

George / Lexical Analyzer

Task: Write a program, using files (both input and output) that reads in a small C/C++ program using a postfix expression and identify the lexemes as to their respective tokens.

Language : C++

Overview: My program lexical analyzes a simple postfix expression from an input file (plain text file) called *list.txt*. It identifies the lexemes present in the postfix expression and outputs them alongside their respective token. The output of the lexemes and tokens are organized in a table. The table is sent to an output file called *output.txt*. My lexical analyzer recognizes a specific subset and keywords of the c++ language.

Description: First, I typed up all the headers for my code and made sure to `#include <string>` because I planned on using the built in string function `.length()`. Also, I *used* `#include <fstream>` to be able to read from /write to files. After, I created an object called *infile* of type `ifstream` in order to read the postfix expression from the input file *list.txt*. Then, I created an object called *outfile* of type `fstream` to write to the file *output.txt*. I used the `.open()` command to open both *list.txt* and *output.txt*.

```
ifstream infile; // infile (object) of type ifstream

infile.open("list.txt"); //read from input file called list.txt // list.txt default
file //input file stream

fstream outfile;

outfile.open("output.txt");
```

Then, I created an if statement that would check if the input file does not exist. If the input file does not exist this would make the stream invalid. The if statement uses a boolean context. If not *infile* then output the error message on to the screen and terminate the program.

```
if(!infile) { // if file does not open //stream variable boolean context cerr <<

    "Failed, could not open file "<< endl; //error message

    exit(1); //termination of program }
```

So, what are my lexemes and tokens ?

Before writing my program I already had in mind the lexemes I wanted my code to recognize and what the name of their respective tokens would be. The chart below shows the lexemes and tokens I used.

<u>Lexemes Tokens</u>
int double bool string float variable_name
main void function_name
cout output_stream
a b c d e f g h i z identifier
== equalto_operator
if else while c++_keywords
_ delimiter_op
+ - % / * id_operation
= assignment_op
{ open_bracket
} close_bracket
(open_parenth
) close_parenth
0 1 2 3 4 9 int_literal
; ending_operator

*The program tests for positive single digit numbers (0 -9). Therefore, if a postfix expression contains an integer with more than one digit each digit in the number is read separately. For instance, the number 657 would be recognized as three separate integer literals being 6, 5, and then 7.

*In addition, my program recognizes only lower case single alphabet letters as lexemes. For example, if a postfix expression used the variables *ABCD* or *abc* my program would not recognize it as a lexeme.

Main code:

I declared a string array called alpha that would represent the 26 letters in the alphabet. The lower case letters in the string array represent the different possible identifiers. Then, I declared a string and assigned it the variable name *postfix*. *Postfix* is the string I read from the input file.

```
string alpha[] = {"a","b", "c", "d", "e", "f","g", "h", "i", "j", "k","l", "m", "n",
"o", "p","q", "r", "s","t", "u","v", "w", "x", "y", "z"}; //string array
```

```
string postfix; // treat postfix expression in input file as a string
```

Then I decided to use a while loop. The condition in my while loop is (*infile >> postfix*.) I used the extraction operator so that my filestream (*infile*) reads in my string (*postfix*). With the >> operator I am able to read the postfix expression **word by word** rather than line by line . The space in between each word is used as a delimiter. Therefore, the values in the postfix expression that are separated by a space are considered two different values and read as individual words. I made sure that each individual value in the postfix expression is separated by a space before writing a new one. Then I declared six *if statements* that would check for specific string values. If postfix is equal to (==) a specific string it is outputted to the *output.txt* file on the left hand side of the table and the designated token name is outputted on the right hand side of the table.

```
while (infile >> postfix) { // reading my string word by word // infile acts as cin

    if(postfix == "main" || postfix == "void"){

        outfile << postfix << " |" << " Function_Name" << endl; }
```

After the six if statements I declared a *for loop*. A string is an array of characters and to read each character in the postfix expression I will need to read each index of that array. Thus, I decided that I would use a for loop to read through each index in the string. I initialized int i to 0 so the first index read in the string is 0. Next, I made the condition statement *i is less than postfix.length()*. I used the built-in string function .length(). The variable i will start from 0 and go up to the length of the string. I incremented i (*i++*) so each index that is read will start from 0 and it will keep going by 1 until the condition statement is no longer true. The for loop is then followed by nine if statements that check for specific characters (lexemes). These specific char values represent the new group of lexemes we want to recognize in the string. Next, I closed the for loop then the while loop.

```
for(int i =0; i < postfix.length(); i++) { // reads each char in the string

    if(postfix[i] == '{'){

        outfile << postfix[i] << " |" << " Open_bracket'"<< endl; //recognizing a
open bracket

    }
```

Lastly, I called `close()` for infile and outfile.

