# Lexical Analyser for C

Aishwarya Rajan - 13CO109
Anagh Singh - 13CO110

January 28, 2016



# Department of Computer Science and Engineering, NITK

# Contents

# List of Figures

# Listings

**Abstract**

Lexical analyzer for c language is developed using Lex/Flex tools available that analyse C programs. This analyser supports nested comments and returns meaningful errors if there are any. This program also analyses all the comments and strings that it gets as input till the end of the file.

# 1 Introduction

Lexical analysis is the process of analysing a stream of individual characters (normally arranged as lines), into a sequence of lexical tokens (tokenization - for instance of "words" and punctuation symbols that make up source code) to feed into the parser. Roughly the equivalent of splitting ordinary text written in a natural language (e.g. English) into a sequence of words and punctuation symbols. Lexical analysis is often done with tools such as lex, flex and jflex.

Strictly speaking, tokenization may be handled by the parser. The reason why we tend to bother with tokenising in practice is that it makes the parser simpler, and decouples it from the character encoding used for the source code.

## 1.1 Lexical Analysis

### 1.1.1 What is a token?

In computing, a token is a categorized block of text, usually consisting of indivisible characters known as lexemes. A lexical analyzer initially reads in lexemes and categorizes them according to function, giving them meaning. This assignment of meaning is known as tokenization. A token can look like anything: English, gibberish symbols, anything; it just needs to be a useful part of the structured text. Tokens are frequently defined by regular expressions, which are understood by a lexical analyzer such as lex. The lexical analyzer reads in a stream of lexemes and categorizes them into tokens. This is called "tokenizing." If the lexer finds an invalid token, it will report an error. Following tokenizing is parsing. From there, the interpreted data may be loaded into data structures, for general use, interpretation, or compiling.

### 1.1.2 What is a pattern?

There is a set of strings in the input for which the same token is produced as output. This set of strings is described by a rule called a pattern associated with the token. Regular expressions are an important notation

for specifying patterns. For example, the pattern for the Pascal identifier token, id, is:

$$\text{id} \rightarrow \text{letter (letter | digit)*}$$

### 1.1.3 Finite State Automaton

An FSA is usually used to do lexical analysis.

An FSA consists of states, starting state, accept state and transition table. The automaton reads an input symbol and moves the state accordingly. If the FSA reaches the accept state after the input string is read until its end, the string is said to be accepted or recognized. A set of recognized strings is said to be a language recognized by the FSA.

### 1.1.4 Regular Expressions

The regular expressions over an alphabet specify a language according to the following rules. $\epsilon$ is a regular expression that denotes , that is, the set containing the empty string. If $a$ is a symbol in alphabet, then $a$ is a regular expression that denotes $a$, that is, the set containing the string $a$. Suppose $r$ and $s$ are regular expression denoting the languages $L(r)$ and $L(s)$. Then

- $(r)|(s)$ is a regular expression denoting $L(r)UL(s)$.

- $(r)(s)$ is a regular expression denoting $L(r)L(s)$.

- $(r)*$ is a regular expression denoting $(L(r))*$.

- $(r)|(s)$ is a regular expression denoting $L(r)UL(s)$.

## 1.2 Flex Script

FLEX (Fast LEXical analyzer generator) is a tool for generating scanners. In stead of writing a scanner from scratch, we only need to identify the vocabulary of a certain language (e.g. Simple), write a specification of patterns using regular expressions (e.g. DIGIT [0-9]), and FLEX will construct a scanner for it.

*flex* is a tool for generating scanners. A scanner is a program which recognizes lexical patterns in text. The flex program reads the given input files, or its standard input if no file names are given, for a description of a scanner to generate. The description is in the form of pairs of regular expressions and C code, called rules. flex generates as output a C source file, *lex.yy.c* by default, which defines a routine *yylex()*. This file can be compiled and linked with the flex runtime library to produce an executable.

When the executable is run, it analyzes its input for occurrences of the regular expressions. Whenever it finds one, it executes the corresponding C code.

The internals of a Flex script are as follows: 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

The flex input file consists of three sections, separated by a line containing only '%%'.

```
definitions
%%
rules
%%
user code
```

Character classes are expanded immediately when seen in the flex input. This means the character classes are sensitive to the locale in which flex is executed, and the resulting scanner will not be sensitive to the runtime locale. This may or may not be desirable.

