

Lexical Analyzer

Group Members

- Abdullah Ejaz Janjua 2023038
- Muhammad Ahmad Amjad 2023361

Lexical Analyzer Code

The lexical analyzer was implemented using Flex to recognize and categorize various tokens in C++ source code. The analyzer identifies tokens such as keywords, identifiers, operators, separators, literals, and comments while tracking line numbers.

```
%{
#include <stdio.h>
#include <string.h>
int line = 1;
char filename[100];
char ofilename[100];
FILE *inputfile, *outputfile;

%}
%%

"//"*. *  {}
"/*"
{
    int c;
    while((c = input()) != 0)
    {
        if(c == '\n') line++;
        if(c == '*' && (c = input()) == '/') break;
        if(c == 0) break;
    }
}

^#include[ \t]*<[^>]+> { fprintf(outputfile, "Line %d: %s \t ———> Header\n",
```

```

^#define[ ]* { fprintf(outputfile,"Line %d: %s \t ———> Pre-process Directives\n", line, yytext); }
int|float|char|double|void|bool|short|long|signed|unsigned|string { fprintf(outputfile,"Line %d: %s \t ———> Keyword\n", line, yytext); }
if|else|while|for|do|switch|case|break|continue|return|goto|cout|cin|using|namespace { fprintf(outputfile,"Line %d: %s \t ———> Operator\n", line, yytext); }
"++"|"--"|"+="|"-=|"*="|"/="|"%=|"&="|"|"|"^="|"<="|">="|"&&"|"||" { fprintf(outputfile,"Line %d: %s \t ———> Operator\n", line, yytext); }
"<"|">"|"<="|">="|"+"| "-"| "*"| "/"| "="| "<"| ">"| "=="| "!="| "%"| "!"| "&"| "|"| "^"|"~" { fprintf(outputfile,"Line %d: %s \t ———> Operator\n", line, yytext); }

[{}();,\[\]] { fprintf(outputfile,"Line %d: %s \t ———> Separator\n", line, yytext); }

[0-9]+\.[0-9]+ { fprintf(outputfile,"Line %d: %s \t ———> Float\n", line, yytext); }

[0-9]+ { fprintf(outputfile,"Line %d: %s \t ———> Integer\n", line, yytext); }

'(\.|\.[^\\'])' { fprintf(outputfile,"Line %d: %s \t ———> Char Literal\n", line, yytext); }

\".*\" { fprintf(outputfile,"Line %d: %s \t ———> String Literal\n", line, yytext); }

[a-zA-Z_][a-zA-Z0-9_]* { fprintf(outputfile,"Line %d: %s \t ———> Identifier\n", line, yytext); }

\n      { line++; }
[ \t]+  { }
.        { }
%%
yywrap() {}
int main()
{
    printf("Enter cpp file name: "); scanf("%99s", filename);
    for (int i = 0; i < 100; i++)
    {
        if (filename[i] == '.')
        {
            if (filename[i + 1] == 'c' && filename[i + 2] == 'p' && filename[i + 3] == 'p')
            {
                break;
            }
        }
    }
}

```

```

        else
        {
            printf("file extention must be .cpp\n");
            return 1;
        }
    }
}

strcpy(ofilename, filename);
char *dot = strrchr(ofilename, '.'); // returns a pointer to the end of last oc
*dot = '\0';
strcat(ofilename, "_tokens.txt");
inputfile = fopen(filename, "r");
outputfile = fopen(ofilename, "w");
if (!inputfile || !outputfile)
{
    printf("Could not open file\n");
    return 1;
}

yyin = inputfile;
yylex();
fclose(inputfile);
fclose(outputfile);
}

```

Code Explanation

Our lexical analyzer uses Flex pattern matching capabilities to recognize various elements of C++ syntax:

1. Comments Handling:

- Single-line comments (`//`) are recognized and ignored
- Multi-line comments (`/* */`) are processed character by character to ensure proper line counting

2. Preprocessor Directives:

- Recognizes `#include` statements (both angle bracket and quotation formats)
- Identifies `#define` macros
- Identifies namespace declarations

