

### Experiment 04 : Lexical Analyzer

#### 12 TE COMP A Naman Baranwal

**Learning Objective:** Student should be able to design handwritten lexical analyser.

**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 tokenval 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
--------------------	---------------------

$\epsilon$	$\{ \epsilon \}$
$a \in \Sigma$	$S \{ a \}$
$(r1)   (r2)$	$L(r1) \cup L(r2)$
$(r1) (r2)$	$L(r1) L(r2)$
$(r)^*$	$(L(r))^*$
$(r)$	$L(r)$

- $(r)^+ = (r)(r)^* \cdot (r)^? = (r) | \epsilon$
- 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 automata 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  $\rightarrow$  NFA  $\rightarrow$  DFA (two steps: first to NFA, then to DFA)

**Algorithm2:** Regular Expression  $\rightarrow$  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() {
    int a,b;
    printf("Hello");
    getch();
}
```

**Output:**

Preprocessor Directives: #include

Header File: stdio.h

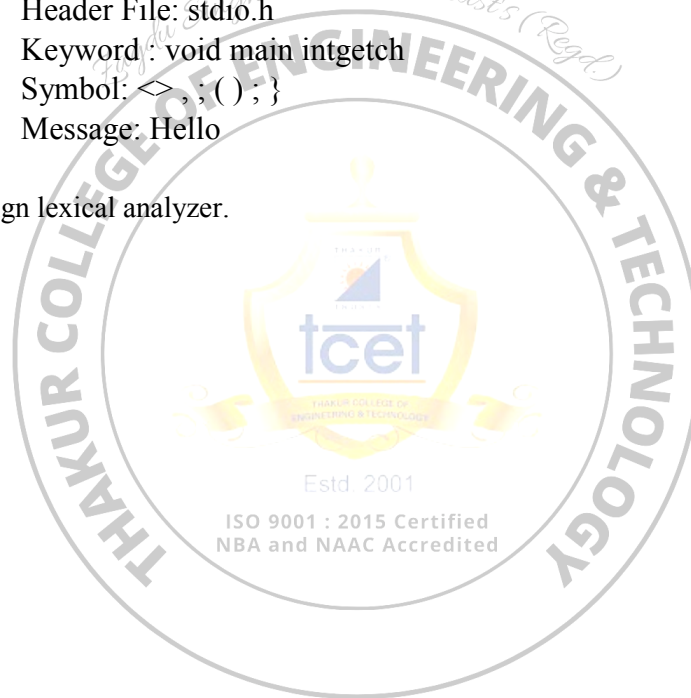
Keyword : void main int getch

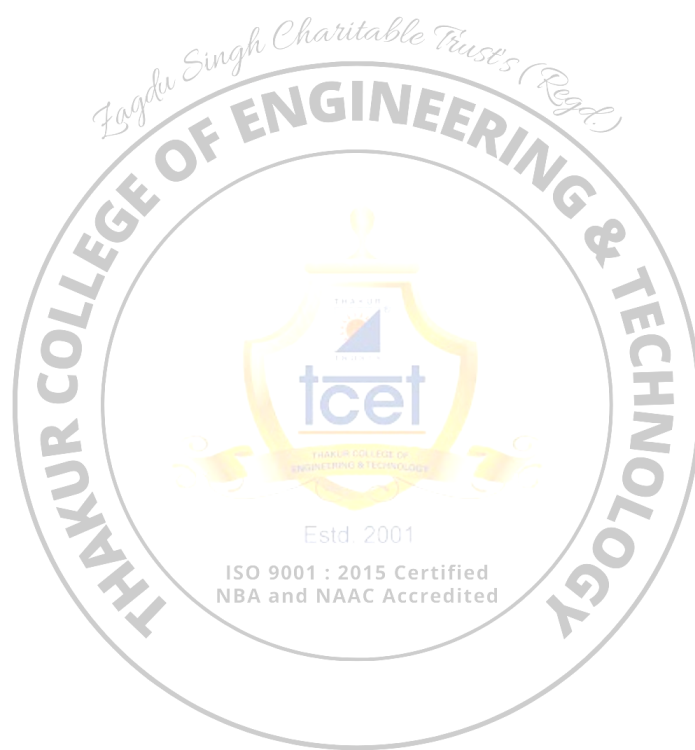
Symbol: <> , ; ( ) ; }

Message: Hello

**Application:** To design lexical analyzer.

**Design:**





### **Result and Discussion:**

**The code for lexical analyzer was written and executed and the results were seen successfully.**

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

LO1: Appreciate the role of lexical analyzer in compiler design

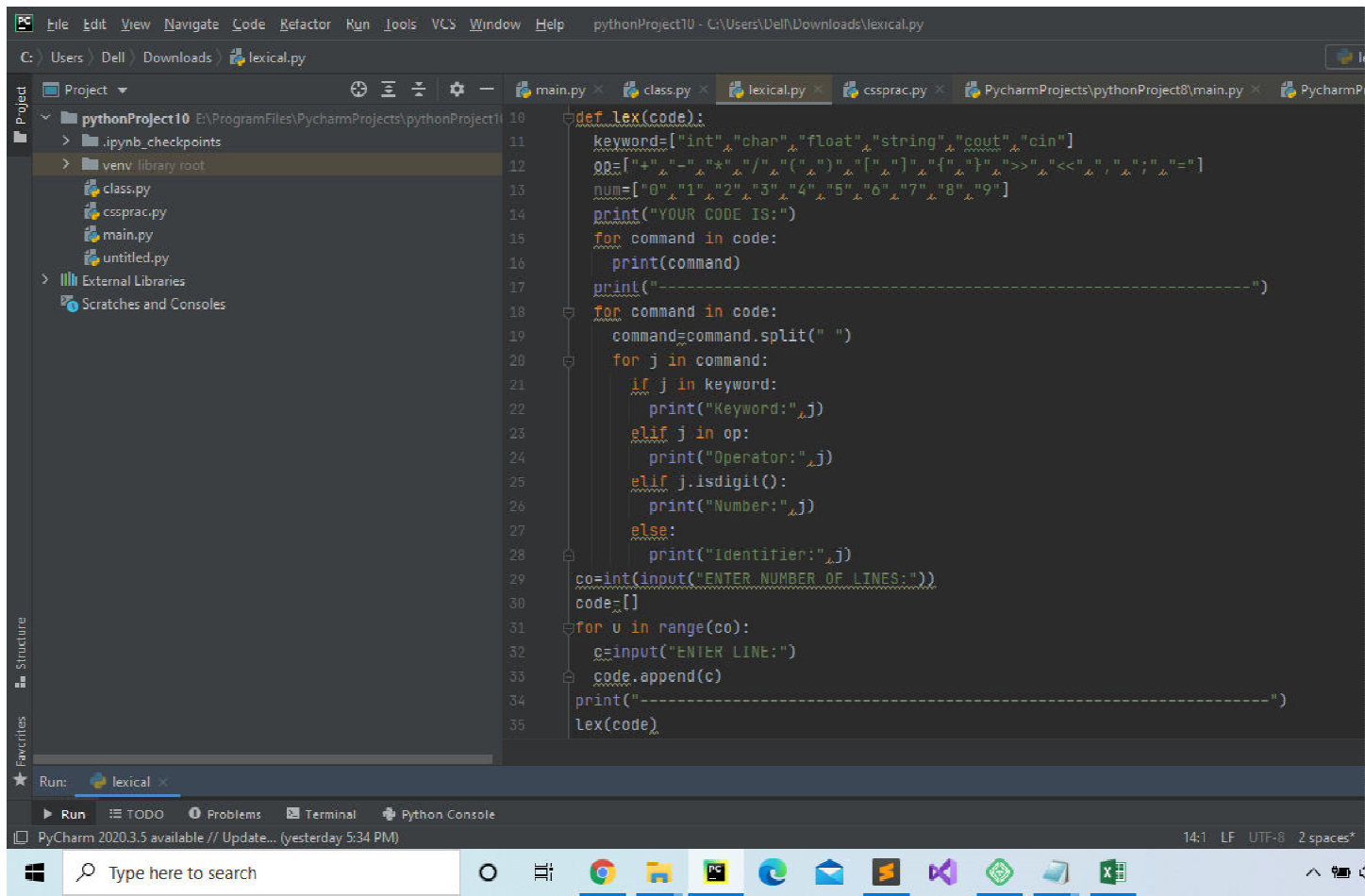
