

SCANNER FOR C LANGUAGE

COMPILER DESIGN PROJECT REPORT

By

ARYAN SRIVASTAVA (RA2011033010013)

SACHIN S (RA2011033010038)

Guided by:

Dr. Sheryl Oliver A

In partial fulfilment for the Course

of

18CSC304J – COMPILER DESIGN

in CSE SPLZ IN SOFTWARE ENGINEERING



FACULTY OF ENGINEERING AND TECHNOLOGY

SRM INSTITUTE OF SCIENCE AND TECHNOLOGY

Kattankulathur, Chengalpattu District



SRM INSTITUTE OF SCIENCE AND TECHNOLOGY KATTANKULATHUR-603203

BONAFIDE CERTIFICATE

Certified that, this lab report titled **SCANNER FOR C-LANGUAGE** is the bonafide work done by:

ARYAN SRIVASTAVA (RA2011033010013)

SACHIN S (RA2011033010038)

who carried out the lab exercises under my supervision.

Certified further, that to the best of my knowledge the work reported herein does not form part of any other work.

Signature
Ms. Sheryl Oliver A
Course Faculty
Assistant Professor

Signature
Head Of Department
CINTEL

Date:5/5/2023

SCANNER FOR C LANGUAGE

ABSTRACT

A compiler is a program that can read a program in one language - the source language – and translate it to an equivalent program in another language - the target language. An important role of the compiler is to detect any errors in the source program during the translation process.

The first phase of a compiler is called **lexical analysis or scanning**. The lexical analyzer reads the stream of characters making up the source program and groups the characters into meaningful sequences called lexemes. For each lexeme, the lexical analyzer produces as output a token that it passes on to the subsequent phase, syntax analysis.

As the first phase of a compiler, 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 lexical analyzer maintains a data structure called as the symbol table. When the lexical analyzer discovers a lexeme constituting an identifier, it enters that lexeme into the symbol table. The lexical analyzer performs certain other tasks besides identification of lexemes. One such task is stripping out comments and whitespace. Another task is correlating error messages generated by the compiler with the source program.

Functionality / use of the scanner in compiler -

- Identifying and categorizing each character in the source code based on its function in the language, such as identifying keywords, identifiers, constants, operators, and punctuations.
- Building a symbol table that maintains information about each identifier, such as its name, type, and scope.
- Generating error messages if the source code contains any lexical errors, such as misspelled keywords or identifiers that are not defined.
- Optimizing the input stream by removing white spaces and comments that are not relevant to the program's functionality.

How a Scanner works-

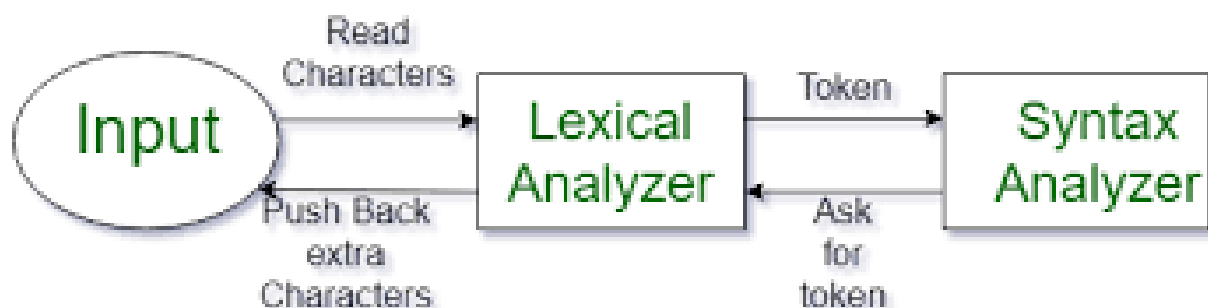
Input preprocessing involves cleaning up the input text by removing comments, whitespace, and other non-essential characters to prepare it for lexical analysis. This ensures that the input text is in a format that can be easily processed.

Tokenization involves breaking down the input text into a sequence of tokens, which are essentially the smallest units of meaning in the text. This is typically done by matching the characters in the input text against a set of patterns or regular expressions that define the different types of tokens.