3. Keywords and Types:

- Basic C++ types: int, float, char, double, void, bool, short, long, signed, unsigned, string
- Extended keyword set: if, else, while, for, do, switch, case, break, continue, return, goto, cout, cin, using, namespace, include, sizeof, struct, class, public, private, protected, virtual, static, const, new, delete, this, try, catch, throw, true, false, default, typedef, template, inline, friend, extern, enum, register, operator, mutable, volatile, nullptr, main, and, or, not

4. Operators:

- Arithmetic: +, -, *, /, %
- Assignment: =
- Comparison: <, >, ==, <=, >=, !=
- Logical: &&, ||, !
- Bitwise: &, |, ^, ~, <<, >>
- Compound operators: ++, --, +=, -=, *=, /=, %=, &=, |=, ^=, <<=, >>=

5. Separators:

- Brackets: {}, []
- Parentheses: ()
- Other separators: ;, ,

6. Literals:

- Integer literals: Sequences of digits
- Float literals: Digits with decimal point
- Character literals: Single characters in single quotes (with escape sequence support)
- String literals: Text in double quotes

7. Identifiers:

- Variable names, function names, etc. following C++ naming conventions (starts with letter or underscore, followed by letters, digits, or underscores)

8. Whitespace Handling:

- Spaces and tabs are ignored
- Newlines are counted to track line numbers

9. Output Management:

- Creates a separate output file with "_tokens.txt" suffix
- Each token is logged with line number, token value, and token type

The main function prompts the user for a C++ file, verifies its extension, opens it for reading, and passes it to the lexical analyzer. Results are written to an output file for further processing or analysis.

Sample Examples and Outputs

Example 1: Basic C++ Program

Input File: example1.cpp

```
#include <iostream>
using namespace std;

int main() {
    float x = 3.14;
    // This is a comment
    if (x > 0) {
        x = x + 1;
    }
    return 0;
}
```

Output:

```
Line 1: #include <iostream>    ———> Header
Line 2: using                ———> Keyword
```

Line 2: namespace → Keyword
Line 2: std → Identifier
Line 2: ; → Separator
Line 4: int → Type
Line 4: main → Keyword
Line 4: (→ Separator
Line 4:) → Separator
Line 4: { → Separator
Line 5: float → Type
Line 5: x → Identifier
Line 5: = → Operator
Line 5: 3.14 → Float
Line 5: ; → Separator
Line 7: if → Keyword
Line 7: (→ Separator
Line 7: x → Identifier
Line 7: > → Operator
Line 7: 0 → Integer
Line 7:) → Separator
Line 7: { → Separator
Line 8: x → Identifier
Line 8: = → Operator
Line 8: x → Identifier
Line 8: + → Operator
Line 8: 1 → Integer
Line 8: ; → Separator
Line 9: } → Separator
Line 10: return → Keyword
Line 10: 0 → Integer
Line 10: ; → Separator
Line 11: } → Separator

Example 2: More Complex C++ Program

Input File: example2.cpp

```
#include <iostream>
#include <string>
```

```

#define MAX

using namespace std;

/* This is a
   multi-line comment
   spanning several lines */

double calculateArea(double length, double width) {
    return length * width;
}

int main() {
    char grade = 'A';
    string message = "Hello World";
    int numbers[5] = {1, 2, 3, 4, 5};

    for (int i = 0; i < 5; i++) {
        if (numbers[i] <= 3) {
            // Skip small numbers
            continue;
        }
        cout << numbers[i] << endl;
    }

    return 0;
}

```

Output:

```

Line 1: #include <iostream>    ———> Header
Line 2: #include <string>     ———> Header
Line 3: #define              ———> Pre-process Directives
Line 3: MAX                  ———> Identifier
Line 3: 10                   ———> Integer
Line 5: using                ———> Keyword
Line 5: namespace           ———> Keyword
Line 5: std                  ———> Identifier

