

## Prepared By

Omar Abdelaziz  
200038794

## Under Supervision

DR :Nehal abselsalam  
ENG :Isalm

```
// مكتبة للإدخال والإخراج (cout, cin) <include <iostream#
// مكتبة فيها دوال للتعامل مع الحروف <include <cctype#
// مكتبة للتعامل مع النصوص <include <string#
// (std::string) (isalpha, isdigit زي)

// عشان نقدر نستخدم cout و string ;using namespace std
// من غير ما نكتب ::std

/* المتغيرات العامة */
// بتحدد نوع الحرف الحالي (حرف، رقم، ;int charClass
// غير معروف)
// الكلمة اللي البرنامج بيجمعها دلوقتي (الـ ;string lexeme
// token)
// الحرف الحالي اللي بنحلله ;char nextChar
// الكود الخاص بالرمز الحالي (مش مستخدم ;int token
// هنا بس ممكن في تطويرات)
// نوع التوكن اللي تم تحليله ;int nextToken
// السطر اللي هندخله ونعالجه ;string input
// مؤشر للمكان الحالي في string input ;int inputIndex = 0

/* تعريف أنواع الحروف */
// تعريف للحرف الأبجدي define LETTER 0#
```

```

// تعريف للرقم define DIGIT 1#
// تعريف للحرف الغير معروف define UNKNOWN 99#

```

/\* تعريف التوكنز \*/

```

// رقم صحيح define INT_LIT 10#
// اسم متغير define IDENT 11#
// علامة يساوي define ASSIGN_OP 20#
// علامة جمع + define ADD_OP 21#
// علامة طرح - define SUB_OP 22#
// علامة ضرب * define MULT_OP 23#
// علامة قسمة / define DIV_OP 24#
// قوس مفتوح ) define LEFT_PAREN 25#
// قوس مغلق ( define RIGHT_PAREN 26#
// نهاية الملف أو نهاية السطر define EOF_TOKEN -1#
// فاصلة منقوطة ; define SEMICOLON 27#

```

/\* تعريف الدوال \*/

```

// دالة بتضيف الحرف الحالي لـ lexeme void addChar();
// دالة بتجيب الحرف الجاي من input void getChar();
// دالة بتجاهل الفراغات void getNonBlank();
// دالة التحليل الرئيسي (lexical analyzer) int lex();
// دالة بتحدد نوع الرموز زي + و * و = int lookup(char ch)

```

```

} void addChar

```

```

// لو lexeme لسه صغير } if (lexeme.length() <= 98)

```

```

// ضيف الحرف الحالي ليه ;lexeme += nextChar

```

```

} else {

```

```

// لو ;cout << "Error - lexeme is too long" << endl

```

الكلمة طويلة جداً

```

{

```

```

    {

        } ()void getChar
        // لو لسه في حروف } if (inputIndex < input.length())
        في النص
        // خذ الحرف التالي ;nextChar = input[inputIndex++]
        وزود المؤشر
        // لو الحرف حرف if (isalpha(nextChar))
        ;charClass = LETTER
        // لو رقم else if (isdigit(nextChar))
        ;charClass = DIGIT
        else
        // أي حاجة تانية ;charClass = UNKNOWN
        } else {
        // لو وصلنا لنهاية ;charClass = EOF_TOKEN
        النص
        // نهاية النص ;nextChar = '\0
        {
        {

        } ()void getNonBlank
        // طالما الحرف فراغ while (isspace(nextChar))
        // اتخطاه وخذ اللي بعده ;()getChar
        {

        } int lookup(char ch)
        } switch (ch)
        case '(': addChar(); nextToken = LEFT_PAREN;
        // قوس مفتوح ;break

```

```

case ')': addChar(); nextToken = RIGHT_PAREN;
           // قوس مغلق ;break
;case '+': addChar(); nextToken = ADD_OP; break
           + //
;case '-': addChar(); nextToken = SUB_OP; break
           - //
    case '*': addChar(); nextToken = MULT_OP;
           * // ;break
;case '/': addChar(); nextToken = DIV_OP; break
           / //
case ';': addChar(); nextToken = SEMICOLON;
           ; // ;break
case '=': addChar(); nextToken = ASSIGN_OP;
           = // ;break
default: addChar(); nextToken = EOF_TOKEN;
           // أي رمز ثاني غير معروف ;break
           {
               // رجع نوع التوكن ;return nextToken
           }

           } ()int lex
           ;"" = lexeme
           // نبدأ بكلمة فاصية
           // نتخطى أي فراغات ;()getNonBlank

           } switch (charClass)
           :case LETTER
           حرف
           lexeme لـ ضيفه // ;()addChar
           // خذ اللي بعده ;()getChar