### 1.2.1   Input Matching

When the generated scanner is run, it analyzes its input looking for strings which match any of its patterns. If it finds more than one match, it takes the one matching the most text (for trailing context rules, this includes the length of the trailing part, even though it will then be returned to the input). If it finds two or more matches of the same length, the rule listed first in the flex input file is chosen.

Once the match is determined, the text corresponding to the match (called the token) is made available in the global character pointer yytext, and its length in the global integer yyleng. The action corresponding to the matched pattern is then executed, and then the remaining input is scanned for another match.

If no match is found, then the default rule is executed: the next character in the input is considered matched and copied to the standard output. Thus, the simplest valid flex input is:

```
                    %%
```

### 1.2.2 Generated Scanner

The output of flex is the file lex.yy.c, which contains the scanning routine
yylex(), a number of tables used by it for matching tokens, and a number
of auxiliary routines and macros. By default, yylex() is declared as follows:

```
int yylex()
    {
    ... various definitions and the actions in here ...
    }
```

Whenever yylex() is called, it scans tokens from the global input file yyin
(which defaults to stdin). It continues until it either reaches an end-of-file
(at which point it returns the value 0) or one of its actions executes a return
statement. If yylex() stops scanning due to executing a return statement in
one of the actions, the scanner may then be called again and it will resume
scanning where it left off. The scanner writes its ECHO output to the
yyout global (default, stdout), which may be redefined by the user simply
by assigning it to some other FILE pointer.

## 1.3 C Program

This section contains the C progam which we provide as input to our scan-
ner.

```
 1   /* C Program  to  Swap  two  numbers
 2    ****  Using  temporary  variable
 3   */
 4   #include <stdio.h>
 5
 6   int main()
 7   {
 8       int x, y, temp;
 9
10       printf("Enter  the  value  of  x  and  y\n");
11       scanf("%d%d", &x, &y);
12       x=2;
13       printf("Before  Swapping\nx = %d\ny = %d\n",x,y);
14
15       temp = x;
16       x     = y;
17       y     = temp;
18
19       printf("After  Swapping\nx = %d\ny = %d\n",x,y);
20
21       return 0;
22   }
```

This is another program that we provide as input:

```
1   /* C Program to Swap two numbers
2    */Using temporary variable*/
3   */
4   #include <stdio.h>
5
6   int main()
7   {
8       int x, y, temp;
9
10      printf("Enter the value of x and y\n");
11      scanf("%d%d", &x, &y);
12
13      printf("Before Swapping\nx = %d\ny = %d\n",x,y);
14
15      temp = x;
16      x    = y;
17      y    = temp;
18
19      printf("After Swapping\nx = %d\ny = %d\n",x,y);
20
21      return 0;
22  }
```

# 2   Design of Program

## 2.1   Code

This section contains our actual scanner. The code for it is:

```
1   %{
2   #include <stdio.h>
3
4   #include <stdlib.h>
5
6   #include <malloc.h>
7
8   #include <string.h>
9
10  int comment=0,bracCount=0;
11
12  char* ident[50]; //Symbol Table
13  char* constant[50]; //Constants Table
14  int it=0, ct=0;
15
16  #define OPERATOR      1
17  #define LCURLY        2
18  #define RCURLY        3
19  #define COMMENT       4
20  #define SEMICOLON     5
21  #define KEYWORD       6
22  #define IDENTIFIER    7
23  #define INTEGER       8
24  #define STRING        9
25  #define FUNC          10
26  #define FLOAT         11
27  #define PUNCTUATOR    12
28
```

```
29   %}
30
31   /* Definitions for finite Automata */
32
33   digit [0−9]
34
35   alpha [a−zA−Z_]
36
37   comment (\/\/.*)
38
39   comstart (\/\*)
40
41   comend (\*\/)
42
43   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"
44
45   relop >|<|<=|>=|!=
46
47   delimiter [ \t\v\n]+
48
49   %%
50
51
52
53   ^#([−a−zA−Z0−9.]|{relop}|{delimiter})* {
54                            fprintf(yyout,"\n%s\t\t\tPREPROCESSOR DIRECTIVE
                                \t\t\t\t",yytext);
55
56                            }
57
58   {comment}              {
59              if(comment<=0) fprintf(yyout,"\n//\t\t\t\t\t LINE
                   COMMENT \t\t\t\t\t");
60
61              }
62
63   {comstart}                {
64
65                         comment++;
66
67                         if(comment>1)
68
69                           fprintf(yyout,"\nERROR :: Nested Comment
                               found.");
70
71                         else
72
73                           fprintf(yyout,"\n/*\t\t MultiLine Comment
                               Begins\t\t\t");
74
75                       }
76
77
78   {comend}              {
79
80                         if(comment>0)
81
82                          {
```