```

Line 5: ; → Separator
Line 11: double → Type
Line 11: calculateArea → Identifier
Line 11: (→ Separator
Line 11: double → Type
Line 11: length → Identifier
Line 11: , → Separator
Line 11: double → Type
Line 11: width → Identifier
Line 11:) → Separator
Line 12: { → Separator
Line 13: return → Keyword
Line 13: length → Identifier
Line 13: * → Operator
Line 13: width → Identifier
Line 13: ; → Separator
Line 14: } → Separator
Line 16: int → Type
Line 16: main → Keyword
Line 16: (→ Separator
Line 16:) → Separator
Line 17: { → Separator
Line 18: char → Type
Line 18: grade → Identifier
Line 18: = → Operator
Line 18: 'A' → Char Literal
Line 18: ; → Separator
Line 19: string → Type
Line 19: message → Identifier
Line 19: = → Operator
Line 19: "Hello World" → String Literal
Line 19: ; → Separator
Line 20: int → Type
Line 20: numbers → Identifier
Line 20: [→ Separator
Line 20: 5 → Integer
Line 20:] → Separator
Line 20: = → Operator

Line 20: { —→ Separator
Line 20: 1 —→ Integer
Line 20: , —→ Separator
Line 20: 2 —→ Integer
Line 20: , —→ Separator
Line 20: 3 —→ Integer
Line 20: , —→ Separator
Line 20: 4 —→ Integer
Line 20: , —→ Separator
Line 20: 5 —→ Integer
Line 20: } —→ Separator
Line 20: ; —→ Separator
Line 22: for —→ Keyword
Line 22: (—→ Separator
Line 22: int —→ Type
Line 22: i —→ Identifier
Line 22: = —→ Operator
Line 22: 0 —→ Integer
Line 22: ; —→ Separator
Line 22: i —→ Identifier
Line 22: < —→ Operator
Line 22: 5 —→ Integer
Line 22: ; —→ Separator
Line 22: i —→ Identifier
Line 22: ++ —→ Compound Operator
Line 22:) —→ Separator
Line 23: { —→ Separator
Line 24: if —→ Keyword
Line 24: (—→ Separator
Line 24: numbers —→ Identifier
Line 24: [—→ Separator
Line 24: i —→ Identifier
Line 24:] —→ Separator
Line 24: <= —→ Operator
Line 24: 3 —→ Integer
Line 24:) —→ Separator
Line 25: { —→ Separator
Line 27: continue —→ Keyword

```
Line 27: ;    ———> Separator
Line 28: }    ———> Separator
Line 29: cout  ———> Keyword
Line 29: <<   ———> Operator
Line 29: numbers ———> Identifier
Line 29: [    ———> Separator
Line 29: i    ———> Identifier
Line 29: ]    ———> Separator
Line 29: <<   ———> Operator
Line 29: endl  ———> Identifier
Line 29: ;    ———> Separator
Line 30: }    ———> Separator
Line 32: return ———> Keyword
Line 32: 0     ———> Integer
Line 32: ;    ———> Separator
Line 33: }    ———> Separator
```

Design Considerations and Implementation Details

Pattern Matching Strategy

Our lexical analyzer follows a "longest match" strategy, where it tries to match the longest possible pattern at each step. The patterns are arranged in a specific order to ensure correct token identification. For example, keywords and identifiers have patterns that could potentially overlap, so keywords are checked first.

Comment Handling

For handling multi-line comments, we implemented a character-by-character scanning approach to ensure accurate line counting. This approach also helps in properly identifying the end of a comment with the "*/" sequence.

Conclusion

This project successfully demonstrates the application of formal language theory concepts, specifically regular expressions, in creating a lexical analyzer for C++ programs. The analyzer can identify and categorize various tokens as required, providing a solid foundation for the front-end of a compiler.

The implementation effectively uses Flex to define patterns for different token types, manages state through line counting, and handles complex structures like comments. The output format provides clear information about each token's type, value, and line number, which can be useful for subsequent phases of compilation.

Through this project, we gained practical experience in formal language applications and a better understanding of the lexical analysis phase of compiler design.