```

```

while (charClass == LETTER || charClass ==
                                } DIGIT)
    طالما الحروف أو أرقام،      ;()addChar
                                ضيفهم
                                // وكمل      ;()getChar
                                {
                                // عرفنا إن دي متغير      ;nextToken = IDENT
                                ;break
                                :case DIGIT
                                ;()addChar
                                ;()getChar
                                } while (charClass == DIGIT)
                                أرقام
                                ;()addChar
                                ;()getChar
                                {
                                // رقم صحيح      ;nextToken = INT_LIT
                                ;break
                                :case UNKNOWN
                                // رموز زي +، =
                                ... ، *
                                // نحدد نوعها      ;lookup(nextChar)
                                // نجيب الحرف اللي بعده      ;()getChar
                                ;break
                                :case EOF_TOKEN
                                // لو خالصنا كل حاجة      ;nextToken = EOF_TOKEN
                                // نكتب نهاية الملف      ;"lexeme = "EOF
                                ;break

```

```

    {
        // نطبع نوع التوكن والكلمة اللي لقيناها
        cout << "Next token is: " << nextToken << ", Next
            ;lexeme is " << lexeme << endl
        // نرجع نوع التوكن
    }

} ()int main
    // السطر اللي عايزين نحلله
    ;input = "sum = 9 + 2 * 5
    // نبدأ من أول السطر
    ;inputIndex = 0

    // نقرأ أول حرف
    ;()getChar
    } do
    // نحلل التوكن اللي بعده
    ;()lex
    ;(while (nextToken != EOF_TOKEN {
        // لحد ما نوصل
        // للنهاية

    // نهاية البرنامج
    ;return 0
    }

```

## 1. Introduction

This program is a basic implementation of a Lexical Analyzer (or scanner), which is the first phase of a compiler. It reads the input string and breaks it into tokens such as identifiers, integers, and operators. This allows the next phase of the compiler to understand the structure of the program.

### 1.1. Phases of Compiler

The phases of a compiler include:

- Lexical Analysis: Breaks input into tokens.

- Syntax Analysis: Checks the structure using grammar rules.
- Semantic Analysis: Validates meaning (e.g., type checking).
- Intermediate Code Generation: Converts into intermediate representation.
- Code Optimization: Improves performance.
- Code Generation: Produces final machine code.
- Symbol Table Management: Tracks variable and function info.
- Error Handling: Reports and manages errors throughout compilation.

## **2. Lexical Analyzer**

The lexical analyzer reads the source input (in this case, a string like 'sum = 9 + 2 \* 5;') and splits it into tokens.

Each token is categorized, such as:

- IDENT: Identifier like 'sum'
- INT\_LIT: Integer literal like 9 or 5
- ADD\_OP, MULT\_OP: Operators like '+' or '\*'

Whitespace is ignored, and characters are processed one at a time using functions like `getChar()` and `lookup()`.

## **3. Software Tools**

For lexical analysis, developers can use tools like:

- Lex / Flex: For token generation.
- Custom C++/Java code: As shown in this example.
- IDEs and debuggers for testing code logic.

### **3.1. Computer Program**

A computer program is a sequence of instructions written to perform a specific task. In this case, the

program reads an expression and categorizes parts of it into tokens. It's written in C++ and simulates lexical analysis behavior.

### **3.2. Programming Language**

C++ is used to write this lexical analyzer. It's a high-level programming language that supports procedural and object-oriented paradigms. C++ allows for control over memory and efficient character processing, which is useful in compiler development.

