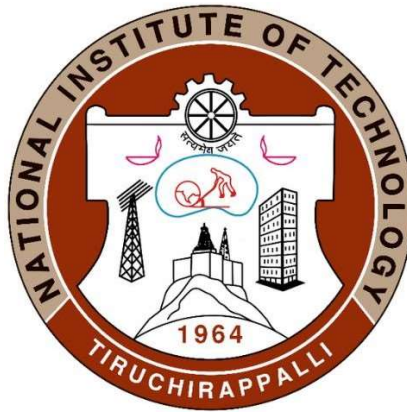


NATIONAL INSTITUTE OF TECHNOLOGY, TIRUCHIRAPPALLI



CSPC62 COMPILER DESIGN TOPIC: C++ Compiler LAB REPORT-1 Sub Topic : Lexical Analyzer

DONE BY:

S.No	Name	RollNo
1.	Hema Sai Dorababu K	106121065
2.	Chetan Reddy T	106121139
3.	Harshith Babu B	106121029

Introduction:

The development of programming languages has been pivotal in shaping the landscape of modern computing. Among these languages, C++ stands out as a powerful and versatile tool for software development, renowned for its efficiency, performance, and object-oriented features. However, the process of transforming human-readable C++ code into machine-executable instructions requires a critical intermediary: the compiler.

This project report delves into the intricate world of compiler design, focusing specifically on the creation of a C++ compiler. Our aim is to provide a comprehensive understanding of the underlying principles, methodologies, and challenges involved in building such a fundamental tool.

Throughout this report, we will explore the various phases of compiler construction, starting from lexical analysis and syntax parsing, to semantic analysis, intermediate code generation, optimization, and ultimately, code generation. Each phase plays a crucial role in translating high-level C++ code into efficient machine code that can be executed by a computer.

QUESTION-1:

Components of the Development Language:

1.Keywords:

Keywords, also known as reserved words, are predefined tokens with special meanings in the programming language. These words are reserved by the language and cannot be used for any other purpose such as naming variables or functions. Examples of keywords in C++ include `int`, `double`, `if`, `else`, `for`, `while`, `class`, `struct`, `return`, `namespace`, `public`, `private`, and `virtual`, among others. Keywords typically represent language constructs, control flow statements, data types, access specifiers, and other fundamental elements of the language.

2.Identifiers:

Identifiers are user-defined names used to represent various program elements such as variables, functions, classes, and namespaces. An identifier can consist of letters, digits, and underscores, with the first character being a letter or an underscore.

Identifiers must follow certain rules:

- They cannot be a keyword or contain a keyword.
- They cannot contain spaces or special characters (except underscores).
- They are case-sensitive (e.g., `myVariable`, `MyVariable`, and `myvariable` are considered different identifiers).
- Examples of identifiers include variable names (`int age`, `double balance`), function names (`void calculateInterest()`), class names (`class Car`, `struct Point`), and namespace names (`namespace MathUtils`), among others.

C++ supports various data types such as integers (`int`, `short`, `long`), floating-point numbers (`float`, `double`), characters (`char`), Boolean (`bool`), and user-defined types (structures, classes, enumerations, unions).

3. Control Structures:

Control structures in C++ facilitate decision-making and looping in programs. These include if-else statements, switch-case statements, while loops, do-while loops, and for loops.

4. Functions:

Functions are blocks of code that perform specific tasks. In C++, functions can be defined to accept parameters and return values. C++ also supports function overloading, allowing multiple functions with the same name but different parameter lists.

5. Classes and Objects:

C++ is an object-oriented programming (OOP) language, and classes are the building blocks of object-oriented design. A class defines a blueprint for creating objects that encapsulate data and behavior. Objects are instances of classes and can interact with each other through member functions and variables.

