

**MINISTRY OF EDUCATION AND RESEARCH OF REPUBLIC OF MOLDOVA**  
**TECHNICAL UNIVERSITY OF MOLDOVA**  
**FACULTY OF COMPUTERS, INFORMATICS AND MICROELECTRONICS**  
**DEPARTMENT OF SOFTWARE ENGINEERING AND AUTOMATICS**

## **Laboratory Work 3:**

### **Lexer & Scanner.**

---

**Course: Formal Languages & Finite Automata**

**Student: Gusev Roman**

**Group: FAF-222**

**Professor: Cojuhari Irina**

**University Assistant: Crețu Dumitru**

**Chișinău, 2024**

## Content

<b>Theory</b> . . . . .	<b>3</b>
<b>Objectives</b> . . . . .	<b>3</b>
<b>Implementation</b> . . . . .	<b>3</b>
<b>Conclusions</b> . . . . .	<b>22</b>
<b>Bibliography</b> . . . . .	<b>28</b>

## Theory

- **Definitions:**

- **Lexer** - The lexical analyzer defines how the contents of a file are broken into tokens, which is the basis for supporting custom language features [1];
- **Token** - A lexical token is a string with an assigned and thus identified meaning.
- **Lexeme** - A lexeme is only a string of characters known to be of a certain kind.

## Objectives

- Understand what lexical analysis [2] is;
- Get familiar with the inner workings of a lexer/scanner/tokenizer;
- Implement a sample lexer and show how it works.

## Implementation Description

- For the start, I decided to implement a Lexer for the PBL DSL that me and the team I work with decided to develop. During the work on this Laboratory Work, I decided to implement an example Lexer that I found in some guides on YouTube [3]. Specifically, I followed the structure of the presented in the video Lexer, but for the extension I used only my knowledge.
- The Lexer I developed pass through the entire Input symbol by symbol until the End of File, and transform each combination of the Input into Lexemes. The mathematical expressions are not required to be separated by Whitespace and may be written in both, for example "2+2" and "2 + 2", forms, and will be tokenized correctly. For the "words" type Lexemes, they are required to be separated by whitespace. This is done in order to greedily tokenize them and eliminate cases in inputs such as "Formula1", that may be tokenized like:

```
1 FORMULA_TOKEN : Formula
2 NUMBER: 1
```

and tokenize that as a single word, that might give theoretically output as:

```
1 IDENTIFIER: Formula1
```

This was done in order to have a defined structure for the Tokens, that will lead to a whitespace-sensitive Language.

- First thing developed was the Tokens themselves. I have several Tokens that I decided to hold in Enum classes, and that have been used in the Lexer to classify them.

```

1  from enum import Enum
2
3  class OperatorType(Enum):
4      MULTIPLY = "*"
5      DIVISION = "/"
6      ADDITION = "+"
7      SUBTRACTION = "-"
8      POWER = "**"
9
10 class ComparatorType(Enum):
11     EQUALS = "=="
12     NOT_EQUALS = "!="
13     LESS_THAN = "<"
14     GREATER_THAN = ">"
15     LESS_THAN_OR_EQUALS = "<="
16     GREATER_THAN_OR_EQUALS = ">="
17
18 class FileType(Enum):
19     CSV = "csv"
20     TEXT = "txt"
21     JSON = "json"
22     EXCEL = "excel"
23     CONSOLE = "console"
24
25 class ImageType(Enum):
26     JPG = "jpg"
27     PNG = "png"
28
29 class PlotType(Enum):
30     GRAPH = "graph"
31     BAR = "bar"
32     PIE = "pie"
33     PLOT = "plot"

```

```

34     HIST = "hist"
35
36 class VariableType(Enum):
37     FORMULA = "Formula"
38     DATA = "Data"
39     DATASET = "dataset"
40     NAME = "name"
41
42 class ExportToType(Enum):
43     EXPORT_TO_IMAGE = "ExportToImage"
44     EXPORT_TO_FILE = "ExportToFile"
45
46 class VisualizationType(Enum):
47     VISUALIZE_DATA = "VisualData"
48     VISUALIZE_FORMULA = "VisualFormula"

```

- At the same time, I started with the Token class, that will describe what type of Lexeme I have in input. Here I defined the Constructor for the Token, where it gets the type/kind of the Token (type is a reserved keyword in python therefore I used kind), and the line of the input - whole input of the User. Also, it gets the string, that is the Token itself. But there may be a case when the String returned is empty, so it is replaced with "EOF" - End of File Token:

```

1 class Token:
2     def __init__(self, kind, line):
3         self.kind = kind
4         self.line = line
5         self.locale, self.string = line.new_locale()
6
7         # If retrieved string is empty, then set it as "EOF",
8         # otherwise - stays the same retrieved string
9         self.string = self.string or "EOF"
10
11     ...

