# Parser

Course: Formal Languages & Finite Automata

Author: Ernu Catalina

## Theory:

### Overview of Lexical Analysis:

Lexical analysis, a foundational phase in the process of compiling or interpreting a programming language, involves extracting lexical tokens from a string of characters. This process, performed by a lexer (or tokenizer/scanner), converts the input source code into meaningful units called tokens. Tokens are categorized based on language-specific rules, and each token type represents a distinct element such as keywords, identifiers, literals, or punctuation symbols.

### Lexemes vs. Tokens:

While lexemes are raw substrings obtained from the input text, tokens are their categorized forms. For instance, a lexeme could be a sequence of characters like "if" or "123", whereas tokens would classify these lexemes into categories such as KEYWORD or NUMBER.

### Historical Context:

The concept of lexical analysis dates back to the early days of programming languages in the 1950s and 1960s. Early compilers, like the one for FORTRAN, included simple lexical analyzers to identify keywords and identifiers. Over time, techniques and tools such as lex (a lexical analyzer generator) and yacc (Yet Another Compiler Compiler) were developed, which significantly advanced the field of lexical analysis and parsing.

### Importance of Lexical Analysis:

Lexical analysis simplifies the parsing process by breaking down the input into manageable tokens. This separation allows for more efficient and organized parsing and analysis of the source code, which is crucial for the development of robust compilers and interpreters.

## Objectives:

Get familiar with parsing, what it is and how it can be programmed [1]. Get familiar with the concept of AST [2]. In addition to what has been done in the 3rd lab work do the following: In case you didn't have a type that denotes the possible types of tokens you need to: Have a type TokenType (like an enum) that can be used in the lexical analysis to categorize the tokens. Please use regular expressions to identify the type of the token. Implement the necessary data structures for an AST that could be used for the text you have processed in the 3rd lab work. Implement a simple parser program that could extract the syntactic information from the input text.

# Lexer

**1. Input Processing:**

The lexer starts by taking in a piece of text as input.

It reads this input character by character, starting from the beginning.

**2. Tokenization Process:**

As it reads each character, the lexer decides what type of token it represents.

It identifies tokens based on specific patterns or rules defined for the programming language.

For example, if the lexer encounters a sequence of letters, it recognizes it as an identifier token. Similarly, if it finds digits or arithmetic operators, it categorizes them as number or operator tokens, respectively.

```
class TokenType(enum.Enum):
    EOF = -1
    DEF = -2
    IDENTIFIER = -4
    NUMBER = -5
    OPERATOR = -6
```

**3. Token Extraction:**

Once the lexer identifies a token, it extracts it from the input text.

This extraction involves collecting all characters that belong to the token until a character is encountered that does not fit the token's pattern.

For instance, if the lexer is extracting an identifier token, it continues collecting letters and digits until it reaches a non-alphanumeric character.

```
def advance(self):
    self.pos += 1
    if self.pos >= len(self.input_text):
        self.current_char = None
    else:
        self.current_char = self.input_text[self.pos]
```

The advance method moves the current position forward by one character. If the end of the input is reached, it sets current_char to None.

```
def get_identifier(self):
    result = ''
    while self.current_char is not None and
(self.current_char.isalnum() or self.current_char == '_'):
        result += self.current_char
        self.advance()
    return Token(TokenType.IDENTIFIER, result)
```

The get_identifier method collects characters to form identifiers (e.g., variable names). It continues until a non-alphanumeric character or underscore is encountered.

**4. Output:**
After processing each token, the lexer returns it along with its type to the caller.
This process continues until the lexer reaches the end of the input text.

# Parser:

**1. Abstract Syntax Trees (ASTs):**
The parser takes the tokens generated by the lexer and constructs an Abstract Syntax Tree (AST).
An AST is a hierarchical representation of the syntactic structure of the input text.
Each node in the tree corresponds to a syntactic construct, such as a function definition or an arithmetic operation.

```
def get_next_token(self):
    while self.current_char is not None:
        if self.current_char.isspace():
            self.skip_whitespace()
            continue
        if self.current_char.isalpha() or self.current_char == '_':
            return self.get_identifier()
        if self.current_char.isdigit() or self.current_char == '.':
            return self.get_number()
        if self.current_char in {'+', '-', '*', '/'}:
            op = self.current_char
            self.advance()
            return Token(TokenType.OPERATOR, op)
        if self.current_char == '=':
            self.advance()
            return Token(TokenType.DEF, '=')
        self.advance()
    return Token(TokenType.EOF, None)
```

This is the main method for retrieving the next token. It handles different character types by delegating to specific methods (e.g., get_identifier, get_number). Special characters are directly converted into tokens, and it skips over unrecognized characters.

**2. Parsing Process:**
The parser reads the tokens provided by the lexer and interprets their sequence to understand the structure of the input text.
It follows the grammar rules of the programming language to organize the tokens into meaningful constructs.
For example, if the parser encounters a series of tokens representing a function definition, it constructs a corresponding node in the AST to capture this definition.

**3. Hierarchical Structure:**
As the parser processes tokens, it organizes them into a hierarchical structure that reflects the nested relationships present in the input text.
For example, an if-else statement within a function body would result in nested nodes within the AST to represent this structure accurately.

**4. Semantic Information:**
In addition to capturing syntactic structure, the parser may also annotate nodes in the AST with semantic information.
This information could include type annotations, variable bindings, or other metadata that aids in later stages of compilation or interpretation.

**5. Output:**
Once the parsing process is complete, the parser returns the fully constructed AST representing the input text's syntactic structure.
This AST can then be used for various purposes, such as code generation, optimization, or further analysis.

**Code Explanation:**
The following Python code defines a lexer that processes input text and breaks it down into tokens. The lexer recognizes different token types such as identifiers, numbers, operators, and special keywords like DEF and EXTERN.

**Explanation:**
- TokenType Enum: Defines various token types like EOF, DEF, IDENTIFIER, NUMBER, OPERATOR, etc.
- Token Class: Represents a token with a type and value.
- Lexer Class: Responsible for tokenizing the input text. It includes methods to advance through the text, skip whitespace, extract identifiers, extract numbers, and get the next token.
- Lexer Methods:
- advance(): Moves to the next character in the input text.
- skip_whitespace(): Skips whitespace characters.
- get_identifier(): Extracts identifiers and keywords.
- get_number(): Extracts numeric values.
- get_next_token(): Main method to iterate through the input text and return tokens based on recognized patterns.

# Conclusions

In conclusion, the lexer implemented in this laboratory effectively breaks down input text into tokens, laying the foundation for further parsing and interpretation. By categorizing characters into tokens such as identifiers, numbers, operators, and keywords, the lexer enables efficient syntactic analysis and subsequent semantic processing. The methods in the lexer systematically iterate through the input text, identify each token type, and extract relevant information, ensuring thorough and accurate tokenization.

This lexer plays a crucial role in language processing pipelines, serving as the initial stage in transforming raw text into structured data. Its ability to recognize and categorize tokens is essential for developing robust compilers, interpreters, and other language processing tools. With a well-designed lexer, language processing tasks can proceed smoothly, facilitating the effective implementation of various programming languages and domain-specific languages.