9

```
83
84                              comment−−;
85
86                              fprintf(yyout,"\n*/\t\tMultiLine Comment
                                    Ends\t\t\t");
87
88                          }
89
90                      else
91
92                          fprintf(yyout,"\n\tERROR: */ occurs before
                                    : /* ");
93
94                  }
95
96   {keyword}                              {if(comment<=0) fprintf(yyout,"\n%s
         \t\tKEYWORD\t\t",yytext);}
97
98   {alpha}+\(                             {if(comment<=0) fprintf(yyout,"\n%s
         \t \tFUNCTION\t\t",yytext);}
99
100  {alpha}({alpha}|{digit})*          {
101                      if(comment<=0)
102                      {
103                        fprintf(yyout,"\n%s\t\tIDENTIFIER\t\t ",yytext)
                                    ;
104                        ident[it++]=yytext;
105
106                      }
107                  }
108
109  {digit}+                        {
110                      if(comment<=0) {
111                        fprintf(yyout,"\n%s\t\tNUMERIC CONSTANT\t\t ",
                                    yytext);
112                        constant[ct++]=yytext;
113                      }
114                  }
115
116  {digit}+\.{digit}+              {if(comment<=0) {fprintf(yyout,"\n%
         s \t\t FLOAT CONSTANT\t\t ",yytext);
117                      constant[ct++]=yytext;
118                  }}
119
120  \"(\\.|[^\\"])*\"               {if(comment<=0) {fprintf(yyout,"\n
         %−39s STRING LITERAL\t\t ",yytext);
121                      //constant[ct++]=yytext;
122                  }}
123
124  \"(\\.|[^\\"])*                 {if(comment<=0) printf("ERROR: End
         of string not found.\n");}
125
126  "{"                            {if(comment<=0) {fprintf(yyout,"\n%s\t\
         tOPEN BRACKET\t\t ",yytext);bracCount++;}}
127
128  "}"                            {if(comment<=0) {fprintf(yyout,"\n%s\t\
         tCLOSE BRACKET\t\t ",yytext);bracCount−−;}}
129
130  ","|";"|":"|"("|")"|"."        {if(comment<=0) fprintf(yyout,"\n%s \t\
         tPUNCTUATOR\t\t ",yytext);}
131
```

```
132    "&"|"!"|"~"|"+"|"-"|"*"  {if(comment<=0) fprintf(yyout,"\n%s \t\
          tOPERATOR\t\t ",yytext);}
133    "/"|"%"|"<"|">"|"^"|"|"  {if(comment<=0) fprintf(yyout,"\n%s \t\
          tOPERATOR\t\t ",yytext);}
134    "?"|"="|{relop}|"..."     {if(comment<=0) fprintf(yyout,"\n%s \t\
          tOPERATOR\t\t ",yytext);}
135    ">>="|"<<="|"+="|"-="     {if(comment<=0) fprintf(yyout,"\n%s \t\
          tOPERATOR\t\t ",yytext);}
136    "*="|"/="|"%="|"&="        {if(comment<=0) fprintf(yyout,"\n%s \t\
          tOPERATOR\t\t ",yytext);}
137    "^="|"|="|">>"|"<<"|"+"  {if(comment<=0) fprintf(yyout,"\n%s \t\
          tOPERATOR\t\t ",yytext);}
138    "--"|"->"|"&&"|"||"|"<="  {if(comment<=0) fprintf(yyout,"\n%s \t\
          tOPERATOR\t\t ",yytext);}

140    {delimiter} {;}

142    %%

144    /* MAIN PROGRAM */

146    int main(int argc,char **argv){

148    if(argv[1]==NULL)

150      printf("Error opening file. Usage ./a.out <filename>\n");

152    else{

154        yyin=fopen(argv[1],"r");

156        printf("\n\t\t TOKEN IDENTIFICATION:\n Lexeme\t\t\tToken\n");

158        yylex();

160        if(comment!=0)

162          printf("\nERROR: Comment does not end\n");

164        if(bracCount!=0)

166          printf("\nERROR: Bracket mismatch\n");

168        printf("\n");

170        printf("\n it:%d, ct:%d\n************Symbol Table************\
          n",it,ct);
171        int index,k,j,flag;
172        char x[50][50];
173         for(index=0;index<it;index++){
174            for(k=0;k<strlen(ident[index]);k++)
175            {
176                if (!isalnum(*(ident[index]+k))){
177                    x[index][k]='\0';
178                    break; }
179                x[index][k]=*(ident[index]+k);
180            }
181            flag=0;
182            for(j=0;j<index;j++)
183                if(strcmp(x[j],x[index])==0) {flag=1; break;}

185            if(!flag)
```

