

Laboratory Work Nr.3 Lexer Scanner

Cretu Cristian

March 18, 2024

Theory

The term lexer comes from lexical analysis which, in turn, represents the process of extracting lexical tokens from a string of characters. There are several alternative names for the mechanism called lexer, for example tokenizer or scanner. The lexical analysis is one of the first stages used in a compiler/interpreter when dealing with programming, markup or other types of languages. The tokens are identified based on some rules of the language and the products that the lexer gives are called lexemes. So basically the lexer is a stream of lexemes. Now in case it is not clear what's the difference between lexemes and tokens, there is a big one. The lexeme is just the byproduct of splitting based on delimiters, for example spaces, but the tokens give names or categories to each lexeme. So the tokens don't retain necessarily the actual value of the lexeme, but rather the type of it and maybe some metadata.

Objectives

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

Implementation description

Token Class

A class that represents a token in the language I try to tokenize. I overload the built in functions so that I can easily manipulate tokens.

```
1 class Token:
2     """
3     Represents a token in a lexer.
4
5     Attributes:
```

```

6         - type (str): The type of the token.
7         - value (str): The value of the token.
8     """
9
10    def __init__(self, type, value):
11        self.type = type
12        self.value = value
13
14    def __repr__(self):
15        return f"{self.type}"
16
17    def __str__(self):
18        return f"{self.type}"
19
20    def __eq__(self, other):
21        return self.type == other.type and self.value ==
other.value
22
23    def __ne__(self, other):
24        return not self.__eq__(other)
25
26    def __hash__(self):
27        return hash((self.type, self.value))
28
29    def __lt__(self, other):
30        return self.type < other.type and self.value < other
.value

```

Tokenizer Class

This class is responsible for tokenizing a line of code and matching it to tokens while preserving the identifiers in the code.

```

1 from Token import Token
2 import re
3
4
5 class Tokenizer:
6     """
7     A class that tokenizes a given line of code based on
8     predefined tokens.
9
10    Attributes:
11        - tokens (list): A list of Token objects
12        representing the predefined tokens.
13
14    Methods:
15        - tokenize(line): Tokenizes the given line of code
16        and returns a list of tokens found.

```

```

14         - get_tokens(): Returns the list of predefined
tokens.
15         - print_tokens(): Prints the list of predefined
tokens.
16         """
17
18     def __init__(self):
19         self.tokens = [
20             Token("int", r"\bint\b"),
21             Token("float", r"\bfloat\b"),
22             Token("string", r"\bstring\b"),
23             Token("bool", r"\bbool\b"),
24             Token("true", r"\btrue\b"),
25             Token("false", r"\bfalse\b"),
26             Token("if", r"\bif\b"),
27             Token("else", r"\belse\b"),
28             Token("while", r"\bwhile\b"),
29             Token("for", r"\bfor\b"),
30             Token("return", r"\breturn\b"),
31             Token("break", r"\bbreak\b"),
32             Token("continue", r"\bcontinue\b"),
33             Token("function", r"\bfun\b"),
34             Token("print", r"\bprint\b"),
35             Token("lparen", r"("),
36             Token("rparen", r")"),
37             Token("lbrace", r"{"),
38             Token("rbrace", r"}"),
39             Token("lbracket", r "["),
40             Token("rbracket", r "]" ),
41             Token("comma", r","),
42             Token("semicolon", r";"),
43             Token("colon", r":"),
44             Token("dot", r"\."),
45             Token("plus", r"\+"),
46             Token("minus", r"\-"),
47             Token("multiply", r"\*"),
48             Token("divide", r"/"),
49             Token("modulus", r"%"),
50             Token("assign", r"="),
51             Token("equal", r"=="),
52             Token("not_equal", r"!="),
53             Token("greater", r">"),
54             Token("less", r"<"),
55             Token("greater_equal", r">="),
56             Token("less_equal", r"<="),
57             Token("quote", r"'"),
58             Token("single_quote", r"\'"),
59             Token("and", r"&&"),
60             Token("or", r"\|\|"),
61             Token("not", r"!"),