LO2: Define role of lexical analyzer.

**Course Outcomes:** Upon completion of the course students will be able to design handwritten lexical analyzer using HL programming language.

**Conclusion: The code for lexical analyzer was written and executed and the results were seen successfully.**

For Faculty Use

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



PC File Edit View Navigate Code Refactor Run Tools VCS Window Help pythonProject10 C:\Users\Dell\Downlo

C:\Users> Dell> Downloads> lexical.py

Project

- pythonProject10 E:\ProgramFiles\PycharmProjects\pythonProject10
  - .ipynb\_checkpoints
  - venv library root
  - class.py
  - cssprac.py
  - main.py
  - untitled.py
- External Libraries
- Scratches and Consoles

main.py class.py lexical.py

```
9
10 def lex(code):
11     keyword=["int","char","float"
12     op=["+","-","*","/","(","")"
13     num=["0","1","2","3","4","5"
14     print("YOUR CODE IS:")
15     for command in code:
16         print(command)
17     print("-----")
18     for command in code:
19         command=command.split(" ")
```

Run: lexical

E:\ProgramFiles\PycharmProjects\pythonProject10\venv\Scripts\python.exe C:/User  
ENTER NUMBER OF LINES:2  
ENTER LINE:2+8  
ENTER LINE:4\*7  
-----  
YOUR CODE IS:  
2+8  
4\*7  
-----  
Identifier: 2+8  
Identifier: 4\*7  
  
Process finished with exit code 0  
|

Run TODO Problems Terminal Python Console

PyCharm 2020.3.5 available // Update... (yesterday 5:34 PM)

Type here to search