6. Headers:

In C++, headers are files that contain declarations and definitions used by other files in a program. They typically have a .h extension, although it's not a strict requirement. Here's some key information about headers in C++:

Include Directives: Headers are included in C++ source files using the `#include` directive. For example:

```
#include <iostream> // Includes the iostream header
```

Contents of Headers: Headers may contain function prototypes, class declarations, constants, macros, template definitions, and inline function implementations.

Standard Library Headers: C++ provides a standard library, often referred to as the Standard Template Library (STL), which contains many pre-defined headers for common functionalities. Examples include:

- `iostream`: Input/output operations like `cout` and `cin`.
- `vector`: Dynamic array implementation.
- `string`: String manipulation functions.
- `algorithm`: Algorithms like sorting and searching.
- `cmath`: Mathematical functions like `sqrt` and `sin`.
- `bits/stdc++.h`: to include all the libraries

User-Defined Headers: Programmers can create their own headers to organize and modularize their code. These headers typically contain declarations for classes, functions, and other entities defined in corresponding source files.

7. Relational Operators:

Relational operators in C++ are symbols used to compare the relationship between two operands. They evaluate to a boolean value (true or false) based on whether the specified relationship holds true. Here are the relational operators in C++:

- Equality (==): Compares if two operands are equal.
- Inequality (!=): Compares if two operands are not equal.
- Greater Than (>): Checks if the left operand is greater than the right operand.
- Less Than (<): Checks if the left operand is less than the right operand.
- Greater Than or Equal To (>=): Checks if the left operand is greater than or equal to the right operand.
- Less Than or Equal To (<=): Checks if the left operand is less than or equal to the right operand.

8.Punctuation:

Braces, brackets, and semicolons are punctuation symbols used in C++ to define structure, delineate blocks of code, and terminate statements. Here's a brief explanation of each:

❖ Braces {}:

Braces are used to define the beginning and end of code blocks in C++. They are commonly used in functions, control flow statements (such as if-else statements, loops), and to define the body of classes, structs, and namespaces.

Example:

```
if (condition) {
    // Code block
}
```

❖ Brackets []:

Brackets have multiple uses like Array indexing:- Used to access elements of an array, in Array declaration:- Used to declare arrays, Used in lambda expressions for capture lists.

Example:

```
int arr[5]; // Declaration of an integer array with 5 elements
int x = arr[2]; // Accessing the third element of the array
```

❖ Semicolons ;:

Semicolons are used to terminate statements in C++. They are placed at the end of each statement to indicate the end of the statement.

Example:

```
int x = 5; // Statement assigning the value 5 to variable x
```

❖ Parentheses ():

Parentheses have various uses in C++:

Function calls:- Used to enclose the arguments passed to a function.

Expression grouping:- Used to specify the order of evaluation in expressions.

Control flow statements:- Used to enclose conditions in if, while, for statements.

Example:

```
int result = add(3, 5); // Function call with arguments 3 and 5
```

These punctuation symbols are fundamental in C++ syntax and play a crucial role in defining the structure and behavior of C++ programs.

❖ Commas:

Commas (,) are punctuation symbols used in C++ for various purposes. Here's a brief explanation of their uses:

Separating Items in Lists:- Commas are commonly used to separate items in lists, such as function arguments, variable declarations, and elements in initializer lists.

Example:

```
int a, b, c; // Declaring multiple variables
```

```
void foo(int x, int y) { // Function body }
```

```
int arr[] = {1, 2, 3, 4, 5}; // Initializing an array
```

In Function Calls:-

Commas are used to separate arguments in function calls.

Example:

```
int sum = add(3, 5); // Function call with two arguments
```

In Initialization:-

Commas are used to separate elements in initialization lists for arrays, structs, and classes.

Example:

```
int arr[] = {1, 2, 3, 4, 5}; // Initializing an array
```

```
struct Point {
```

```
    int x; int y;
```

```
};
```

```
Point p = {10, 20}; // Initializing a struct
```

In For Loops:-

Commas are used in the initialization, condition, and iteration expressions of for loops.

