FOR SCIENCE & TECHNOLOGY
College of Information
Technology



جامعة مصر للعلوم والتكنولوجيا كلية تكنولوجيا المعلومات



# LEXICAL ANALYZER

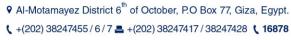


# **Prepared By**

Student Name: Mohamed Hanafy Student ID: 200043742

# **Under Supervision**

Name of Doctor: Nahal AbdelSalam Name of T. A.: Fares Imad Al-din



FOR SCIENCE & TECHNOLOGY
College of Information
Technology





## 1. Introduction

## • 1.1. Phases of Compiler

Explain the different phases of a compiler, including lexical analysis, syntax analysis, semantic analysis, optimization, and code generation. You can highlight the role of the lexical analyzer in breaking down source code into tokens, making the process of parsing and compiling easier.

## 2. Lexical Analyzer

- A lexical analyzer, or lexer, is responsible for reading the input source code and dividing it into a series of tokens. These tokens represent the basic building blocks (keywords, operators, identifiers, etc.) of the source code.
- Discuss the types of tokens that your lexical analyzer identifies, including commands (like int, float, etc.), operators (+, -, etc.), and literals (numeric or decimal literals).
- Mention that the lexical analyzer works by iterating through the characters in the source code, classifying them into different categories based on predefined patterns.

## 3. Software Tools

## • 3.1. Computer Program

Your C++ program implements the lexical analyzer. The program reads the input string, processes each character, and categorizes it into tokens.

## • 3.2. Programming Language

The lexical analyzer is implemented using C++, a general-purpose programming language known for its efficiency and control over system resources. C++ offers powerful features like object-oriented programming, which allows you to use structures like enum class and struct to model tokens and their types efficiently.

FOR SCIENCE & TECHNOLOGY
College of Information
Technology





## 4. Implementation of a Lexical Analyzer

• This section explains the design and implementation of your lexical analyzer.

#### Data Structures

- You use an enum class to define the possible token types (TokenType), such as commands, operators, literals, etc.
- You define a Token structure to store the type and value of each token.

## Functionality

• The tokenize function is the core of your lexical analyzer. It takes an input string, processes each character, and classifies it into tokens based on certain conditions (such as whether a character is a letter, digit, or operator).

## **o** Handling Different Token Types

 Explain how your program distinguishes between keywords (commands), operators, and literals.

#### Tokenization Process

 Describe the process of reading through the input and identifying different token types, starting from recognizing whitespace and moving through various character checks (e.g., isdigit, isalpha, etc.).

## 5. References

• Include all the sources you referenced while working on the project. This might include textbooks, research papers, online tutorials, or documentation that helped in the implementation of the lexical analyzer.



FOR SCIENCE & TECHNOLOGY
College of Information
Technology





## **Additional Notes for the Report:**

- **Figures and Tables:** You might want to include a flowchart or diagram to illustrate the tokenization process, as well as a table showing examples of different token types (commands, operators, literals) and their classifications.
- **Example Code Output:** You can show a sample input and output of your program to illustrate how the lexical analyzer works.
- **Challenges Faced:** Briefly mention any challenges you faced while implementing the lexer (e.g., handling multiple token types, handling invalid tokens) and how you overcame them.

## • Table: Token Types Identified by the Lexical Analyzer

Token Value	<b>Token Type</b>
int,float,string	COMMAND
if,else	COMMAND
return	COMMAND
while	COMMAND
for	COMMAND
cin,cout	COMMAND
A,b,c,	LETTERS
1,2,3,	NUM_LITERAL
1.2,2.5,	DECIMAL_LITERAL
+	PLUS_OP
++	INC_OP
-///	MINUS_OP
	DEC_OP
*	TIMES_OP

FOR SCIENCE & TECHNOLOGY
College of Information
Technology



Token Value	Token Type
/	DIV_OP
٨	POWER_OP
%	MOD_OP
=	ASSIGN_OP
<	LE_COMPARISON
>	RE_COMPARISON
(	L_CIRCLE
)	R_CIRCLE
[	L_SQUARE
]	R_SQUARE
{	L_CURLY
}	R_CURLY
;	SEMICOLON
anything else	INVALID

This table represents all the tokens that the lexical analyzer in your C++ code can identify based on the current implementation. It includes:

- Commands (reserved keywords like int, float, if, etc.),
- Letters (identifiers such as variables or function names),
- Number Literals (integer and decimal values),
- Operators (arithmetic, assignment, comparison, etc.),
- Punctuation (brackets, braces, parentheses, and semicolons),
- Invalid tokens (characters that don't match any expected pattern).