```

- After that, Token class has 2 other methods in it - "\_\_repr\_\_" that will be the representation of the Token in a visual human-readable form for data structures such as Lists, because they call specifically

this method. Tokens will be present in a list, therefore this method is required.

- At the same time, 2-nd method - "mark", will be used to make the visual representation of the Incorrect Tokens during the evaluation. This method will be explained in depth further.

```
1 class Token:
2     ...
3     def __repr__(self):
4         return f"{self.kind}: '{self.string}'"
5
6     def mark(self):
7         self.line.mark(self)
```

- Next Class developed for the Lexer implementation was "SourceLine", that is used to hold the input from the user and will hold also a lot of useful functionalities that will be called by the Lexer class and will track the Tokens based on their index in the input. First of all, I developed the Constructor, that will hold the Input from the user, will hold some locale variables, a list of 2 integers, that will be the position of the START of the Token and END of the Token, that was done because some Tokens are multi-symbol Tokens, and Marked, that will hold the wrong portion of the Token, in case that this Token is invalid:

```
1 class SourceLine:
2     def __init__(self, line):
3         self.line = line # Actual line from the User Input
4         self.locale = [0, 0] # Track START and END of a Token
5         self.marked = [len(line), -1]
6         ...
```

- After that, I have a method to check if the Input from the User was processed entirely. This is done by the comparison check between the locale variable and the length of the input - if the SourceLine locale passed the length of the string input or is equal to this length, then the input was tokenized entirely:

```
1 class SourceLine:
2     ...
3     # Checks if line was processed entirely
4     def finished(self):
5         return self.locale[1] >= len(self.line)
```

6

...

- After that, I developed the method to actually take and adjust the state of the variables in order to progress in the Token components further. This method will peek the next symbol and will increment the locale variables - specifically the END position if the Input is not processed fully, otherwise will stay the same. Also, this method returns the symbol that was peeked:

```

1 class SourceLine:
2     ...
3     # Adjusts the states and takes one symbol at a time
4     def take(self):
5         symbol = self.next()
6         self.locale[1] += 0 if self.finished() else 1
7
8         return symbol
9     ...

```

- Here is presented a method to Ignore some Symbols from the Input, such as Whitespaces or others. This is done by setting the START Position/Index of the Token to END Position/Index - they will become the same and will be after the Ignored Symbols:

```

1 class SourceLine:
2     ...
3     # Ignores current Locale for example spaces
4     def ignore(self):
5         self.locale[0] = self.locale[1]
6     ...

```

- After that, I have a method to retrieve the entire Multi-Symbol Token. This is done by String Slicing. This method will Slice the Input String and will retrieve the Token between START and END Indices:

```

1 class SourceLine:
2     ...
3     # Takes MultiCharacter Symbol
4     def taken(self):
5         return self.line[self.locale[0]:self.locale[1]]
6     ...

```

- 
- Here I present the method that will update the Locale variables in order to assign it correctly those START and END Positions for the new Token that will be created. This method was featured in the Constructor for the Token class. This is done by the assignment of the Locale variables and the actual Token string when the Token is constructed itself. After that the locales that are in the class SourceLine are reset to the END of the current Token:

```
1 class SourceLine:
2     ...
3     def new_locale(self):
4         locale, taken = self.locale.copy(), self.taken()
5         self.locale[0] = self.locale[1]
6         return locale, taken
7     ...
```

- Here I have 2 methods that are related to Error handling, i.e., will show where the error in the String is. Given the Token, that has its own Locale variables, it will compare them with the marked variables. Marked are initially set outside the Token range, and for the START Position of the mistake Symbol, it will take the minimum value between START of the Token and the outside bound of the marked variable, and for the END - it is the maximum value. After that, it will append to a string called "marks", that will hold the marks, and in a for loop it will add the "^\n" marks that will indicate where the mistake is between the boundaries of the marked values:

```
1 class SourceLine:
2     ...
3     def mark(self, token):
4         self.marked[0] = min(token.locale[0], self.marked[0])
5         self.marked[1] = max(token.locale[1], self.marked[1])
6
7     # Creates Error Marks
8     def get_marks(self):
9         marks = ""
10
11         for i in range(len(self.line) + 1):
12             between = self.marked[0] <= i < self.marked[1]
```



```

13         marks += "^" if between or self.marked[0] == i else " "
14
15     return marks