```

```

62         Token("identifier", r"[a-zA-Z_][a-zA-Z0-9_]*"),
63         Token("int_literal", r"\d+"),
64         Token("float_literal", r"\d+\.\d+"),
65         Token("string_literal", r'".*"'),
66     ]
67
68     def tokenize(self, line):
69         """
70         Tokenizes the given line of code and returns a list
71         of tokens found.
72
73         Args:
74             - line (str): The line of code to be tokenized.
75
76         Returns:
77             - list: A list of Token objects representing the
78             tokens found in the line of code.
79         """
80         tokens_found = []
81         index = 0
82         while index < len(line):
83             match_found = False
84             for token in self.tokens:
85                 pattern = re.compile(token.value)
86                 match = pattern.match(line, index)
87                 if match and match.start() == index:
88                     tokens_found.append((token, match.group(
89 )))
90                     index = match.end()
91                     match_found = True
92                     break
93             if not match_found:
94                 index += 1
95         return tokens_found
96
97     def get_tokens(self):
98         """
99         Returns the list of predefined tokens.
100
101         Returns:
102             - list: A list of Token objects representing the
103             predefined tokens.
104         """
105         return self.tokens
106
107     def print_tokens(self):
108         """
109         Prints the list of predefined tokens.
110         """
111         print(self.tokens)

```

Lexer Class

This class is a wrapper around the Tokenizer class and it encapsulates the Tokenizer class to be easier to use.

```
1 class Lexer:
2     """
3     The Lexer class is responsible for tokenizing input
4     lines from a file.
5
6     Args:
7         - file_path (str): The path to the input file.
8
9     Attributes:
10        - lines (FileReader): An instance of the FileReader
11        class to read lines from the file.
12        - tokens (Tokenizer): An instance of the Tokenizer
13        class to tokenize the lines.
14
15    Methods:
16        - tokenize(): Tokenizes each line from the file and
17        prints the resulting tokens.
18    """
19
20    def __init__(self, file_path: str):
21        self.lines = FileReader(file_path=file_path)
22        self.tokens = Tokenizer()
23
24    def tokenize(self):
25        """
26        Tokenizes each line from the file and prints the
27        resulting tokens.
28        """
29        self.tokens = [self.tokens.tokenize(line) for line
30                        in self.lines.get()]
31
32    def get_tokens(self):
33        """
34        Returns the list of predefined tokens.
35
36        Returns:
37            - list: A list of Token objects representing the
38            predefined tokens.
39        """
40        self.tokenize()
41        return self.tokens
42
43    def print_tokens(self):
44        """
45        Prints the list of predefined tokens.
```

```

39         """
40         self.tokenize()
41         for line in self.tokens:
42             print(line)

```

Main File

```

1 from Lexer import Lexer
2
3 lxr = Lexer(file_path="/Users/cristiancretu/Documents/
    UniCode/LFA/Lexer/example.flexksx")
4 lxr.print_tokens()

```

0.1 Example Usage

I run the main file over this code example:

```

1 print("Hello");
2 print("Hi world");
3 for(i=0;i<5;i+=1) {print("Hello world")}

```

Which gives me the following output:

```

1 [(print, 'print'), (lparen, '('), (quote, '"'), (identifier,
    'Hello'), (quote, '"'), (rparen, ')'), (semicolon, ';')]
2 [(print, 'print'), (lparen, '('), (quote, '"'), (identifier,
    'Hi'), (identifier, 'world'), (quote, '"'), (rparen, ')'),
    (semicolon, ';')]
3 [(for, 'for'), (lparen, '('), (identifier, 'i'), (assign, '='),
    (int_literal, '0'), (semicolon, ';'), (identifier, 'i'),
    (less, '<'), (int_literal, '5'), (semicolon, ';'), (
    identifier, 'i'), (plus, '+'), (assign, '='), (
    int_literal, '1'), (rparen, ')'), (lbrace, '{'), (print,
    'print'), (lparen, '('), (quote, '"'), (identifier, '
    Hello'), (identifier, 'world'), (quote, '"'), (rparen, ')')
    ), (rbrace, '}')]

```

Where each line is a list of token objects, that when printed give a representation of their type and their value.

Conclusions / Screenshots / Results

The program successfully converted my custom written code into tokens that can be further used for future implementations. It used Regular Expressions used from a Python library to implement the identification of certain symbol patterns that represent tokens in order to achieve the desired tokenization.