11

```
186              printf("—>%s\n",x[index]);}
187        printf("\n*************************************\n");
188
189        printf("\n************Constants  Table********\n");
190         for(index=0;index<ct;index++){
191
192      for(k=0;k<strlen(constant[index]);k++)
193           {
194                if((*(constant[index]+k))==';'){
195                   x[index][k]='\0';
196                   break;  }
197              x[index][k]=*(constant[index]+k);
198           }
199           printf("—>%s\n",x[index]);
200  }
201      printf("\n*************************************\n");
202
203      }
204
205  }
206
207  int yywrap()
208  {
209  return(1);
210  }
```

## 2.2 Explanation

Given lex.y file is code for our lexical analyser which is takes a C program as input and scan and analyses that code and works as given below.

### 2.2.1 Alphabet keywords

They are already predefined in their respective section.

### 2.2.2 To keep track of code blocks

This will take a section which is between/before or after comments (excludes comments).

### 2.2.3 To resolve nested comments

We resolved this situtation by keeping integer track of the depth of comments.

### 2.2.4 String errors

These are solved by checking last character of given string to check if it has ended or not.

# 3 Test Cases

| S.No | Test case | Steps | Expected Results | Status |
|------|-----------|-------|------------------|--------|
| 1 | Identifier | Enter an identifier<br>Enter an identifier<br>Enter an identifier<br>Enter an identifier | Var1 VALID IDENTIFIER<br>1dr INVALID IDENTIFIER<br>_var VALID IDENTIFIER<br>V1:b INVALID IDENTIFIER | Working |
| 2 | Integer Constant(int) | Enter an integer<br>Enter an integer<br>Enter an integer<br>Enter an integer<br>Enter an integer<br>Enter an integer | 3 VALID INTEGER<br>2.0 INVALID INTEGER<br>23,333 INVALID INTEGER<br>-55 VALID INTEGER<br>+10 VALID INTEGER<br>12938365675 INVALID INTEGER<br>09 INVALID INTEGER | Working |
| 3 | String Constant | Enter a string<br>Enter a string<br>Enter a string<br>Enter a string | "Horse" VALID STRING<br>""Horse" INVALID STRING<br>Horse INVALID STRING<br>"100, Rs" VALID STRING | Working |
| 4 | Character Constant | Enter a character<br>Enter a character<br>Enter a character<br>Enter a character<br>Enter a character<br>Enter a character<br>Enter a character | 'A' VALID CHARACTER<br>"n" INVALID CHARACTER<br>'kg' INVALID CHARACTER<br>'' VALID CHARACTER<br>'V́ALID CHARACTER<br>'$iN$' INVALID CHARACTER<br>'9' VALID CHARACTER | Working |
| 5 | Relational Operators | Enter an R-operator<br>Enter an R-operator<br>Enter an R-operator<br>Enter an R-operator<br>Enter an R-operator<br>Enter an R-operator | >VALID R-OPERATOR<br>>= VALID R-OPERATOR<br>= INVALID R-OPERATOR<br>== VALID R-OPERATOR<br>!= VALID R-OPERATOR<br>- INVALID R-OPERATOR | Working |
| 6 | Real Constant | Enter a real constant<br>Enter a real constant<br>Enter a real constant<br>Enter a real constant<br>Enter a real constant<br>Enter a real constant<br>Enter a real constant<br>Enter a real constant<br>Enter a real constant | 5.6 VALID CONSTANT<br>3.4.5 INVALID CONSTANT<br>34 INVALID CONSTANT<br>34 456 INVALID CONSTANT<br>34,477 INVALID CONSTANT<br>-364.8 VALID CONSTANT<br>+ie INVALID CONSTANT<br>+3.0 VALID CONSTANT<br>314159E-5L<br>VALID CONSTANT | Working |