Example: `for (int i = 0, j = 10; i < 5; ++i, --j) { // Loop body }`

In Enumerations:-

Commas are used to separate enumerators in enumerations.

Example: `enum Color { RED, GREEN, BLUE};`

Question-2:

Regular Expressions and Definitions:

```
delim [ \t]
ws {delim}+
letter [a-zA-Z]
digit [0-9]
id {letter}({letter}|{digit})*
relop (<|=|>|=|!=|<|>)
leftshift (<<)
rightshift (>>)
plus [+]
minus [-]
mult [*]
div [/]
num ({digit}+)
float ({num}\.{num})
arithmeticop ({plus}|{minus}|{mult}|{div})
increment {plus}{plus}
decrement {minus}{minus}
assignop =
string (\\"(\\.|[^\\""])*\\"|\'(\\.|[^\\"'])*\')
keyword (if|else|const|while|for|int|float|return|void|main|char|"long
long"|double|short|long|unsigned|signed|define|struct|enum|typedef|sizeof|static|regi
ster|auto|break|case|continue|default|do|goto|switch|cout|cin|endl|bool|using|namespa
ce|std|include|iostream|vector|map|set|queue|stack|push_back|pop_back|pop|push|top|fr
ont|priority_queue)
inval (({digit}+{id}))
```

LEXICAL ANALYZER

Question-3:

Code Snippet:

```
%{
    #include <ctype.h>
    #include <stdio.h>
    #include <string.h>
    #include <stdlib.h>
    #include <stdbool.h>
    #include <math.h>
    extern int yylex();
    extern char *yytext;
    char *prevtoken;
    char *tokentobeadded;
    int idCount = 0;
    int lineno = 0;
}%
/* regular definitions */
delim [ \t]
ws {delim}+
letter [a-zA-Z]
digit [0-9]
id {letter}({letter}|{digit})*
relop (<|=|>|=|!=|<|>)
leftshift (<<)
rightshift (>>)
plus [+]
minus [-]
mult [*]
div [/]
num ({digit}+)
float ({num}\.{num})
arithmeticop ({plus}|{minus}|{mult}|{div})
increment {plus}{plus}
decrement {minus}{minus}
assignop =
string (\\"(\\.|[^\\""])*\\"|\'(\\.|[^\\"'])*\')
keyword (if|else|const|while|for|int|float|return|void|main|char|"long
long"|double|short|long|unsigned|signed|define|struct|enum|typedef|sizeof|static|regi
ster|auto|break|case|continue|default|do|goto|switch|cout|cin|endl|bool|using|namespa
ce|std|include|iostream|vector|map|set|queue|stack|push_back|pop_back|pop|push|top|fr
ont|priority_queue)
inval (({digit}+{id}))
```

```

%%
"bits/stdc++.h" {printf("HEADER: %s, line: %d\n", yytext, lineno+1);}
"\n" {lineno++;printf("\nNEWLINE: %s, line: %d", yytext, lineno+1);}
"/*.* { /* ignore comments */ }
"/*.*"*/" { /* ignore comments */ }
{keyword} {
    printf("KEYWORD: %s    line: %d\n", yytext, lineno+1);
}
{id} {
    prevtoken = "id";
    printf("ID: %s, line: %d\n", yytext, lineno+1);
}
{ws} { /* ignore whitespace */ }
{leftshift} {
    printf("LeftShift: %s, line: %d\n", yytext, lineno+1);
}
{rightshift} {
    printf("RightShift: %s, line: %d\n", yytext, lineno+1);
}
{relop} {
    printf("RELOP: %s, line: %d\n", yytext, lineno+1);
}
{inval} {printf("\nInvalid token: %s, line: %d\n", yytext, lineno+1);}
{increment} {printf("INCREMENT: %s, line: %d\n", yytext, lineno+1);}
{decrement} {printf("DECREMENT: %s, line: %d\n", yytext, lineno+1);}
{float} {printf("FLOAT: %s, line: %d\n", yytext, lineno+1);}
{num} {
    if(tokentobeadded && strcmp(tokentobeadded, "-") == 0) {
        printf("Number:-%s, line: %d\n", yytext, lineno+1);
        prevtoken = NULL;
        tokentobeadded = NULL;
    }else{
        printf("Number: %s, line: %d\n", yytext, lineno+1);
    }
}
{arithmeticop} {
    if(strcmp(prevtoken, "assign") == 0 && strcmp(yytext, "-") == 0) {
        tokentobeadded = "-";
    }else{
        printf("ARITHMETICOP: %s, line: %d\n", yytext, lineno+1);
    }
}
"." {printf("DOT: %s, line: %d\n", yytext, lineno+1);}
"#" {printf("Hash: %s, line: %d\n", yytext, lineno+1);}
"(" {printf("LPAREN: %s, line: %d\n", yytext, lineno+1);}
")" {printf("RPAREN: %s, line: %d\n", yytext, lineno+1);}
"[" {printf("LBRACKET: %s, line: %d\n", yytext, lineno+1);}
"]" {printf("RBRACKET: %s, line: %d\n", yytext, lineno+1);}
{" {printf("LBRACE: %s, line: %d\n", yytext, lineno+1);}
"}" {printf("RBRACE: %s, line: %d\n", yytext, lineno+1);}
"=" {prevtoken = "assign";printf("ASSIGNOP: %s, line: %d\n", yytext, lineno+1);}