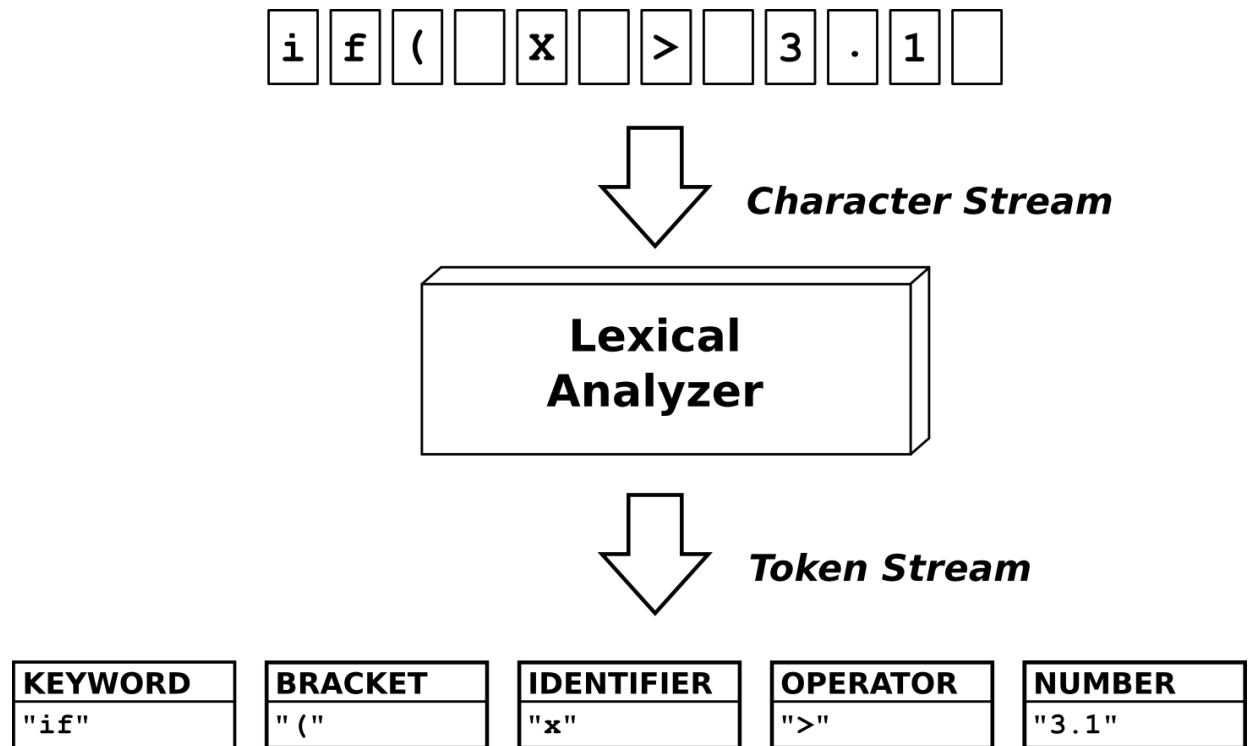
Token classification is the process of identifying the type of each token generated during tokenization. This can involve identifying keywords, identifiers, operators, punctuation symbols, and other token types as separate categories.

Token validation involves checking each token generated during tokenization to ensure that it is valid according to the rules of the programming language. For example, it might check that a variable name is a valid identifier, or that an operator has the correct syntax.

Output generation is the final stage of the lexical analysis process, where the lexer generates the output, which is typically a list of tokens. This list of tokens can then be passed on to the next stage of compilation or interpretation.



Tokenization of a text:



HARDWARE AND SOFTWARE REQUIREMENTS

Hardware components used in this process are -

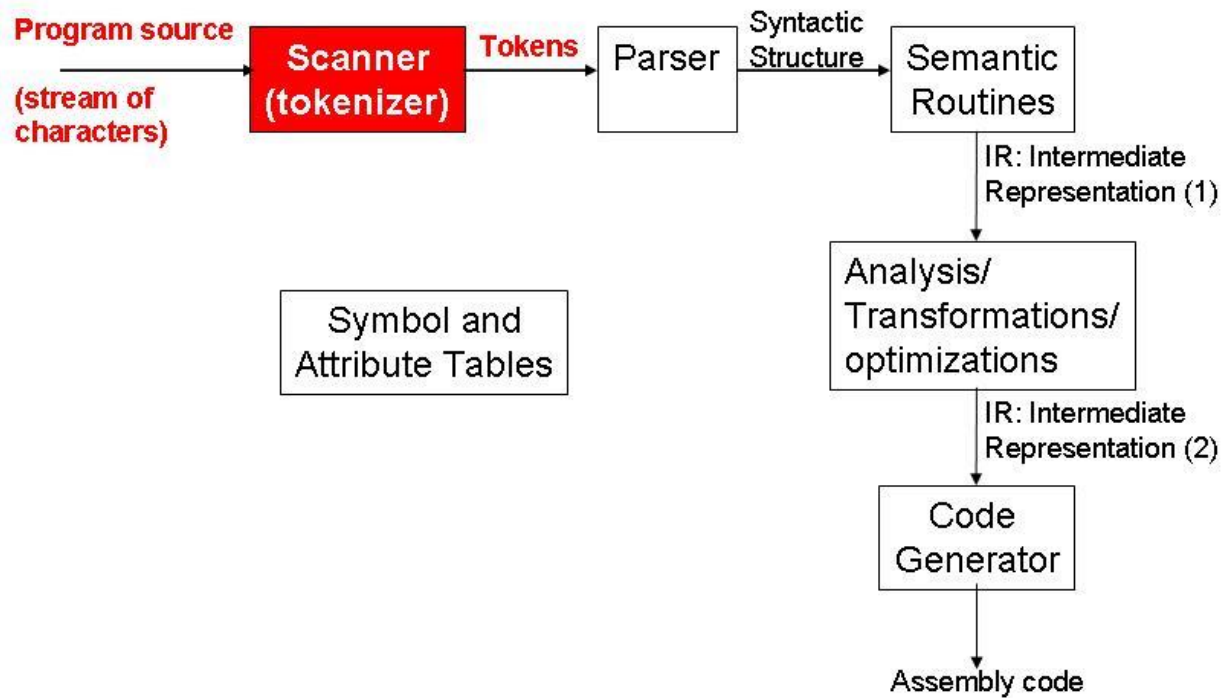
- Input device: This component is responsible for receiving the input source code written in the C language, such as a keyboard or a file storage device.
- Processor: The processor is responsible for executing the scanner program, which performs the lexical analysis of the input source code and generates a stream of tokens as output.
- Memory: The scanner program requires temporary memory to store the current lexeme being analyzed and the current state of the finite-state machine used to recognize the input source code's tokens.

Software components used in this process are -

- Scanner program: The scanner program is responsible for performing the lexical analysis of the input source code, which involves reading the input character by character, identifying and categorizing each character into a token, and generating a stream of tokens as output.
- Symbol table: The symbol table is a data structure used by the scanner to store information about each identifier found in the input source code, such as its name, type, and scope.
- Error handling module: The error handling module is responsible for detecting and reporting any lexical errors, such as misspelled keywords or undefined identifiers, in the input source code.
- Regular expressions: Regular expressions are used by the scanner to recognize patterns in the input source code and categorize them into tokens.

Input Stream -> Buffer -> Lexical Analyzer (Scanner) -> Token Stream

BLOCK DIAGRAM



METHODOLOGY USED

We have used Flex to perform lexical analysis on a subset of the C programming language.

Flex is a lexical analyzer generator that takes in a set of descriptions of possible tokens and produces a C file that performs lexical analysis and identifies the tokens.

This document is divided into the following sections:

- **Functionality:** Contains a description of our Flex program and the variety of tokens that it can identify and the error handling strategies.
- **Symbol table and Constants table:** Contains an overview of the architecture of the symbol and constants table which contain descriptions of the lexemes identified during lexical analysis.
- **Code Organisation:** Contains a description of the files used for lexical analysis
- **Source Code:** Contains the source code used for lexical analysis

The entire code for lexical analysis is broken down into 3 files: lexer.l, tokens.h and symboltable.h

FILE	CONTENTS
lexer.l	A lex file containing the lex specification of regular expressions
tokens.h	Contains enumerated constants for keywords, operator, special symbols, constants and identifiers.
symboltable.h	Contains the definition of the symbol table and the constants table and also defines functions for inserting into the hash table and displaying its contents.

