# Lexer & Scanner

**Course: Formal Languages & Finite Automata**

**Author: Popov Nichita**

---

## Theory

In the realm of computer science, lexical analysis, or tokenization, plays a pivotal role in the interpretation and compilation of programming languages. This process, facilitated by a lexer or scanner, involves dissecting a string of characters into tokens, which are categorized units of meaning that conform to the syntax of the language. Unlike lexemes, which are the raw textual fragments identified by the lexer, tokens are enriched with a layer of semantic meaning, categorizing each fragment under a specific type or class based on predefined language rules. This distinction is crucial for the subsequent phases of parsing and semantic analysis in compilers and interpreters, enabling a structured and efficient processing of source code.

## 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

The implementation focuses on constructing a simple yet functional lexer capable of tokenizing a basic arithmetic expression. The project is structured into several core components:

- `Lexer` Class: Class acts as the orchestrator of the lexical analysis process, meticulously analyzing a string of characters from the input source code to identify meaningful segments known as tokens. This class is meticulously designed to navigate through the input text character by character, applying a set of predefined rules and patterns to differentiate between various token types such as integers, arithmetic operators, and parentheses.

Key functionalities of the `Lexer` class include:

1) Character Navigation: It maintains a pointer to the current character in the input text, allowing for sequential analysis and token identification.
2) Whitespace Handling: Efficiently ignores spaces, tabs, and newlines, which are irrelevant for token formation but essential for human-readable code.
3) Token Recognition: Utilizes pattern matching and state transitions to categorize character sequences into tokens. This involves distinguishing numeric values, operators, and structural characters.

4) Error Detection: Implements a robust error handling mechanism to identify and report unrecognized characters, ensuring the lexer processes only valid syntax elements.

Method `get_next_token` used for tokenization:

This method is crucial for lexical analysis. It sequentially reads the input text and identifies tokens by their patterns:

It skips whitespace, ensuring tokens are separated correctly. It recognizes integers by aggregating consecutive digits. It identifies arithmetic operators (+, -) and parentheses (), providing structure to expressions. It raises an error for unrecognized characters, ensuring only valid tokens are processed. By tokenizing input character by character, it lays the foundation for parsing and further processing.

```python
def get_next_token(self):
    """Tokenize the input text by recognizing tokens one at a time."""
    while self.current_char is not None:
        if self.current_char.isspace():  # Skip whitespace
            self.advance()
            continue

        if self.current_char.isdigit():  # Handle integers
            return Token(TokenType.INTEGER, self.integer())

        if self.current_char == '+':  # Plus token
            self.advance()
            return Token(TokenType.PLUS)

        if self.current_char == '-':  # Minus token
            self.advance()
            return Token(TokenType.MINUS)

        if self.current_char == '(':  # Left parenthesis
            self.advance()
            return Token(TokenType.LPAREN)

        if self.current_char == ')':  # Right parenthesis
            self.advance()
            return Token(TokenType.RPAREN)

        self.error()  # If none match, raise an error

    # Once we're done with all characters, return an EOF token
    return Token(TokenType.EOF)
```

```python
    def error(self):
        """Raise an exception if an invalid character is encountered."""
        raise Exception('Invalid character')
```

- **Token** Class: This class encapsulates the concept of a token, representing the smallest unit of meaningful data in the source code. Each token is characterized by its type (such as INTEGER, PLUS, MINUS) and, when applicable, its value. This class serves as a foundational element for parsing, where the token's type and value become crucial in understanding the syntactic and semantic structure of the code.

Features of the `Token` class include:

1) Type and Value Association: Links a specific type of token with its corresponding value, facilitating a more nuanced interpretation of the source code during subsequent analysis phases.
2) Debugging and Testing Support: Offers a human-readable representation of token instances, simplifying the debugging process and verification of the lexer's output.

`Token` Class:

```python
class Token:
    # A token consists of a type and an optional value.
    def __init__(self, type, value=None):
        self.type = type   # The type of the token (from TokenType)
        self.value = value   # The value of the token (relevant for integers)

    # Representation of the Token instance for debugging and testing.
    def __repr__(self):
        return f"Token({self.type}, {repr(self.value)})"
```

- **TokenType** Class: The `TokenType` class, leveraging Python's Enum, defines an exhaustive enumeration of all possible token types that the lexer can recognize. This approach ensures a well-organized and extensible framework for adding or modifying token types, thereby accommodating different or more complex language features.

Advantages of using the `TokenType` class include: 1) Clear Semantic Distinctions: Each enumerated type distinctly identifies a specific category of token, enhancing code readability and maintainability. 2) Ease of Extension: New token types can be seamlessly integrated into the lexer's logic, supporting the evolution of the language's syntax.

`TokenType` Class:

```python
from enum import Enum, auto

# Define the types of tokens that our lexer can recognize.
class TokenType(Enum):
```

```python
    INTEGER = auto()   # Represents an integer
    PLUS = auto()      # Represents a plus sign '+'
    MINUS = auto()     # Represents a minus sign '-'
    LPAREN = auto()    # Represents a left parenthesis '('
    RPAREN = auto()    # Represents a right parenthesis ')'
    EOF = auto()       # Represents the end of the input
```

- `main.py` Entry Point: Demonstrates the lexer's capabilities by tokenizing a sample arithmetic expression. The program iterates through each token extracted by the lexer, printing it to the console, thereby showcasing the lexer's ability to break down and categorize parts of the expression.

## Results

For example I took a simple expression as a string input:

`3 + 5 - ( 7 - 2 )`

And as the uotput I got:

```
Token(TokenType.INTEGER, 3)
Token(TokenType.PLUS, None)
Token(TokenType.INTEGER, 5)
Token(TokenType.MINUS, None)
Token(TokenType.LPAREN, None)
Token(TokenType.INTEGER, 7)
Token(TokenType.MINUS, None)
Token(TokenType.INTEGER, 2)
Token(TokenType.RPAREN, None)
```

## Conclusion

This project underscores the fundamental role of lexical analysis in the processing of programming languages. Through the practical task of implementing a lexer, we gain a deeper understanding of how compilers and interpreters begin to interpret and compile code. While the scope of this lexer is limited to basic arithmetic expressions, the concepts and techniques explored are applicable to a wide range of programming languages and contexts. This endeavor not only achieves its educational objectives but also paves the way for further exploration into the intricacies of language processing and compiler construction.

## References

- https://llvm.org/docs/tutorial/MyFirstLanguageFrontend/LangImpl01.html
- https://en.wikipedia.org/wiki/Lexical_analysis
- https://www.youtube.com/watch?v=dQw4w9WgXcQ