FOR SCIENCE & TECHNOLOGY
College of Information
Technology



جامعة مصر للعلوم والتكنولوجيا كلية تكنولوجيا المعلومات

## • Code

```
#include <iostream>
using namespace std;
enum class TokenType {
  COMMAND, LETTERS, NUM LITERAL, DECIMAL LITERAL,
  PLUS OP, MINUS OP, TIMES OP, DIV OP, POWER OP, MOD OP, ASSIGN OP,
  INC_OP, DEC_OP, LE_COMPARISON, RE_COMPARISON,
  L_SQUARE, R_SQUARE, L_CIRCLE, R_CIRCLE, L_CURLY, R_CURLY, SEMICOLON,
  INVALID
};
struct Token {
  TokenType type;
  string value;
};
bool isCommand(const string& word) {
  return (word == "int" || word == "float" || word == "if" || word == "else" ||
      word == "return" || word == "while" || word == "for" || word == "cin" ||
      word == "cout" || word == "string");
}
const int MAX TOKENS = 100;
int tokenize(const string& input, Token tokens[], int maxTokens) {
  int count = 0;
  size ti = 0;
  while (i < input.length() && count < maxTokens) {</pre>
    char c = input[i];
    if (isspace(c)) {
      i++;
      continue;
    } else if (isalpha(c)) {
      string word;
      while (i < input.length() && (isalnum(input[i]) || input[i] == ' ')) {
        word += input[i];
        i++;
      tokens[count++] = {isCommand(word) ? TokenType::COMMAND :
TokenType::LETTERS, word};
    } else if (isdigit(c)) {
```

FOR SCIENCE & TECHNOLOGY
College of Information
Technology



```
string number;
      bool hasDecimal = false;
      while (i < input.length() && (isdigit(input[i]) | | (input[i] == '.' && !hasDecimal))) {</pre>
         if (input[i] == '.') hasDecimal = true;
         number += input[i];
         i++;
      }
      tokens[count++] = {hasDecimal ? TokenType::DECIMAL LITERAL :
TokenType::NUM LITERAL, number};
    else if (c == '+') {
      if (i + 1 < input.length() && input[i + 1] == '+') {
         tokens[count++] = {TokenType::INC OP, "++"};
         i += 2;
      } else {
         tokens[count++] = {TokenType::PLUS OP, "+"};
      }
    else if (c == '-') {
      if (i + 1 < input.length() && input[i + 1] == '-') {
         tokens[count++] = {TokenType::DEC_OP, "--"};
         i += 2;
      } else {
         tokens[count++] = {TokenType::MINUS OP, "-"};
         i++;
      }
    else if (c == '*') {
      tokens[count++] = {TokenType::TIMES OP, "*"};
      i++:
    } else if (c == '/') {
      tokens[count++] = {TokenType::DIV OP, "/"};
      i++;
    } else if (c == '^') {
      tokens[count++] = {TokenType::POWER OP, "^"};
      i++;
    else if (c == '%') {
      tokens[count++] = {TokenType::MOD OP, "%"};
      j±+:
    } else if (c == '=') {
       tokens[count++] = {TokenType::ASSIGN OP, "="};
      i++;
```