CODE

```
File Edit Selection View Go Run Terminal Help LEX.cpp - CD_Project - Visual Studio Code

EXPLORER
CD_PROJECT
  abctxt
  LEX.cpp
  LEX.exe
  OutputFile.txt

LEX.cpp ...
LEX.cpp ...
1 #include <bits/stdc++.h>
2 #include <regex>
3 #include <time.h>
4 #include <iterator>
5 #include <windows.h>
6 #define deb(x) cout<<#x<<" = "<<<<endl
7
8 using namespace std;
9
10 map<string,string> Make_Regex_Map(){
11     map<string,string> my_map {
12         {"\\s|\\W|\\(|\\)|\\{|\\}|\\[|\\]|\\^|\\$", "Special Symbol"},
13         {"int|char|float|bool|cin|cout|main|using|namespace|std", "Keywords"},
14         {"\\#include|define", "Pre-Processor Directive"},
15         {"\\|ostream|\\|stdio|\\|string", "library"},
16         {"\\*|\\+|\\>|\\<|\\<|\\>", "operator"},
17         {"[0-9]+", "Integer"},
18         {"[*\\include][*ostream][*int][*main][*cin][*cout][*];[*][*][*B ;cin][*a-z]+", "Identifier"},
19         {"[A-Z]+", "Variable"},
20         {"[", ""},
21     };
22     return my_map;
23 }
24
25 map<size_t,pair<string,string>> Match_language (map<string,string> patterns,string str){
26
27     map< size_t, pair<string,string> > lang_matches;
28
29     for ( auto i = patterns.begin(); i != patterns.end(); ++i )
30     {
31         regex compare(i->first);
32         auto words_begin = sregex_iterator( str.begin(), str.end(), compare );
33         auto words_end = sregex_iterator();
34         //MAKING PAIRS OF [STRING OF REGEX 'compare' : 'pattern']
35         for ( auto it = words_begin; it != words_end; ++it )
36             lang_matches[ it->position() ] = make_pair( it->str(), i->second );
37     }
38 }
```

The screenshot shows the Visual Studio Code interface with the Explorer sidebar on the left. The Explorer shows a project named 'CD_PROJECT' with files 'abc.txt', 'LEX.cpp', 'LEX.exe', and 'OutputFile.txt'. The main editor window displays the C++ code for 'LEX.cpp'. The code defines a function 'tell_lexeme' that maps operators to tokens and a 'main' function that reads a source code file, fetches tokens, and uses a map to match them. The status bar at the bottom indicates the file is 'LEX.cpp', line 21, column 18, with 4 spaces, UTF-8 encoding, and CRLF line endings.

```

39     }
40
41     string tell_lexeme(string op){
42         if(op=="*") return "MUL";
43         else if(op=="+") return "ADD";
44         else if(op==">>") return "INS";
45         else if(op=="<<") return "EXTR";
46         else if(op==">") return "RSHFT";
47         else if(op=="<") return "LSHFT";
48     }
49     int main()
50     {
51         ofstream fout;
52         cout<<endl<<endl<<endl;
53         cout.fill(' ');
54         cout.width(100);
55         fout.open("OutputFile.txt");
56         char c;
57         string filename;
58
59
60         cout<<"ENTER THE SOURCE CODE FILE NAME: Example \"abc.txt\" \n";
61         cin>>filename;
62         fstream fin(filename, fstream::in);
63         string str;
64         //Fetching Source code in string type 'str'
65         if(fin.is_open()){
66             while(fin>> noskipws>>c)
67                 str+=c;
68
69             //Making a map which will define the regex in source code to its pattern in my language.
70             map<string,string> patterns =Make_Regex_Map();
71
72             //DECLARING MAP 'lang_matches' from 'patterns' map which will pair up the patterns
73             from the ['Source Code':Defined Pattern] via a Regex named 'compare'. */
74             map<size_t, pair<string,string> > lang_matches = Match_Language(patterns,str);
75

```

[illegible]

```
145     }
146   }
147 }
148
149 string command= " ";
150
151 while(command != "EXIT"){
152   cout.fill(' ');
153   cout.width(40);
154   cout<<"\n\n\t PRESS TYPE 'EXIT' TO CLOSE WINDOW.\n\t NOTE: AN OUTPUT FILE WILL BE GENERATED IN THE SAME FOLDER AS 'output.txt' \n";
155   cin.width(40);
156   cin>>command;
157
158   if(command == "exit"||command == "EXIT"|| command == "Exit")
159     break;
160
161   else{
162     cout.fill(' ');
163     cout.width(40);
164     cout<<"Please enter correct word.";
165     cin.width(10);
166     cin>>command;
167   }
168
169 }
170
171
172 else{
173   cout.fill(' ');
174   cout.width(40);
175   cout<<"\n FILE NOT FOUND!\n\n";
176 }
177 return 0;
178 }
```