```



```

";" {printf("Delimiter: %s, line: %d\n", yytext, lineno+1);}
", " {printf("COMMA: %s, line: %d\n", yytext, lineno+1);}
{string} {printf("string: %s, line: %d\n", yytext, lineno+1);}
. {printf("Invalid token: %s, line: %d\n", yytext, lineno+1);}
%%
int main() {
    FILE *fp = fopen("test.cpp", "r");
    yyin = fp;
    yylex();
    return 0;
}
int yywrap(){ return 1; }

```

Question-4:

Error/Ambiguity Handling:

- When “-” is generated as a token then there are two possibilities, they are

- “=” comes in front of the “-”
- A variable comes in front of the “-”.

Regarding the above ambiguity we use below if else conditions to check the possibility of existence of one of the situations and make the use of “-” correctly.

```

{num} {
    if(tokenbeadded && strcmp(tokenbeadded, "-") == 0) {
        printf("Number: %s, line: %d\n", yytext, lineno+1);
        prevtoken = NULL;
        tokenbeadded = NULL;
    }else{
        printf("Number: %s, line: %d\n", yytext, lineno+1);
    }
}
{arithmeticop} {
    if(strcmp(prevtoken, "assign") == 0 && strcmp(yytext, "-") == 0) {
        tokenbeadded = "-";
    }else{
        printf("ARITHMETICOP: %s, line: %d\n", yytext, lineno+1);
    }
}

```

- If there is any invalid identifier which starts with number will be identified as invalid tokens by using below code snippet.

```

inval (({digit}+{id})) // regular expression for invalid token
{inval} {printf("\nInvalid token: %s, line: %d\n", yytext, lineno+1);} // prints
invalid token.

```

Question-5:

Sample Program:

```
#include <bits/stdc++.h>
int main()
{
    int a123 = -4;
    a123 = a - 4;
    a123 += -4;
    b = 4.56;
    long long z = 10;
    // cout<<"Hello World";
    /* cout<<"Jai Baalayya"*/
    cout << "Hai ,how"
         " are you?";
    for (int i = 0; i < 10; i++)
    {
        cout << i << endl;
    }
    while (1)
    {
        cout << "Hello World" << 'c';
    }
    return 0;
}
```

Generated Tokens Output:

Hash: #, line: 1

KEYWORD: include line: 1

RELOP: <, line: 1

HEADER: bits/stdc++.h, line: 1

RELOP: >, line: 1

NEWLINE:

, line: 2KEYWORD: int line: 2

KEYWORD: main line: 2

LPAREN: (, line: 2

RPAREN:), line: 2

NEWLINE:

, line: 3LBRACE: {, line: 3

NEWLINE:

, line: 4KEYWORD: int line: 4

ID: a123, line: 4

ASSIGNOP: =, line: 4

Number:-4, line: 4

Delimiter: ,, line: 4

NEWLINE:

, line: 5ID: a123, line: 5

ASSIGNOP: =, line: 5

ID: a, line: 5

ARITHMETICOP: -, line: 5

Number: 4, line: 5

Delimiter: ,, line: 5

NEWLINE:

, line: 6ID: a123, line: 6

ARITHMETICOP: +, line: 6

ASSIGNOP: =, line: 6

Number:-4, line: 6

Delimiter: ,, line: 6

NEWLINE:

, line: 7ID: b, line: 7

ASSIGNOP: =, line: 7

FLOAT: 4.56, line: 7

Delimiter: ,, line: 7

NEWLINE:

, line: 8KEYWORD: long long line: 8

ID: z, line: 8

ASSIGNOP: =, line: 8

Number: 10, line: 8

Delimiter: ,, line: 8

NEWLINE:

, line: 9

NEWLINE:

, line: 10

NEWLINE:

, line: 11KEYWORD: cout line: 11

LeftShift: <<, line: 11

string: "Hai ,how", line: 11

NEWLINE:

, line: 12string: " are you?", line: 12

Delimiter: ,, line: 12

NEWLINE:

, line: 13KEYWORD: for line: 13

LPAREN: (, line: 13

KEYWORD: int line: 13

ID: i, line: 13

ASSIGNOP: =, line: 13

Number: 0, line: 13

Delimiter: ,, line: 13

ID: i, line: 13

RELOP: <, line: 13

Number: 10, line: 13

Delimiter: ,, line: 13

ID: i, line: 13

INCREMENT: ++, line: 13

RPAREN:), line: 13

NEWLINE:

, line: 14LBACE: {, line: 14

NEWLINE:

, line: 15KEYWORD: cout line: 15

LeftShift: <<, line: 15

ID: i, line: 15

LeftShift: <<, line: 15

KEYWORD: endl line: 15

Delimiter: ;, line: 15

NEWLINE:

, line: 16RBACE: }, line: 16

NEWLINE:

, line: 17KEYWORD: while line: 17

LPAREN: (, line: 17

Number: 1, line: 17

RPAREN:), line: 17

NEWLINE:

, line: 18LBACE: {, line: 18

NEWLINE:

, line: 19KEYWORD: cout line: 19

LeftShift: <<, line: 19

string: "Hello World", line: 19

LeftShift: <<, line: 19

string: 'c', line: 19

Delimiter: ,, line: 19

NEWLINE:

, line: 20RBACE: }, line: 20

NEWLINE:

, line: 21KEYWORD: return line: 21

Number: 0, line: 21

Delimiter: ,, line: 21

NEWLINE:

, line: 22RBACE: }, line: 22

Question-6:

C++ programs:

1. Without errors (already executed in question 4)

```
#include <bits/stdc++.h>
int main()
{
    int a123 = -4;
    a123 = a - 4;
    a123 += -4;
    b = 4.56;
    long long z = 10;
    // cout<<"Hello World";
    /* cout<<"Jai Baalayya"*/
    cout << "Hai ,how" " are you?";
    for (int i = 0; i < 10; i++)
    {
        cout << i << endl;
    }
    while (1)
    {
        cout << "Hello World" << 'c';
    }
    return 0;
```

```
}
```

