# A Brief of the code:

This Python script is a simple **lexical analyzer (scanner)**, which takes a source code file as input and breaks its content down into small, meaningful parts called **tokens**. Tokens represent basic elements like numbers, symbols, identifiers (like variable names), and special reserved words.

Here's how it works step-by-step:

- **Token Types Definition**:
  The script first defines constants representing different kinds of tokens, such as:

  - **Identifiers** (variable names, e.g., count, total)

  - **Constants** (integer numbers like 42, floating-point numbers like 3.14)

  - **Operators** (like +, -, *, /)

  - **Symbols** (like (, ), ;, :)

  - **Reserved words** (like if, then, else)

Then it has a **Scanner class** responsible for reading through the input text:

- It goes through the text one character at a time.

- Skips spaces and unnecessary whitespace.

- Recognizes and classifies each sequence as either an identifier, number, operator, or a reserved word.

- Returns tokens one by one, such as:

  - LexicalToken(IDENTIFIER, "x")

  - LexicalToken(SYMBOL_ADD, "+")

  - LexicalToken(INTEGER_CONSTANT, "42")

The script ends by reading a file provided by the user, scanning it, and producing tokens from the source code.

In short, it converts raw source code text into a structured, tokenized format that's easier for further processing by a parser or interpreter.

**Simplified Summary of the Script:**

This Python script is a simple **lexer (or scanner)**, which reads the source code from a file and breaks it down into meaningful components known as **tokens**. Tokens represent the smallest units that have meaning in a programming language, such as keywords (if, else), identifiers (x, counter), numbers (42, 3.14), and symbols (+, ;, =).

Its main purpose is to process raw text code into structured tokens so that later stages (like a parser or interpreter) can more easily analyze the code.

---

**Brief Explanation of How the Code Works:**

**1. Token Definitions**
The script first defines a set of tokens, each represented by a constant (like IDENTIFIER, INTEGER_CONSTANT, SYMBOL_ADD). These represent different elements like numbers, symbols, keywords, and special characters. For convenience, each token has an integer identifier assigned automatically using Python's range.

**Example:**

- INTEGER_CONSTANT token represents numeric values like 123.

- IDENTIFIER token would represent variable names (e.g., x or value).

- Operators and symbols (+, -, ;, (, ), etc.) each have their own tokens.

**Example Token Definitions:**

IDENTIFIER, INTEGER_CONSTANT, SYMBOL_ADD = range(3)

**Reserved Words Mapping**

Reserved words (language keywords such as if, then, else) are mapped explicitly to their own token types in a dictionary:

RESERVED_WORDS = {

  'if': KEYWORD_IF,

  'then': KEYWORD_THEN,

  'else': KEYWORD_ELSE

}

---

**2. Scanner (Lexer)**

The Scanner class is responsible for reading the input source code string character by character, recognizing token patterns, and returning tokens.

**How the Scanner Works:**

- It moves character by character through the input source code.

- Skips over any whitespace.

- Recognizes tokens based on characters (like ;, (, +), numbers, or alphabetic strings (like if or variable names).

**Example:** If the input code is:

x = 5 + 3;

The lexer would produce tokens like:

LexicalToken(IDENTIFIER, 'x'),

LexicalToken(SYMBOL_ASSIGN, '='),

LexicalToken(INTEGER_CONSTANT, '7'),

LexicalToken(SYMBOL_ADD, '+'),

LexicalToken(INTEGER_CONSTANT, '3'),

LexicalToken(SYMBOL_SEMICOLON, ';')

This token list would be passed on to a parser for further processing (though a parser is not explicitly shown in this snippet).

**Token Identification Logic:**

- **Single-character tokens** (;, =, +, -, (, )) are directly recognized and converted into tokens immediately.

- **Two-character tokens** (<=, >=, <>) are handled by checking two characters sequentially.

- **Identifiers and Keywords**: Starts with a letter or underscore and continues with letters, numbers, or underscores. Checks against reserved words to decide if it's a special keyword token or an identifier token.

- **Numbers**: Parses numeric constants, supporting integers or possibly floats (though float handling is partially shown).

### Example: Recognizing Identifiers & Keywords

- if becomes a keyword token (KEYWORD_IF)

- myVar is recognized as an identifier.

### Example: Recognizing Numbers

x = 42;

- x becomes IDENTIFIER

- = becomes SYMBOL_ASSIGN

- 42 recognized as INTEGER_CONSTANT

### Main Execution Flow

When run, the script:

1. Reads source code from a file (specified by command line argument).

2. Creates a Scanner object initialized with this source code.

3. Tokenizes the input continuously until it reaches the end (EOF token).

4. (Implicitly, though not fully shown here) These tokens could then be passed to a parser for syntax analysis or execution.

### How to Use the Script (Example):

Suppose you have a file named example.code with contents:

x = 5 + 3;

Run the lexer like this:

Bash:

python lexer.py example_code.txt

The scanner would output tokens:

LexicalToken(IDENTIFIER, 'x')

LexicalToken(SYMBOL_ASSIGN, '=')

LexicalToken(INTEGER_CONSTANT, '7')

LexicalToken(SYMBOL_ADD, '+')

LexicalToken(INTEGER_CONSTANT, '3')

LexicalToken(SYMBOL_SEMICOLON, ';')

LexicalToken(END_OF_FILE, 'EOF')

**Brief Recap**

In short, this script:

- Reads a programming language text input from a file.

- Converts the text into tokens that represent meaningful elements.

- Provides these tokens for further processing (such as parsing into an executable structure).

This scanner is fundamental in language processing: without tokenizing, the code cannot be parsed, interpreted, or compiled.