```

- After that, I have developed a new class for Error handling, i.e., will be the class for the actual `LanguageError`, of type `RuntimeError`, that will get the marks of the Token where error was obtained and will output them, with the actual name of the Error that will be assigned when subclasses of that `LanguageError` class will be created:

```

1 class LanguageError(RuntimeError):
2     def __init__(self, component, message):
3         component.mark()
4         self.line = component.line
5         self.message = message
6
7     def __str__(self):
8         return f"{self.line.get_marks()}\n{type(self).__name__} : {
          self.message}"

```

- For the actual Lexer, I created a new python file with the same name - `Lexer`. First of all, I imported all the previously created files and also created a class for the `LexerError`, that is a subclass of the `LanguageError`, that will be used to create more specific Errors further:

```

1 from Error import LanguageError
2 from Token import Token
3 from Tokens import FileType, VariableType, ExportToType, ImageType,
   PlotType, VisualizationType
4
5 class LexerError(LanguageError):
6     pass
7
8 ...

```

- After that, I created the class for the Lexer, that has a Constructor in the beginning, that will hold the actual Input of the User, that was converted to the `SourceLine` class - this will provide the extensive use of the methods that were developed in that class:

```

1 ...

```

```

2 class Lexer:
3     def __init__(self):
4         self.line = None
5     ...

```

- Here is another method from Lexer class, that returns a new Token with specified Kind/Type and with the input String from the User. This method will be used the most amongst other methods, because by this method, new Tokens will be created from Lexemes, based on some criteria of matching patterns further when analyzing the Input String:

```

1 ...
2 class Lexer:
3     ...
4     def new_token(self, kind):
5         return Token(kind, self.line)
6     ...

```

- Here is another method from Lexer class, that is called once to start the process of Tokenization. First, I take the Input String and assign it to the stored Input String variable that is in the scope of a Lexer object. Afterward, I create a list of Tokens, that will store all the Tokens. In a while loop, I iterate over the Input String until it is finished, i.e., until the End of File. Also, I call the method to ignore the Whitespaces. At the end of each iteration, the list of tokens is appended with new Token that. And at the end, I add the last Token - "EOF":

```

1 ...
2 class Lexer:
3     ...
4     def make_tokens(self, line):
5         self.line = line
6         tokens = []
7
8         while not self.line.finished():
9             self.ignore_spaces()
10            if self.line.finished():
11                break

```

```

12         tokens += [self.make_token()]
13
14
15         tokens += [self.new_token("EOF")]
16         return tokens
17     ...

```

- Here is the method to ignore Whitespaces. It is an iteration based method, that will go over the Input String until it encounters a non-whitespace character. In the loop, will be called 2 methods - to take new element (will update the variables in the scope of the SourceLine object) and will ignore each whitespace:

```

1     ...
2     class Lexer:
3         ...
4         def ignore_spaces(self):
5             while self.line.next().isspace():
6                 self.line.take()
7                 self.line.ignore()
8         ...

```

- Here is the method to make Tokens. Each next Symbol from the Input will be checked to match any of the Token Type. Specifically, in this code snippet is shown how Non-Word Lexemes are matched to a Token Type, such as: LPAREN or RPAREN Tokens (parenthesis Tokens) or NUMBER Tokens and so on. They all call their own specific method of Tokenization:

```

1     ...
2     class Lexer:
3         ...
4         def make_token(self):
5             if self.line.next() in "()":
6                 return self.make_parenthesis()
7
8             if self.line.next() in "[]":
9                 return self.make_brackets()

```

```

10
11     if self.line.next() == ";":
12         return self.make_semicolon()
13
14     if self.line.next() == ",":
15         return self.make_colon()
16
17     if self.line.next() in "<=>":
18         return self.make_comparison()
19
20     if self.line.next() in "|_^":
21         return self.make_punctuator()
22
23     if self.line.next() in "~+-=\%":
24         return self.make_operator()
25
26     if self.line.next() in "0123456789.":
27         return self.make_number()
28
29     if self.line.next() in "\#":
30         return self.make_comment_line()
31
32     if self.line.next() in "/":
33         return self.make_comment_block()
34
35     if self.line.next() == "'":
36         return self.make_path()
37     ...