| | | | | |
|---|---|---|---|---|
| 7 | Keywords | Enter a keyword<br>Enter a keyword<br>Enter a keyword<br>Enter a keyword<br>Enter a keyword<br>Enter a keyword<br>Enter a keyword | double VALID KEYWORD<br><br>main INVALID KEYWORD<br>case VALID KEYWORD<br>INT INVALID KEYWORD<br>while VALID KEYWORD<br>yes! INVALID KEYWORD<br>auto VALID KEYWORD | Working |
| 8 | Nested Comments | Enter a comment<br>Enter a comment<br>Enter a comment | /*hello /*world*/ VALID<br>/*hi /*hello*/*/ INVALID<br>/*yes*/*/ INVALID | |
| 9 | Delimiters | Enter a delimiter<br>Enter a delimiter<br>Enter a delimiter<br>Enter a delimiter<br>Enter a delimiter<br>Enter a delimiter<br>Enter a delimiter<br>Enter a delimiter<br>Enter a delimiter | { VALID<br>; VALID<br>: VALID<br>, VALID<br># VALID<br>[ VALID<br>) VALID<br>   INVALID<br>. INVALID | Working |
| 10 | Unary Operators | Enter a U-operator<br>Enter a U-operator<br>Enter a U-operator<br>Enter a U-operator<br>Enter a U-operator<br>Enter a U-operator<br>Enter a U-operator<br>Enter a U-operator<br>Enter a U-operator | ++ VALID<br>+ VALID<br>sizeof VALID<br>* VALID<br>& VALID<br>   VALID<br>! VALID<br>/ INVALID<br>. INVALID | Working |
| 11 | Binary Operators | Enter a B-Operator<br>Enter a B-Operator<br>Enter a B-Operator<br>Enter a B-Operator<br>Enter a B-Operator<br>Enter a B-Operator<br>Enter a B-Operator<br>Enter a B-Operator | + VALID<br>- VALID<br>/ VALID<br>* VALID<br>—— VALID<br>++ INVALID<br>! INVALID<br>{ INVALID | Working |
| 12 | Ternary Operator | Enter a T-Operator<br>Enter a T-Operator | ?: VALID<br>!: INVALID | Working |

# 4  Testing With Errors



Figure 1: Errors in Nested Comments

# 5  Implementation

For implementation of this given code we used following technique:

## 5.1  To keep track of blocks

```
int comment=0,bracCount=0
comment (\/\/.*)  //line comment
comstart (\/\*)   //multi line comment starts
comend (\*\/)  //multi line comment ends
```