2. Code with errors

```
#include <bits/stdc++.h>
int main()
{
    int a123 = -4;
    a123 = a - 4;
    a123 += -4;
    b = 4.56;
    long long z = 10;
    // cout<<"Hello World";
    /* cout<<"Jai Baalayya"*/
    cout << "Hai ,how" " are you?";
    for (int i = 0; i < 10; i++)
    {
        cout << i << endl;
    }

    while (1)
    {
        cout << "Hello World" << 'c';
    }
    Int 123ab = 5;
    return 0;
}
```

Question-7:

Output of the above erroneous code:

Hash: #, line: 1

KEYWORD: include line: 1

RELOP: <, line: 1

HEADER: bits/stdc++.h, line: 1

RELOP: >, line: 1

NEWLINE:

, line: 2KEYWORD: int line: 2

KEYWORD: main line: 2

LPAREN: (, line: 2

RPAREN:), line: 2

NEWLINE:

, line: 3 LBRACE: {, line: 3

NEWLINE:

, line: 4 KEYWORD: int line: 4

ID: a123, line: 4

ASSIGNOP: =, line: 4

Number: -4, line: 4

Delimiter: ;, line: 4

NEWLINE:

, line: 5 ID: a123, line: 5

ASSIGNOP: =, line: 5

ID: a, line: 5

ARITHMETICOP: -, line: 5

Number: 4, line: 5

Delimiter: ;, line: 5

NEWLINE:

, line: 6 ID: a123, line: 6

ARITHMETICOP: +, line: 6

ASSIGNOP: =, line: 6

Number: -4, line: 6

Delimiter: ;, line: 6

NEWLINE:

, line: 7 ID: b, line: 7

ASSIGNOP: =, line: 7

FLOAT: 4.56, line: 7

Delimiter: ;, line: 7

NEWLINE:

, line: 8KEYWORD: long long line: 8

ID: z, line: 8

ASSIGNOP: =, line: 8

Number: 10, line: 8

Delimiter: ,, line: 8

NEWLINE:

, line: 9

NEWLINE:

, line: 10

NEWLINE:

, line: 11KEYWORD: cout line: 11

LeftShift: <<, line: 11

string: "Hai ,how", line: 11

string: " are you?", line: 11

Delimiter: ,, line: 11

NEWLINE:

, line: 12KEYWORD: for line: 12

LPAREN: (, line: 12

KEYWORD: int line: 12

ID: i, line: 12

ASSIGNOP: =, line: 12

Number: 0, line: 12

Delimiter: ,, line: 12

ID: i, line: 12

RELOP: <, line: 12

Number: 10, line: 12

Delimiter: ,, line: 12

ID: i, line: 12

INCREMENT: ++, line: 12

RPAREN:), line: 12

NEWLINE:

, line: 13LBACE: {, line: 13

NEWLINE:

, line: 14KEYWORD: cout line: 14

LeftShift: <<, line: 14

ID: i, line: 14

LeftShift: <<, line: 14

KEYWORD: endl line: 14

Delimiter: ;, line: 14

NEWLINE:

, line: 15RBACE: }, line: 15

NEWLINE:

, line: 16

NEWLINE:

, line: 17KEYWORD: while line: 17

LPAREN: (, line: 17

Number: 1, line: 17

RPAREN:), line: 17

NEWLINE:

, line: 18LBACE: {, line: 18

NEWLINE:

, line: 19KEYWORD: cout line: 19

LeftShift: <<, line: 19

string: "Hello World", line: 19

LeftShift: <<, line: 19

string: 'c', line: 19

Delimiter: ;, line: 19

NEWLINE:

, line: 20
RBRACE: }, line: 20

NEWLINE:

, line: 21
ID: Int, line: 21

Invalid token: 123ab, line: 21

ASSIGNOP: =, line: 21

Number: 5, line: 21

Delimiter: ;, line: 21

NEWLINE:

, line: 22
KEYWORD: return line: 22

Number: 0, line: 22

Delimiter: ;, line: 22

NEWLINE:

, line: 23
RBRACE: }, line: 23

NEWLINE:

, line: 24