```

- Here is the method to make Tokens. Each next Symbol from the Input will be checked to match any of the Token Type. Specifically, in this code snippet is shown how Word Lexemes are matched to a Token Type, such as: FORMULA Token (keyword - Formula) or VISULIZATION\_TYPE Token (VisualData or VisualFormula). They all call their own specific method of Tokenization. In the end, if a Lexeme is not a Token, then raise a LexerError, with "Unrecognized Symbol!" message:

```

1  ...
2  class Lexer:
3      ...
4      def make_token(self):
5          if self.line.next().isalpha():
6              if self.line.next() == "F":
7                  return self.make_keyword_token("Formula")
8              elif self.line.next() == "D":
9                  return self.make_keyword_token("Data")
10             elif self.line.next() == "R":
11                 return self.make_read_from()
12             elif self.line.next() == "E":
13                 return self.make_export_to()
14             elif self.line.next() == "V":
15                 return self.make_visualization()
16             else:
17                 return self.make_ID()
18
19         self.line.take()
20         raise LexerError(self.new_token("?"), "Unrecognized Symbol!")
21     ...

```

- Next, I will list some of the methods of Tokenization, not all, because essentially they are pretty similar between each other, and the main difference is that some of them perform more checks and takes more symbols.
- Here I provided the Tokenization of VISUALIZATION\_TYPE Tokens. First of all, I iterate over the Lexeme until it is finished, and check if it is in the VisualizationType Enum class. If it is, then create a new Token of this Type. Otherwise - raise a LexerError with the message that maybe User meant one of 2 available VisualizatioType Tokens:

```

1  ...
2  class Lexer:
3      ...

```

```

4     def make_visualization(self):
5         while self.line.next().isalnum() and not self.line.finished
6             ():
7             self.line.take()
8
9         if self.line.taken() in VisualizationType:
10             return self.new_token("VISUALIZATION_TYPE")
11
12         raise LexerError(self.new_token("VISUALIZATION_TYPE"), "
13             Maybe you meant 'VisualData' or 'VisualFormula' instead?
14             ")

```

- Here I provided the Tokenization of COMPARISON\_OPERATORS Tokens. First of all, I take the first character of the Lexeme and check if the next element is an Equal Symbol, then it might be ">=" "<=" or "==" Comparison Operator, and create a new COMPARISON\_OPERATOR Token. Otherwise, I check if the current symbol is one of ">" or "<", and if it is not, then I create a ASSIGN\_OPERATOR Token, because the symbol is "=", in the other case, then it is still a COMPARISON\_OPERATOR Token:

```

1     ...
2     class Lexer:
3         ...
4         def make_comparison(self):
5             operator = self.line.take() # Take the first character of
6             the comparison operator
7
8             # Check if the next character is "=" to determine if it's a
9             complete comparison operator
10            if self.line.next() == "=":
11                operator += self.line.take()
12                return self.new_token("COMPARISON_OPERATOR")
13
14            # If the next character is not "=", treat the standalone
15            character as an assignment operator

```

```

13         if operator == ">" or operator == "<":
14             return self.new_token("COMPARISON_OPERATOR")
15         else:
16             return self.new_token("ASSIGN_OPERATOR")
17         ...

```

- Here I provided the Tokenization of COMMENT\_BLOCK Tokens. In the Language I had, there are two Commenting possibilities:

1. **COMMENT\_LINE**: using `"`, that will comment the symbols after this symbol,
2. **COMMENT\_BLOCK**: using `"/"` and `"/"` that will comment all from inside.