## 5.2  Keyword Definitions

```
"auto"|"break"|"case"|"char"|"const"|"continue"|"default"|"do"|"dou
ble"|"else"|"enum"|"extern"|"float"|"for"|"goto"|"if"|"int"|"long"|
"register"|"return"|"short"|"signed"|"sizeof"|"static"|"struct"|"sw
```

```
itch"|"typedef"|"union"|"unsigned"|"void"|"volatile"|"while"
 relop >|<|<=|>=|!=    //Relational Operators
whitespace [ \t]+            //Whitespace
RE for a PreProcessor Statement
^#([azAZ09.]|{relop}|{whitespace})*
```
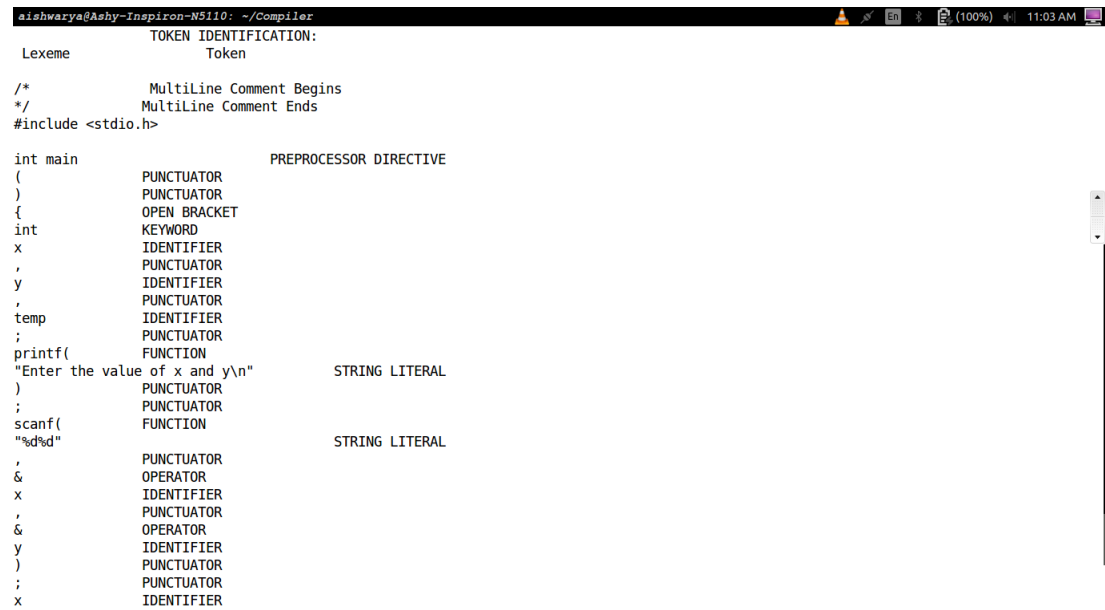
## 5.3  Checking for comment related errors

```
{comment}         {if(comment<=0) fprintf(yyout,"\n//Line Comment ")}
{comstart}        {comment++
                  if(comment>1)
                  fprintf(yyout,"\nERROR: Nested Comment")
                  else
                  fprintf(yyout,"\n/*ML Comment Sarts ")
        }
{comend}          {   if(comment>0)
                  {
                  comment
                  fprintf(yyout,"\n*/ML Comment Ends ")
                  }
                  else
                  fprintf(yyout,"\n\tERROR: */ found before /* ")
                  }
```

# 6 Results



```
aishwarya@Ashy-Inspiron-N5110: ~/Compiler                          En    (100%)   11:03 AM
                TOKEN IDENTIFICATION:
  Lexeme                Token

/*              MultiLine Comment Begins
*/              MultiLine Comment Ends
#include <stdio.h>

int main                        PREPROCESSOR DIRECTIVE
(               PUNCTUATOR
)               PUNCTUATOR
{               OPEN BRACKET
int             KEYWORD
x               IDENTIFIER
,               PUNCTUATOR
y               IDENTIFIER
,               PUNCTUATOR
temp            IDENTIFIER
;               PUNCTUATOR
printf(         FUNCTION
"Enter the value of x and y\n"          STRING LITERAL
)               PUNCTUATOR
;               PUNCTUATOR
scanf(          FUNCTION
"%d%d"                                  STRING LITERAL
,               PUNCTUATOR
&               OPERATOR
x               IDENTIFIER
,               PUNCTUATOR
&               OPERATOR
y               IDENTIFIER
)               PUNCTUATOR
;               PUNCTUATOR
x               IDENTIFIER
```

Figure 2: Successful Run Part 2

17

Figure 3: Successful Run Part 2

# References

[1] En.wikibooks.org, (2016). Compiler Construction/Lexical analysis - Wikibooks, open books for an open world. [online] Available at: https://en.wikibooks.org/wiki/Compiler _Construction/Lexical _analysis [Accessed 28 Jan. 2016].

[2] Cs.fsu.edu, (2016). [online] Available at: http://www.cs.fsu.edu/ xyuan/cop5621/lect2_lexical.ppt [Accessed 28 Jan. 2016].

[3] dragonbook.stanford.edu/lecture-notes/03-Lexical-Analysis.pdf

[4] cecs.wright.edu/ tkprasad/courses/cs780/L3Lexing.pdf

[5] quex.sourceforge.net/

[6] www.cs.nyu.edu/courses/spring11/G22.2130-001/lecture4.pdf