RESULTS

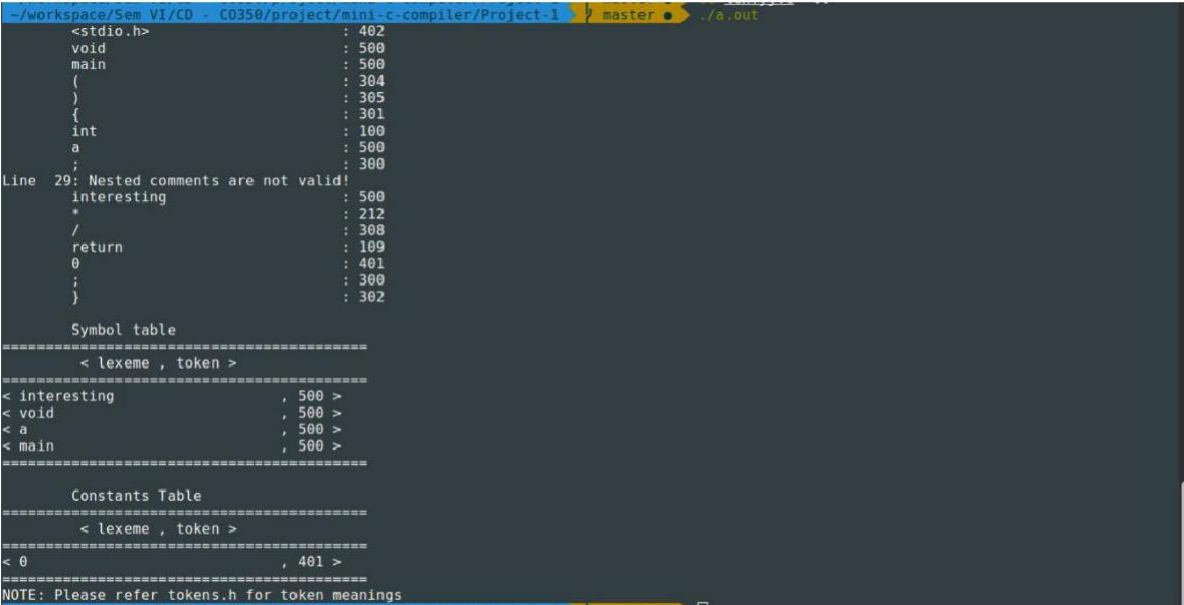
Test Case 1

- Test for single line comments
- Test for multi-line comments
- Test for single line nested comments
- Test for multiline nested comments

The output in lex should remove all the comments including this one
*/

```
#include<stdio.h>
```

```
void main(){  
    // Single line comment  
  
    /* Multi-line comment  
    Like this */  
  
    /* here */ int a; /* "int a" should be untouched */  
  
    // This nested comment // This comment should be removed should be removed  
  
    /* To make things /* nested multi-line comment */ interesting */  
  
    return 0;  
}
```



```
~/workspace/Sem VI/CD - C0350/project/mini-c-compiler/Project-1 master • ./a.out  
<stdio.h> : 402  
void : 500  
main : 500  
( : 304  
) : 305  
{ : 301  
int : 100  
a : 500  
; : 300  
Line 29: Nested comments are not valid!  
interesting : 500  
* : 212  
/ : 300  
return : 100  
0 : 401  
; : 300  
} : 302  
  
Symbol table  
=====  
< lexeme , token >  
=====  
< interesting , 500 >  
< void , 500 >  
< a , 500 >  
< main , 500 >  
=====  
  
Constants Table  
=====  
< lexeme , token >  
=====  
< 0 , 401 >  
=====  
NOTE: Please refer tokens.h for token meanings
```

test-case-1

Test Case 2

- Test for multi-line comment that doesn't end till EOF

The output in lex should print as error message when the comment does not terminate
It should remove the comments that terminate

*/

```
#include<stdio.h>
```

```
void main(){
```

```
    // This is fine
```

```
    /* This as well
```

```
    like we know */
```

```
    /* This is not fine since
```

```
    this comment has to end somewhere
```

```
    return 0;
```

```
}
```

```
~/workspace/Sem VI/CD - C0350/project/mini-c-compiler/Project-1 ? master • ./a.out
<stdio.h> : 402
void : 500
main : 500
( : 304
) : 305
{ : 301
Line 21: Unterminated comment

Symbol table
=====
< lexeme , token >
=====
< void , 500 >
< main , 500 >
=====

Constants Table
=====
< lexeme , token >
=====
NOTE: Please refer tokens.h for token meanings
~/workspace/Sem VI/CD - C0350/project/mini-c-compiler/Project-1 ? master •
```

test-case-2

Test Case 3

- Test for string
- Test for string that doesn't end till EOF
- Test for invalid header name