First of all, I take the first character of the Lexeme and check if the next one is `"*` and then take all the elements until it finds `"/"` symbols. In case that the comment block was not closed, it will raise an error with the message that says that the comment block should be enclosed.

```

1     ...
2     class Lexer:
3         ...
4         def make_comment_block(self):
5             self.line.take()
6             if self.line.next() == "*":
7                 self.line.take()
8                 while True:
9                     if self.line.next() == "*":
10                        self.line.take()
11                        if self.line.next() == "/":
12                            self.line.take()
13                            return self.new_token("COMMENT_BLOCK")
14                        elif self.line.finished():
15                            raise LexerError(self.new_token("COMMENT_BLOCK"),
16                                                "Block comment not terminated")
17                        else:
18                            self.line.take()
19             else:
20                 return self.make_operator()

```

- Here I provided the Tokenization of Keyword Tokens. In the Language I had, there are multiple keywords:

1. **FILE\_TYPE**: csv, txt, json, excel, console,
2. **IMAGE\_TYPE**: png, jpg,
3. **VARIABLE\_TYPE**: Data, Formula, dataset, name,
4. **PLOT\_TYPE**: graph, pie, bar, hist, plot,
5. **ID**: any string that starts with a letter.

First of all, I take the whole Lexeme, and check if they are in one of the mentioned Token Types Enum classes. If they are not, then it is an ID Token = Identifier.

```

1  ...
2  class Lexer:
3      ...
4      def make_ID(self):
5          while self.line.next().isalnum() and not self.line.finished
6              ():
7              self.line.take()
8
9          if self.line.taken().lower() in FileType:
10             return self.new_token("FILE_TYPE")
11
12         if self.line.taken().lower() in ImageType:
13             return self.new_token("IMAGE_TYPE")
14
15         if self.line.taken().lower() in PlotType:
16             return self.new_token("PLOT_TYPE")
17
18         if self.line.taken() in VariableType:
19             return self.new_token("VARIABLE_TYPE")
20
21         if self.line.taken() == "range":
22             return self.new_token("RANGE")

```



22  
23  
24

```
return self.new_token("ID")
```

```
...
```

- There are more Tokens and, therefore, Tokenization Methods, but I listed some of the most important and most unique ones, because there are also methods that are only 2 lines of code length and are similar with the ones above.
- In the Main class, I decided to give User the opportunity to choose between 3 methods of Input:
  1. Manual line by line (will check each line instantly):

```
1  ...
2  def main():
3      ...
4      elif type_input.lower() == "l":
5          print("Enter lines of code. Type 'exit' to finish.")
6          while True:
7              line = input("> ")
8              if line == "exit":
9                  break
10
11             try:
12                 source_line = SourceLine(line)
13                 tokens = lexer.make_tokens(source_line)
14
15                 print("TOKENS:")
16                 for token in tokens:
17                     print(token)
18
19                 json_object = json.dumps(convert(tokens), indent
20                                         =4)
21
22                 with open("./Laboratory-Work-3-Lexer/
23                           ExamplePrograms/Tokenized_Manual_Input.json",
24                           "w") as outfile:
```

```

22         outfile.write(json_object)
23     except LanguageError as error:
24         print(error)
25     ...

```

2. Manual all lines at once (will check the lines altogether):

```

1     ...
2 def main():
3     ...
4     if type_input.lower() == "c":
5         lines = ""
6         print("Enter lines of code. Type 'exit' to finish.")
7         while True:
8             line = input("> ")
9             if line == "exit":
10                break
11
12            lines += line + "\n"
13
14        try:
15            source_line = SourceLine(lines)
16            tokens = lexer.make_tokens(source_line)
17
18            print("TOKENS:")
19            for token in tokens:
20                print(token)
21
22            json_object = json.dumps(convert(tokens), indent=4)
23
24            with open("./Laboratory-Work-3-Lexer/ExamplePrograms/
25                Tokenized_Manual_Input.json", "w") as outfile:
26                outfile.write(json_object)
27        except LanguageError as error:

```

```

27         print(error)
28     ...

```

3. File (will tokenize contents of a .txt file). Additionally, here I decided to export the results of the Tokenization in a json file that will hold every Token by their Type and information of the placement - Index by the variable of Locales:

```

1     ...
2     def main():
3         ...
4         elif type_input.lower() == "f":
5             file_path = select_file()
6
7             if not file_path:
8                 print("No file selected.")
9                 return
10
11             with open(file_path) as f:
12                 lines = f.read()
13
14             try:
15                 source_line = SourceLine(lines)
16                 tokens = lexer.make_tokens(source_line)
17
18                 print("TOKENS:")
19                 for token in tokens:
20                     print(token)
21
22                 json_object = json.dumps(convert(tokens), indent=4)
23
24                 with open(
25                     f"./Laboratory-Work-3-Lexer/ExamplePrograms/
                        Tokenized_{os.path.splitext(os.path.
                        basename(file_path))[0]}.json",

```

```

26         "w") as outfile:
27             outfile.write(json_object)
28     except LanguageError as error:
29         print(error)
30     ...