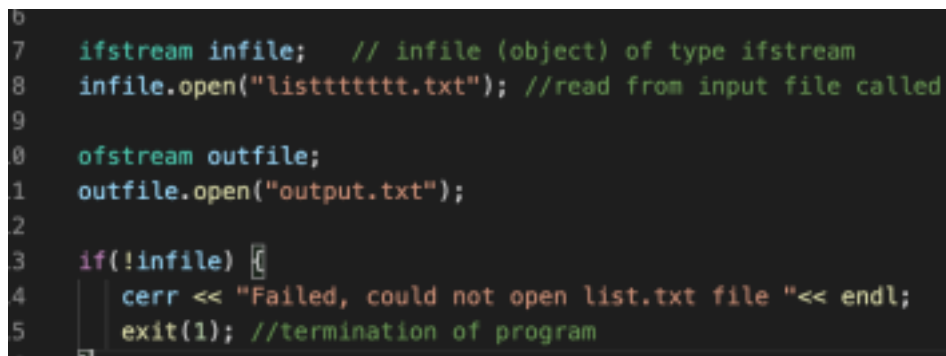
```
infile.close();

outfile.close();

return 0; }
```

Test Results:

1. Testing the error routine if it is not infile. This test uses an invalid file. (*Image 1a- test code , Image 1b test results*)



```
6
7  ifstream infile;    // infile (object) of type ifstream
8  infile.open("listtttttt.txt"); //read from input file called
9
10 ofstream outfile;
11 outfile.open("output.txt");
12
13 if(!infile) {
14     cerr << "Failed, could not open list.txt file "<< endl;
15     exit(1); //termination of program
```

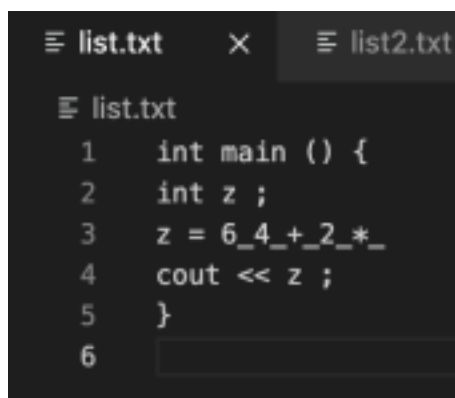
Image 1a



```
Failed, could not open file
```

Image 1b error routine works!

2. Testing the *while loop* to make sure the string is read word by word. I am using the input file *list.txt* with the following postfix expression to test the while loop. (*Image 2a- test code, Image 2b- test results*)



```
list.txt
1  int main () {
2  int z ;
3  z = 6_4_+_2_*_
4  cout << z ;
5  }
6
```

----->list.txt (postfix expression)

Image 2a

Image 2b

3. Testing the first group of *if statements*. The code should be able to recognize the lexemes present in the string and assign it to its appropriate token. (Image 3a- test code, Image 3b-

Image 3a

---->the lexemes & tokens
available for this test run

output.txt

1	-----	
2	Lexemes	Tokens
3	-----	
4	{	Open_parenth
5	}	Close_parenth
6	{	Open_bracket
7	;	Ending_Operator
8	=	Assign_op
9	6	Int_Literal
10	-	Delimiter_op
11	4	Int_Literal
12	-	Delimiter_op
13	+	id_Operation
14	-	Delimiter_op
15	2	Int_Literal
16	-	Delimiter_op
17	*	id_Operation
18	-	Delimiter_op
19	;	Ending_Operator
20	}	Close_bracket
21	-----	

Image 4b for loop and if statements work !