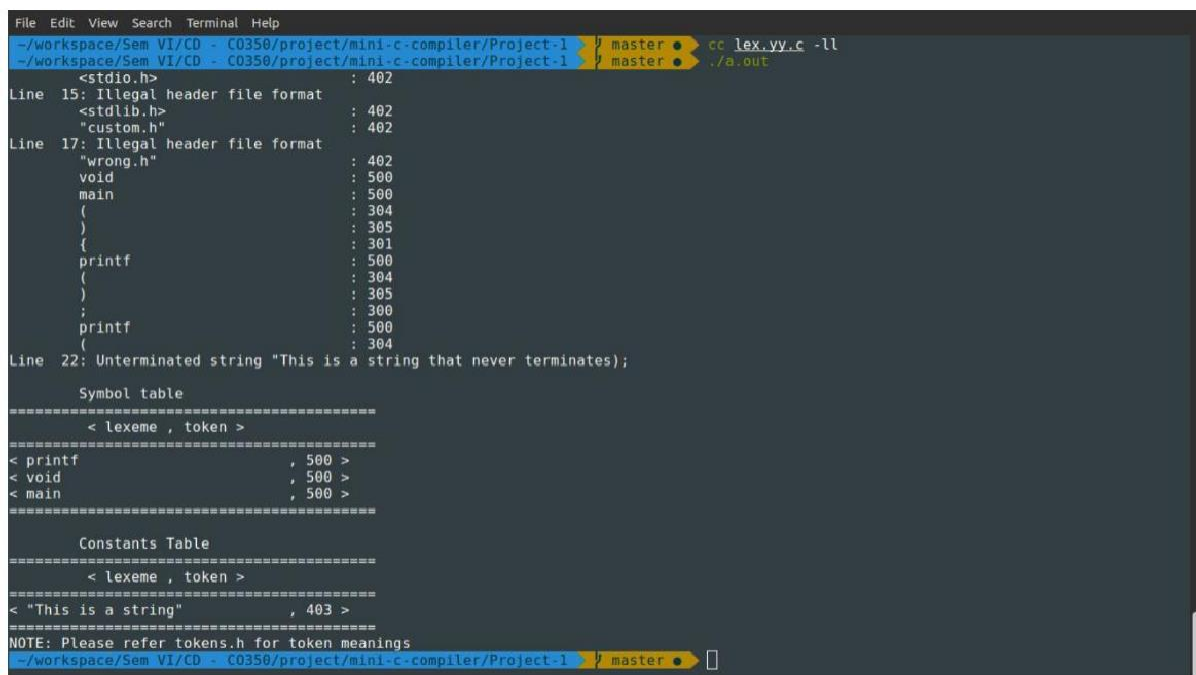
The output in lex should identify the first string correct and display error message that the second one does not terminate

*/

```
#include<stdio.h>
#include <<stdlib.h>
#include "custom.h"
#include ""wrong.h"

void main(){

    printf("This is a string");
    printf("This is a string that never terminates);
}
```



```
File Edit View Search Terminal Help
~/workspace/Sem VI/CD - C0350/project/mini-c-compiler/Project-1 master cc lex.vy.c -ll
~/workspace/Sem VI/CD - C0350/project/mini-c-compiler/Project-1 master ./a.out
<stdio.h> : 402
Line 15: Illegal header file format : 402
<stdlib.h> : 402
"custom.h" : 402
Line 17: Illegal header file format : 402
"wrong.h" : 402
void : 500
main : 500
( : 304
) : 305
{ : 301
printf : 500
( : 304
) : 305
; : 300
printf : 500
( : 304
Line 22: Unterminated string "This is a string that never terminates);

Symbol table
=====
< lexeme , token >
=====
< printf , 500 >
< void , 500 >
< main , 500 >
=====

Constants Table
=====
< lexeme , token >
=====
< "This is a string" , 403 >
=====
NOTE: Please refer tokens.h for token meanings
~/workspace/Sem VI/CD - C0350/project/mini-c-compiler/Project-1 master
```

