

## CS3104 Compiler Lab, Autumn 2025 Assignment 3

Date: 22 August 2025

**Instructions:** Write your name and roll number on top of your source code as comments. Submit two files: your Lex file (`Assignment3.1`) and your C/C++ file (`procltx.c` or `procltx.cpp`). Submit your code at: <http://10.22.10.100/~abyaym/CS3104/submission/>

Building upon your experience from Assignments 1 and 2, you will now develop a **Lexical Analyzer using Lex/Flex** for your calculator language. This assignment introduces you to formal lexical analysis tools and helps you understand how tokens are recognized and processed in a systematic way.

### Problem Statement

Your task is to build a lexical analyzer using Lex/Flex that can tokenize the calculator language from Assignment 2, recognize different types of tokens (numbers, variables, operators, keywords), handle whitespace and comments, generate a token stream with appropriate token types, and provide meaningful error messages for invalid tokens.

### Token Specification

Your lexical analyzer should recognize the following tokens:

- **IDENTIFIER:** Variable names (single letters a-z, A-Z)
- **NUMBER:** Positive integers (one or more digits)
- **ASSIGN:** Assignment operator (=)
- **PLUS, MINUS, MULTIPLY, DIVIDE, MODULO:** Arithmetic operators (+, -, \*, /, %)
- **LPAREN, RPAREN:** Parentheses ((), )
- **QUIT, EXIT:** Exit keywords (quit, exit)
- **NEWLINE:** End of statement (\n)
- **WHITESPACE:** Spaces and tabs (to be ignored)
- **COMMENT:** C-style comments (/\* ... \*/ and // ...)

### Input Format and Example Output

```
1 Input:
2 x = 10 + y * 2 // This is a comment
3 z = (a - b) % 5
4 /* Multi-line
5    comment here */
6 quit
7
8 Expected Token Stream:
9 IDENTIFIER(x) ASSIGN NUMBER(10) PLUS IDENTIFIER(y) MULTIPLY NUMBER(2) NEWLINE
10 IDENTIFIER(z) ASSIGN LPAREN IDENTIFIER(a) MINUS IDENTIFIER(b) RPAREN MODULO
    NUMBER(5) NEWLINE
11 QUIT
```

Listing 1: Sample Input and Expected Token Output

## Core Requirements

- **Lex File (Assignment3.1):** Define regular expressions for all token types, handle comments and whitespace appropriately, provide token codes or symbolic names, and implement error handling for invalid characters.
- **C/C++ Integration File (procltx.c/procltx.cpp):** Process tokens generated by Lex, display token information in a readable format, maintain a simple token counter, and demonstrate integration between Lex and C/C++ code.
- **Token Display:** For each recognized token, display the token type and lexeme (actual text), handle special cases like keywords vs identifiers, and show line numbers for better debugging.

## Sample Implementation Structure

```

1 %{
2 #include <stdio.h>
3 #include <stdlib.h>
4 int line_num = 1;
5 int token_count = 0;
6 %}
7
8 %%
9 [a-zA-Z]          { /* Process IDENTIFIER */ }
10 [0-9]+            { /* Process NUMBER */ }
11 "="              { /* Process ASSIGN */ }
12 "+"              { /* Process PLUS */ }
13 /* Add more patterns here */
14 %%
15
16 int main() {
17     yylex();
18     return 0;
19 }
```

Listing 2: Basic Lex File Structure (Assignment3.1)

## Practice Problems

Complete the following exercises to test your lexical analyzer:

1. **Basic Tokenization:** Test with simple expressions like `a = 5 + 3`. Verify that each token is correctly identified and that the token count is accurate.
2. **Complex Expression:** Process the expression `result = (x + y * 2) / (z - 1) % 10` and ensure proper handling of all operators and parentheses.
3. **Comments and Whitespace:** Test with input containing both single-line and multi-line comments:

```

1 // Variable assignments
2 x = 10 /* this is x */
3 y = 20 + /* inline comment */ 5
4 /* Multi-line comment
5    spanning several lines */
6 z = x * y
7
```

Verify that comments are properly ignored and don't affect tokenization.

#### 4. Nested Comments and Edge Cases: Handle complex nested comment scenarios and edge cases:

```
1 a = 5 /* outer comment /* inner comment */ still outer */ + 3
2 b = 10 // comment with /* symbols inside single-line comment
3 c = /* comment at start */ 20 /* comment at end */
4 // /* Mixed comment styles */
5 d = 15 /* comment with // inside */ * 2
6
```

Ensure your lexer correctly handles nested comments or reports appropriate errors if your implementation doesn't support nesting.

5. **Large Number Handling and Boundary Cases:** Test with various number formats and boundary conditions:

[illegible]

Verify proper tokenization of large numbers and handle potential overflow scenarios gracefully.

### Advanced Features (Optional)

- **Line Number Tracking:** Maintain and display line numbers for each token.
- **Symbol Table Integration:** Extend your tokenizer to work with a simple symbol table.
- **Enhanced Error Reporting:** Provide detailed error messages with line and column numbers.
- **Token Statistics:** Generate statistics about token frequency and types.

## Implementation Guidelines

- Use Lex/Flex regular expressions efficiently and avoid conflicts between patterns.
- Handle edge cases like empty input or files with only comments.
- Ensure your C/C++ code properly integrates with the generated lexer.
- Test with various input combinations including erroneous input.
- Document your token definitions clearly in comments.

## Error Handling Requirements

Your lexical analyzer should handle the following error conditions:

- **Invalid Characters:** Report unknown symbols with line numbers.
- **Unterminated Comments:** Detect and report unclosed multi-line comments.
- **Malformed Numbers:** Handle cases like numbers with invalid characters.

## Compilation Instructions

To compile and run your lexical analyzer:

```
1 # For C implementation:
2 $ flex Assignment3.l
3 $ gcc lex.yy.c procltx.c -o lexer -lfl
4 $ ./lexer < input.txt
5
6 # For C++ implementation:
7 $ flex Assignment3.l
8 $ g++ lex.yy.c procltx.cpp -o lexer -lfl
9 $ ./lexer < input.txt
```

This assignment introduces you to the power of formal lexical analysis tools. The techniques you learn here will be essential for building more sophisticated language processors and understanding how modern compilers work. The systematic approach of Lex/Flex will help you appreciate the advantages of tool-based development over manual parsing.

**Good luck and enjoy exploring lexical analysis!**