5. Testing that all parts of the program are successful. The table should be displayed in the output file. (cpp file attached) (Image 5b- test results)

output.txt

1	-----	
2	Lexemes	Tokens
3	-----	
4	int	Variable_Literal
5	main	Function_Name
6	{	Open_parenth
7	}	Close_parenth
8	{	Open_bracket
9	int	Variable_Literal
10	z	Identifier
11	;	Ending_Operator
12	z	Identifier
13	=	Assign_op
14	6	Int_Literal
15	-	Delimiter_op
16	4	Int_Literal
17	-	Delimiter_op
18	+	id_Operation
19	-	Delimiter_op
20	2	Int_Literal
21	-	Delimiter_op
22	*	id_Operation
23	-	Delimiter_op
24	cout	Output_Stream
25	<<	Output_Operator
26	z	Identifier
27	;	Ending_Operator
28	}	Close_bracket
29	-----	

Image 5b all the lexemes are identified and assigned to their respective token !

// To test all the lexemes and tokens used in my code I created two extra input files.

6. Testing input file called *list2.txt*.

```

list.txt  list2.txt  X
list2.txt
1  int main () {
2  int z ;
3  int b ;
4  z = 6_4+_2*_
5  b = 6_4+_
6  cout << z_b+_ ; }

```

-----> postfix expression in list2.txt

```

list.txt  list2.txt  list3.txt  output.txt X
output.txt
1
2  Lexemes      Tokens
3  -----
4  int          Variable_Literal
5  main         Function_Name
6  (            Open_parenth
7  )            Close_parenth
8  {            Open_bracket
9  int          Variable_Literal
10 z            Identifier
11 ;            Ending_Operator
12 int          Variable_Literal

```

```

list.txt  list2.txt  list3.txt X
list3.txt
1  double y ;
2  y = 9+_4*_2_
3  if ( y == 8 ) {
4      cout << y ;
5      else {
6          y = 5_2+_
7      }
8  }

```

```

27  0            int_Literal
28  -            Delimiter_op
29  4            Int_Literal
30  -            Delimiter_op
31  +            id_Operation
32  -            Delimiter_op
33  cout         Output_Stream
34  <<          Output_Operator
35  -            Delimiter_op
36  -            Delimiter_op
37  +            id_Operation
38  -            Delimiter_op
39  ;            Ending_Operator
40  }            Close_bracket

```

-----> the code recognizes the lexemes and assigns the appropriate tokens.
Testing input file called list3.txt.

-----> postfix expression in list3.txt

---> uses c++ keywords


```

1 -----
2 | Lexemes | Tokens |
3 -----
4 double | Variable_Literal
5 y | identifier
6 ; | Ending_Operator
7 y | identifier
8 = | Assign_op
9 9 | Int_Literal
10 _ | Delimiter_op
11 + | id_Operation
12 _ | Delimiter_op
13 4 | Int_Literal
14 _ | Delimiter_op
15 * | id_Operation
16 _ | Delimiter_op
17 2 | Int_Literal
18 _ | Delimiter_op
19 if | C++_Keywords
20 { | Open_parenth
21 y | identifier
22 == | Equalto_Operator
23 = | Assign_op
24 = | Assign_op
25 0 | Int_Literal
26 ) | Close_parenth
27 { | Open_bracket
28 cout | Output_Stream
29 << | Output_Operator
30 y | identifier
31 ; | Ending_Operator
32 else | C++_Keywords
33 { | Open_bracket
34 y | identifier
35 = | Assign_op
36 5 | Int_Literal
37 _ | Delimiter_op
38 2 | Int_Literal
39 _ | Delimiter_op
40 + | id_Operation
41 _ | Delimiter_op
42 } | Close_bracket
43 } | Close_bracket
44

```

----> the code recognizes the lexemes and assigns the appropriate token!

Conclusion/ Plans: I plan on updating my lexical analyzer so that it can analyze more complex postfix expressions. I plan on making it that it can recognize more than single digit numbers, and variables with more than one letter. I want to create an error routine that recognizes an invalid postfix expression. Another error message I want it to include is “_ delimiter missing_” if there is no delimiters or the proper amount in the postfix expression.

In addition, I want it to be able to recognize a larger subset of the c++ language and more keywords and maybe even recognize some of the C language. In conclusion, I really enjoyed this project and it gave me a lot of more insight as to how a lexical analyzer works.