test-case-3

Test Case 4

Following errors must be detected

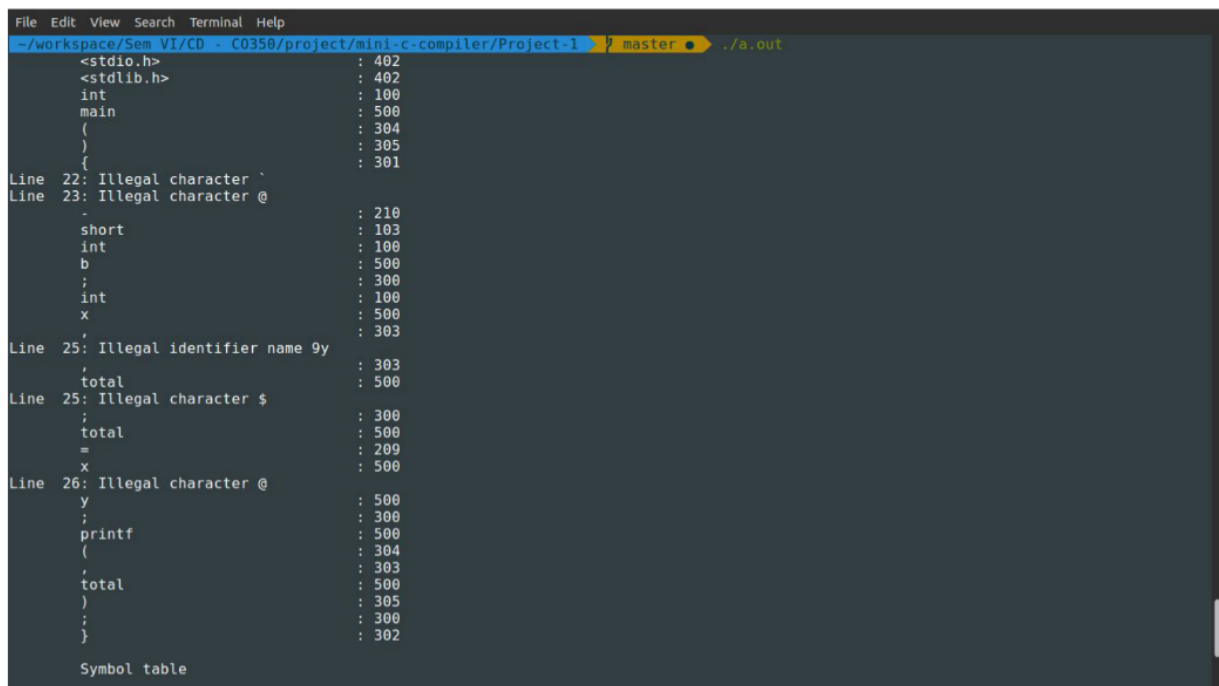
- Invalid identifiers: 9y, total\$
- Invalid operator: @
- Escaped quoted should be part of the string that is identified
- Stray characters: `, @, -

The output should display appropriate errors

```
*/
```

```
#include<stdio.h>
#include<stdlib.h>
```

```
int main()
{
    @ -
    short int b;
    int x, 9y, total$;
    total = x @ y;
    printf ("Total = %d \n \" ", total);
}
```



```
File Edit View Search Terminal Help
~/workspace/Sem VI/CD - C0350/project/mini-c-compiler/Project-1 master ● ./a.out
<stdio.h> : 402
<stdlib.h> : 402
int : 100
main : 500
( : 304
) : 305
{ : 301
Line 22: Illegal character `
Line 23: Illegal character @
- : 210
short : 103
int : 100
b : 500
; : 300
int : 100
x : 500
, : 303
Line 25: Illegal identifier name 9y
total : 500
Line 25: Illegal character $
; : 300
total : 500
= : 200
x : 500
Line 26: Illegal character @
y : 500
; : 300
printf : 500
( : 304
, : 303
total : 500
) : 305
; : 300
} : 302
Symbol table
```

Test-case-4a

```

File Edit View Search Terminal Help
: 303
total : 500
Line 25: Illegal character $
: 300
total : 500
= : 200
x : 500
Line 26: Illegal character @
y : 500
: 300
printf : 500
( : 304
: 303
total : 500
) : 305
: 300
} : 302

Symbol table
=====
< lexeme , token >
=====
< total , 500 >
< printf , 500 >
< x , 500 >
< main , 500 >
< b , 500 >
< y , 500 >
=====

Constants Table
=====
< lexeme , token >
=====
< "Total = %d \n \n " , 403 >
=====
NOTE: Please refer tokens.h for token meanings
~/workspace/Sem VI/CD - C0350/project/mini-c-compiler/Project-1 master