FOR SCIENCE & TECHNOLOGY
College of Information
Technology



```
else if (c == '<') {
      tokens[count++] = {TokenType::LE COMPARISON, "<"};
      i++;
    } else if (c == '>') {
      tokens[count++] = {TokenType::RE COMPARISON, ">"};
      i++;
    } else if (c == '(') {
      tokens[count++] = {TokenType::L CIRCLE, "(");
      i++;
    } else if (c == ')') {
      tokens[count++] = {TokenType::R CIRCLE, ")"};
      i++;
    } else if (c == '[') {
      tokens[count++] = {TokenType::L_SQUARE, "["};
      i++;
    } else if (c == ']') {
      tokens[count++] = {TokenType::R SQUARE, "]"};
      j++;
    else if (c == '{'}) {
      tokens[count++] = {TokenType::L_CURLY, "{"};
      i++;
    } else if (c == '}') {
      tokens[count++] = {TokenType::R_CURLY, "}"};
      i++;
    } else if (c == ';') {
      tokens[count++] = {TokenType::SEMICOLON, ";"};
      i++;
    } else {
      tokens[count++] = {TokenType::INVALID, string(1, c)};
      i++;
    }
  }
  return count;
void printTokens(const Token tokens[], int count) {
  for (int i = 0; i < count; i++) {
    const Token& token = tokens[i];
    cout << "[ " << token.value << " : ";
    switch (token.type) {
      case TokenType::COMMAND: cout << "COMMAND"; break;</pre>
```

FOR SCIENCE & TECHNOLOGY
College of Information
Technology



جامعة مصر للعلوم والتكنولوجيا كلية تكنولوجيا المعلومات

```
case TokenType::LETTERS: cout << "LETTERS"; break;</pre>
      case TokenType::NUM LITERAL: cout << "NUM LITERAL"; break;</pre>
      case TokenType::DECIMAL LITERAL: cout << "DECIMAL LITERAL"; break;
      case TokenType::PLUS OP: cout << "PLUS OP"; break;
      case TokenType::MINUS OP: cout << "MINUS OP"; break;</pre>
      case TokenType::TIMES OP: cout << "TIMES OP"; break;</pre>
      case TokenType::DIV OP: cout << "DIV OP"; break;
      case TokenType::POWER OP: cout << "POWER OP"; break;</pre>
      case TokenType::MOD OP: cout << "MOD OP"; break;</pre>
      case TokenType::ASSIGN OP: cout << "ASSIGN OP"; break;</pre>
      case TokenType::INC OP: cout << "INC OP"; break;</pre>
      case TokenType::DEC_OP: cout << "DEC_OP"; break;</pre>
      case TokenType::LE COMPARISON: cout << "LE COMPARISON"; break;
      case TokenType::RE COMPARISON: cout << "RE COMPARISON"; break;
      case TokenType::L SQUARE: cout << "L SQUARE"; break;</pre>
      case TokenType::R SQUARE: cout << "R SQUARE"; break;</pre>
      case TokenType::L CIRCLE: cout << "L CIRCLE"; break;</pre>
      case TokenType::R CIRCLE: cout << "R CIRCLE"; break;
      case TokenType::L CURLY: cout << "L CURLY"; break;</pre>
      case TokenType::R_CURLY: cout << "R_CURLY"; break;</pre>
      case TokenType::SEMICOLON: cout << "SEMICOLON"; break;</pre>
      case TokenType::INVALID: cout << "INVALID"; break;</pre>
    cout << " ]\n";
  }
}
int main() {
  string input;
  while (true) {
    cout << "Enter source code (type 'exit' to quit): ";
    getline(cin, input);
    if (input == "exit") break;
    Token tokens[MAX_TOKENS];
    int count = tokenize(input, tokens, MAX TOKENS);
    cout << "\nTokens:\n";</pre>
    printTokens(tokens, count);
  return 0;
```

FOR SCIENCE & TECHNOLOGY
College of Information
Technology



جامعة مصر للعلوم والتكنولوجيا كلية تكنولوجيا المعلومات

• Image for output code

```
Enter source code (type 'exit' to quit): +
Tokens:
[ + : PLUS OP ]
Enter source code (type 'exit' to quit): -
Tokens:
[ - : MINUS OP ]
Enter source code (type 'exit' to quit): (
Tokens:
[ ( : L CIRCLE ]
Enter source code (type 'exit' to quit): {
Tokens:
[ { : L CURLY ]
Enter source code (type 'exit' to quit): <
Tokens:
[ < : LE COMPARISON ]
Enter source code (type 'exit' to quit): 1
Tokens:
[ 1 : NUM LITERAL ]
Enter source code (type 'exit' to quit): 3.5
Tokens:
[ 3.5 : DECIMAL LITERAL ]
Enter source code (type 'exit' to quit): exit
```