```

- Also, by the use of Tkinter libray, I provide User to choossem manually a .txt File from the Explorer in order to tokenize its insides:

```

1  ...
2  def select_file():
3      root = tk.Tk()
4      root.withdraw()
5      root.attributes("-topmost", True)
6
7      file_path = filedialog.askopenfilename(initialdir="./
8      Laboratory_Work_3-Lexer/ExamplePrograms",
9                                              title="Select a file")
10
11     return file_path
12 ...

```

- As I mentioned, I made a conversion into a JSON File in order to display the results of the Tokenization. In this method I iterate over the list of Tokens, and separate them by the ": " symbol that delimits the Lexeme from the Type. Additionally, I create a Dictionary that will hold the following structure:

```

1  {
2  Type :
3      {
4      Lexeme :
5          [
6          START_POS,
7          END_POS
8          ]
9      }

```

10

```
}
```

- Here is the method:

1

```
...
```

2

```
def convert(lst):
```

3

```
    res_dict = {}
```

4

```
    for i in range(0, len(lst)):
```

5

```
        key = lst[i].kind
```

6

```
        value = lst[i].string
```

7

```
        locale = lst[i].locale
```

8

```
        if key in res_dict:
```

9

```
            res_dict[key].update({value: locale})
```

10

```
        else:
```

11

```
            res_dict[key] = {value: locale}
```

12

```
    return res_dict
```

## Conclusions / Results

I present here one output for the Laboratory Work nr.3.

First part of the console output is the general information about the laboratory work, student and group:

```
1 Laboratory Work 3 - Lexer
2 Student: Gusev Roman
3 Group: FAF-222
```

After that, user is asked to choose the method of input they want:

```
1 MY LEXER:
2 Choose the type of the input:
3 F - File
4 C - Console all lines together
5 L - Console line-by-line
```

After that is the prompt for user:

```
1 Your choice: ____
```

I will provide example of inputs for each type of Input (Each Line from the console starts with ">" symbol to show where to start):

- Manual line-by-line:

```
1 Enter lines of code. Type 'exit' to finish.
2 > Data data = ReadFrom("/path1/file1")
3 TOKENS:
4 VARIABLE_TYPE: 'Data'
5 ID: 'data'
6 ASSIGN_OPERATOR: '='
7 READ_FROM: 'ReadFrom'
8 LPAREN: '('
9 PATH: '"/path1/file1"'
10 RPAREN: ')'
11 EOF: 'EOF'
```

```

12 > Formula formula 1
13 TOKENS:
14 VARIABLE_TYPE: 'Formula'
15 ID: 'formula'
16 NUMBER: '1'
17 EOF: 'EOF'

```

- Manual line-by-line:

```

1 Enter lines of code. Type 'exit' to finish.
2 > Data data = ReadFrom("/path1/file1")
3 TOKENS:
4 VARIABLE_TYPE: 'Data'
5 ID: 'data'
6 ASSIGN_OPERATOR: '='
7 READ_FROM: 'ReadFrom'
8 LPAREN: '('
9 PATH: '"/path1/file1"'
10 RPAREN: ')'
11 EOF: 'EOF'
12 > Formula formula 1
13 TOKENS:
14 VARIABLE_TYPE: 'Formula'
15 ID: 'formula'
16 NUMBER: '1'
17 EOF: 'EOF'

```

and the JSON Contents (Note that each input will modify the JSON Contents):

```

1 {
2   "VARIABLE_TYPE": {
3     "Formula": [
4       0,
5       7
6     ]

```

```

7      },
8      "ID": {
9          "formula": [
10             8,
11             15
12          ]
13      },
14      "NUMBER": {
15          "1": [
16             16,
17             17
18          ]
19      },
20      "EOF": {
21          "EOF": [
22             17,
23             17
24          ]
25      }
26  }

```

- Manual all lines at once:

```

1  Enter lines of code. Type 'exit' to finish.
2  > Formula 1
3  > Data csv
4  > # alpha
5  > exit
6  TOKENS:
7  VARIABLE_TYPE: 'Formula'
8  NUMBER: '1'
9  VARIABLE_TYPE: 'Data'
10 FILE_TYPE: 'csv'
11 COMMENT_LINE: '# alpha'