```

test-case-4b

/*

Test Case 5

Identifying tokens and displaying symbol and constants table

Following tokens must be detected

- Keywords (int, long int, long long int, main include)
- Identifiers (main, total, x, y, printf),
- Constants (-10, 20, 0x0f, 1234561)
- Strings ("Total = %d \n")
- Special symbols and Brackets ((), {}, ;, ,)
- Operators (+, -, =, *, /, %, --, ++)

The output should display appropriate tokens with their type and also the symbol and constants table

*/

#include<stdio.h>

#include<stdlib.h>

int main()

{

int x, y;

long long int total, diff;

int *ptr;

unsigned int a = 0x0f;

long int mylong = 1234561;

long int i, j;

for(i=0; i < 10; i++){

for(j=10; j > 0; j--){

printf("%d",i);

}

}

x = -10, y = 20;

x=x*3/2;

total = x + y;

diff = x - y;

int rem = x % y;

printf ("Total = %d \n", total);

}

```
File Edit View Search Terminal Help
~/workspace/Sem VI/CD - C0350/project/mini-c-compiler/Project-1 master ● ./a.out
<stdio.h> : 402
<stdlib.h> : 402
int : 100
main : 500
{ : 304
} : 305
{ : 301
int : 100
x : 500
, : 303
y : 500
; : 300
long long : 102
int : 100
total : 500
, : 303
diff : 500
; : 300
int : 100
* : 212
ptr : 500
; : 300
unsigned : 105
int : 100
a : 500
= : 209
0x0f : 400
; : 300
long : 101
int : 100
mylong : 500
= : 209
123456l : 401
; : 300
long : 101
int : 100
i : 500
```

test-case-5a

```
File Edit View Search Terminal Help
int : 100
i : 500
, : 303
j : 500
; : 300
for : 106
( : 304
i : 500
= : 209
0 : 401
; : 300
i : 500
< : 214
10 : 401
; : 300
i : 500
++ : 201
) : 305
{ : 301
for : 106
( : 304
j : 500
= : 209
10 : 401
; : 300
j : 500
> : 215
0 : 401
; : 300
j : 500
-- : 200
) : 305
{ : 301
printf : 500
( : 304
, : 303
i : 500
) : 305
```

test-case-5b


```

File Edit View Search Terminal Help
1      : 500
)      : 305
;      : 300
}      : 302
}      : 302
x      : 500
=      : 209
-10    : 401
,      : 303
y      : 500
=      : 209
20     : 401
;      : 300
x      : 500
=      : 209
x      : 500
*      : 212
3      : 401
/      : 308
2      : 401
;      : 300
total  : 500
=      : 209
x      : 500
+      : 211
y      : 500
;      : 300
diff   : 500
=      : 209
x      : 500
-      : 210
y      : 500
;      : 300
int     : 100
rem     : 500
=      : 209
x      : 500
%      : 213

```

test-case-5c

```

File Edit View Search Terminal Help
;      : 300
}      : 302

Symbol table
=====
< lexeme , token >
=====
< total      , 500 >
< printf     , 500 >
< j          , 500 >
< x          , 500 >
< i          , 500 >
< a          , 500 >
< mylong     , 500 >
< main       , 500 >
< rem        , 500 >
< y          , 500 >
< ptr        , 500 >
< diff       , 500 >
=====

Constants Table
=====
< lexeme , token >
=====
< -10        , 401 >
< 3          , 401 >
< 10         , 401 >
< 20         , 401 >
< 123456l    , 401 >
< "%d"       , 403 >
< 0x0f       , 400 >
< 2          , 401 >
< "Total = %d \n" , 403 >
< 0          , 401 >
=====
NOTE: Please refer tokens.h for token meanings
~/workspace/Sem VI/CD - C0350/project/mini-c-compiler/Project-1 master

```

test-case-5d

CONCLUSION

The lexical analyzer that was created in this project helps us to break down a C source file into tokens as per the C language specifications. Each token (such as identifiers, keywords, special symbols, operators, etc.) has an integer value associated with it.

When we design the parser in the next phase, the parser will call upon the Flex program to give it tokens and the lexical analyzer will return to the parser the integer value associated with the tokens as and when required by the parser. Together with the symbol, the parser will prepare a syntax tree with the help of a grammar that we provide it with. The parser can then logically group the tokens to form meaningful statements and can detect C programming constructs such as arrays, loops, and functions. The parser will also help us identify errors that could not be detected in the lexical analysis phase such as unbalanced parentheses, unterminated statements, missing operators, two operators in a row, etc.