# References

- Compilers – Principles, Techniques, and Tools By Alfred V. Aho, Monica S. Lam, Ravi Sethi, Jeffrey D. Ullman.
- https://www.geeksforgeeks.org/cc-tokens/
- http://www.isi.edu/~pedro/Teaching/CSCI565-Spring11/Practice/SDT-Sample.pdf
- StackOverflow for regex
- http://dinosaur.compilertools.net/lex/index.html