## Experiment 03: Lexical Analyzer

**Learning Objective**: Students should be able to design a handwritten lexical analyzer.

**Tools:** Jdk1.8, Turbo C/C++, Python, Notepad++

**Theory:**

### Design of lexical analyzer
. Allow white spaces, numbers, and arithmetic operators in an expression
. Return tokens and attributes to the syntax analyzer
. A global variable token Val is set to the value of the number
. Design requires that
    - A finite set of tokens be defined
    - Describe strings belonging to each token

### Regular Expressions
• We use regular expressions to describe tokens of a programming language. •
A regular expression is built up of simpler regular expressions (using defining rules)
• Each regular expression denotes a language.
• A language denoted by a regular expression is called as a **regular set**.

### Regular Expressions (Rules)
Regular expressions over alphabet S
Regular Expression Language it denotes
$\varepsilon$ { $\varepsilon$ }
a$\in \Sigma$ S {a}
(r1) | (r2) L(r1) È L(r2)
(r1) (r2) L(r1) L(r2)
(r)* (L(r))*
(r) L(r)
• (r)+ = (r)(r)*
• (r)? = (r) | $\varepsilon$
• We may remove parentheses by using precedence rules.
    * highest
    concatenation next
    | lowest

### How to recognize tokens
Construct an analyzer that will return <token, attribute> pairs
We now consider the following grammar and try to construct an analyzer that will return

**<token, attribute> pairs.**
**relop < | = | = | <> | = | >**
**id letter (letter | digit)\***
**num digit+ ('.' digit+)? (E ('+' | '-')? digit+)?**
**delim blank | tab | newline**
**ws delim+**

Using set of rules as given in the example above we would be able to recognize the tokens. Given a regular expression R and input string x, we have two methods for determining whether x is in L(R). One approach is to use algorithm to construct an NFA N from R, and the other approach is using a DFA.

**Finite Automata**

• A *recognizer* for a language is a program that takes a string x, and answers "yes" if x is a sentence  of that language, and "no" otherwise.

      – We call the recognizer of the tokens as a *finite automaton*.

• A finite automaton can be: *deterministic(DFA)* or *non-deterministic (NFA)* • This means that we may use a deterministic or non-deterministic automaton as a lexical analyzer. • Both deterministic and non-deterministic finite automaton recognizes regular sets. • Which one?

      – deterministic – faster recognizer, but it may take more space

      – non-deterministic – slower, but it may take less space

      – Deterministic automatons are widely used lexical analyzers.

• First, we define regular expressions for tokens; Then we convert them into a DFA to get a lexical analyzer for our tokens.

**Algorithm1:** Regular Expression NFA DFA (two steps: first to NFA, then to DFA) **Algorithm2:** Regular Expression DFA (directly convert a regular expression into a DFA) **Converting Regular Expressions to NFAs**
 • Create transition diagram or transition table i.e. NFA for every expression  • Create a zero state as start state and with an e-transition connect all the NFAs and prepare a  combined NFA.

## Algorithm: for lexical analysis

      1) Specify the grammar with the help of regular expression

      2) Create transition table for combined NFA

      3) read input character

      4) Search the NFA for the input sequence.

      5) On finding accepting state

          i. if token is id or num search the symbol table

              1. if symbol found return symbol id

              2. else enter the symbol in symbol table and return its id.

          ii. Else return token

      6) Repeat steps 3 to 5 for all input characters.

**Input:**

          #include<stdio.h>

```
void main()
{
inta,b;
printf("Hello");
getch();
}
```

```
Preprocessor Directives: #include
Header File: stdio.h
Keyword : void main intgetch
Symbol: <> , ; ( ) ; }
Message: Hello
```

**Application:** To design a lexical analyzer.

**Design:**

```cpp
#include <iostream>
#include <fstream>
#include <vector>
using namespace std;

void helper(string line, vector<int>& ans,
vector<vector<string>>& map) {
 string temp = "";
 for (int i = 0; i < line.size(); i++) {
 if (line[i] == ' ' && temp.size() > 0) {
 bool found = false;
 // keyword
 for (int k = 0; k < map[0].size(); k++) {  if
(temp == map[0][k]) {
 ans[0]++;
 found = true;
 break;
 }
 }
 // identifier
 if (!found && ((temp[0] >= 'a' && temp[0] <=  'z') ||
(temp[0] >= 'A' && temp[0] <= 'Z'))) {  ans[4]++;
 found = true;
 }
 // constant
 if (!found && (temp[0] >= '0' && temp[0] <=  '9')) {
ans[5]++;
 found = true;
 }
 temp = "";
 continue;
 }
 // check for delimiter, operator
```

```cpp
 if (!found) {
 for (int k = 0; k < map[2].size(); k++) {  if
(ch == map[2][k]) {
 ans[2]++;
 found = true;
 break;
 }
 }
 }
 if (!found) temp.push_back(line[i]);
 }
}
int main() {
 fstream input;
 input.open("input.txt", ios::in);

 vector<int> ans(6, 0);
 vector<vector<string>> map = {
    {"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"},
 {"+", "-", "/", "*", "!", "%", "&", "&&", "^", "^^",
"|", "||", "==", "=", ">", "<", ">=", "<="},
 {"(", ")", "[", "]", "{", "}", ",", "."},
 };

 if (input.is_open()) {
 string line;

 while (getline(input, line)) {
 helper(line, ans, map);
 }
```

| | |
|---|---|
| ```
    bool found = false;

  string ch;
  ch.push_back(line[i]);
   for (int k = 0; k < map[1].size(); k++) {
 if (ch == map[1][k]) {
 ans[1]++;
 found = true;
 break;
 }
  }
``` | ```
  input.close();
  }

  cout << "keyword : " << ans[0] << endl;
  cout << "operator : " << ans[1] << endl;
  cout << "delimiter : " << ans[2] << endl;
   cout << "identifier : " << ans[4] << endl;
  cout << "constant : " << ans[5] << endl;

  return 0;
  }
``` |

## Output:



```
int main() {
   int n = 10;
   return 1;
}
```
```
keyword : 3
operator : 1
delimiter : 4
identifier : 2
constant : 0
```

## Result and Discussion:


## Learning Outcomes: The student should have the ability to

LO1: Appreciate the role of lexical analyzer in compiler design

LO2: Define role of lexical analyzer.

**Course Outcomes**: Upon completion of the course students will be able to Illustrate the working of the compiler and handwritten /automatic lexical analyzer.

## Conclusion:


For Faculty Use

| Correction Parameters | Formative Assessment [40%] | Timely completion of Practical [ 40%] | Attendance / Learning Attitude [20%] | |
|---|---|---|---|---|
| Marks Obtained | | | | |