```



```
12 EOF: 'EOF'
```

and the JSON Contents:

```
1 {
2   "VARIABLE_TYPE": {
3     "Formula": [
4       0,
5       7
6     ],
7     "Data": [
8       10,
9       14
10    ]
11  },
12  "NUMBER": {
13    "1": [
14      8,
15      9
16    ]
17  },
18  "FILE_TYPE": {
19    "csv": [
20      15,
21      18
22    ]
23  },
24  "COMMENT_LINE": {
25    "# alpha": [
26      19,
27      26
28    ]
29  },
30  "EOF": {
```

```

31         "EOF": [
32             27,
33             27
34         ]
35     }
36 }

```

- File input:

- Here are the contents of the file:

```

1  # This is Table Data read from TXT File;
2  Data tableData = ReadFrom("/path1/folder1/file.txt");
3  VisualData(console) dataset=(tableData);
4  ExportToFile("/path1/") dataset = (tableData) name = (new.txt);
5  # Operators
6  1.11 / 1
7  1 // 2
8  1.1 * 1
9  1 ** 2
10 1 + 1
11 1 - 1
12 1 == 1
13 2 >= 1
14 3 <= 1
15 1 = 2
16 1 < 2
17 2 > 3
18
19 Formula formulaData = formula[x**2 + x*2 + y];
20 VisualFormula(formulaData) range=(1,2);
21 ExportToImage("/path2/") graph(formulaData) name=(newGraph.png);
22 /* Line 1
23 Data excelData = ReadFrom("/path2/folder1/file.xlsx");
24 Line 3 */

```

– Console Output:

```
1  TOKENS:
2  COMMENT_LINE: '# This is Table Data read from TXT File;'
3  VARIABLE_TYPE: 'Data'
4  ID: 'tableData'
5  ASSIGN_OPERATOR: '='
6  READ_FROM: 'ReadFrom'
7  LPAREN: '('
8  PATH: '/path1/folder1/file.txt'
9  RPAREN: ')'
10 SEMICOLON: ';'
11 VISUALIZATION_TYPE: 'VisualData'
12 LPAREN: '('
13 FILE_TYPE: 'console'
14 RPAREN: ')'
15 VARIABLE_TYPE: 'dataset'
16 ASSIGN_OPERATOR: '='
17 LPAREN: '('
18 ID: 'tableData'
19 RPAREN: ')'
20 SEMICOLON: ';'
21 EXPORT_TO_TYPE: 'ExportToFile'
22 ...
```

and the JSON Contents are presented in the link [4]:

As a conclusion to this Laboratory Work nr.3, I can say that I accomplished the given task, specifically:

- Implement your own Lexer.

Also, I managed to understand better the concept of Lexers, Tokenization, Tokens and Lexemes and their main difference. At the same time, I understood how any Compiler "reads" the code and separates each Lexeme it finds when it analyzes the written code.

## Bibliography

- [1] “Lexer and Parser Definition | IntelliJ Platform Plugin SDK.” n.d. IntelliJ Platform Plugin SDK Help. Accessed March 16, 2024. <https://plugins.jetbrains.com/docs/intellij/lexer-and-parser-definition.html>.
- [2] “Kaleidoscope: Kaleidoscope Introduction and the Lexer — LLVM 19.0.0git Documentation.” n.d. Llm.org. Accessed March 16, 2024. <https://llvm.org/docs/tutorial/MyFirstLanguageFrontend/LangImpl01.html>.
- [3] “How to Create Your Own Programming Language - Episode 2: The Lexer.” n.d. Www.youtube.com. Accessed March 17, 2024. <https://www.youtube.com/watch?v=N103OVKmDR4>.
- [4] “LFA-Laboratory-Works/Laboratory-Work-3-Lexer/ExamplePrograms/Tokenized\_example\_1.Json at Main · Ghenntoggy1/LFA-Laboratory-Works.” n.d. GitHub. Accessed March 17, 2024. [https://github.com/Ghenntoggy1/LFA-Laboratory-Works/blob/main/Laboratory-Work-3-Lexer/ExamplePrograms/Tokenized\\_example\\_1.json](https://github.com/Ghenntoggy1/LFA-Laboratory-Works/blob/main/Laboratory-Work-3-Lexer/ExamplePrograms/Tokenized_example_1.json).