#### **4. Implementation of a Lexical Analyzer**

```
<include <iostream#  
<include <fstream#  
<include <cctype#  
<include <cstdlib#  
<include <string#
```

```
;using namespace std
```

```
/* Global declarations */
```

```
/* Variables */
```

```
;int charClass
```

```
;string lexeme
```

```
;char nextChar
```

```
;int token
```

```
;int nextToken
```

```
;ifstream in_fp
```

```
/* Function declarations */
```

```
;()void addChar
```

```
;()void getChar
```



```

;()void getNonBlank
;()int lex

/* Character classes */
define LETTER 0#
define DIGIT 1#
define UNKNOWN 99#

/* Token codes */
define INT_LIT 10#
define IDENT 11#
define ASSIGN_OP 20#
define ADD_OP 21#
define SUB_OP 22#
define MULT_OP 23#
define DIV_OP 24#
define LEFT_PAREN 25#
define RIGHT_PAREN 26#
define EOF_TOKEN -1#
define SEMICOLON 27#

/* main driver */
} ()int main
/* Open the input data file and process its contents */
    Since we're in an online compiler, let's create a //
        temporary input string
        instead of reading from a file //
;";string input = "sum = 9 + 2 * 5

Create a temporary file for input simulation //

```

```

        ;ofstream temp_file("front.in")
            ;temp_file << input
            ;()temp_file.close

        ;in_fp.open("front.in")

        } if (!in_fp.is_open())
;cout << "ERROR - cannot open front.in" << endl
            ;return 1
        {

            ;()getChar
            } do
            ;()lex
;(while (nextToken != EOF_TOKEN {

            ;()in_fp.close
            ;return 0
        {

lookup - a function to lookup operators and */
/* parentheses and return the token
    } int lookup(char ch)
        } switch (ch)
            :')' case
                ;()addChar
;nextToken = LEFT_PAREN
                ;break
            : '(' case
                ;()addChar

```

```
;nextToken = RIGHT_PAREN
                                ;break
                                : '+' case
                                ;()addChar
;nextToken = ADD_OP
                                ;break
                                : '-' case
                                ;()addChar
;nextToken = SUB_OP
                                ;break
                                : '*' case
                                ;()addChar
;nextToken = MULT_OP
                                ;break
                                : '/' case
                                ;()addChar
;nextToken = DIV_OP
                                ;break
                                : ';' case
                                ;()addChar
;nextToken = SEMICOLON
                                ;break
                                : '=' case
                                ;()addChar
;nextToken = ASSIGN_OP
                                ;break
                                : default
                                ;()addChar
;nextToken = EOF_TOKEN
                                ;break
```

```

        {
            ;return nextToken
        }

/* addChar - a function to add nextChar to lexeme */
        } ()void addChar
    } if (lexeme.length() <= 98)
        ;lexeme += nextChar
        } else {
;cout << "Error - lexeme is too long" << endl
        {
        {

getChar - a function to get the next character of input */
/* and determine its character class
        } ()void getChar
if (in_fp.get(nextChar)) { // Using ifstream's get()
                                method
        if (isalpha(nextChar))
            ;charClass = LETTER
        else if (isdigit(nextChar))
            ;charClass = DIGIT
        else
            ;charClass = UNKNOWN
        } else {
            ;charClass = EOF_TOKEN
nextChar = EOF; // Set nextChar to EOF when
                    end of file is reached
        {
        {

```

```

getNonBlank - a function to call getChar until it */
/* returns a non-whitespace character
    } ()void getNonBlank
while (isspace(nextChar))
    ;()getChar
    {

lex - a simple lexical analyzer for arithmetic */
/* expressions
    } ()int lex
lexeme = ""; // Clear previous lexeme
    ;()getNonBlank

    } switch (charClass)
/* Parse identifiers */
    :case LETTER
        ;()addChar
        ;()getChar
while (charClass == LETTER || charClass ==
        } DIGIT)
        ;()addChar
        ;()getChar
        {
;nextToken = IDENT
        ;break

/* Parse integer literals */
    :case DIGIT
        ;()addChar

```

```

                                ;()getChar
} while (charClass == DIGIT)
                                ;()addChar
                                ;()getChar
                                {
;nextToken = INT_LIT
                                ;break

/* Parentheses and operators */
:case UNKNOWN
;lookup(nextChar)
                                ;()getChar
                                ;break

                                /* EOF */
:case EOF_TOKEN
;nextToken = EOF_TOKEN
;lexeme = "EOF
                                ;break
                                {

cout << "Next token is: " << nextToken << ", Next
;lexeme is " << lexeme << endl
;return nextToken
                                {

```