

Lexical Analyzer Report

Ahmed Badr 900202868

Youssef Khaled 900213467

October 18, 2024

Course Information

Course Name: CSCE4101 - Compiler Design (Fall 2024)

Instructor: Dr. Ahmed Rafea

Introduction

This report presents the development of a lexical analyzer using **Lex/Flex**. The analyzer identifies keywords, identifiers, numbers, operators, and comments while detecting and reporting specific errors. The errors include unclosed comments, invalid identifiers, and malformed numbers. This document discusses the decisions made during development, test cases used, and the program output.

Decisions Made

The following decisions were made during the implementation of the lexical analyzer:

1. **Handling Unclosed Comments:** If the input file contains an unclosed comment (i.e., `'/` without `'`), the program detects this at EOF and reports the error with the relevant line number and position.
2. **Handling Invalid Characters:** If a character not belonging to the language's alphabet appears, the scanner reports it as an invalid character, indicating its position in the input.
3. **Premature Token Termination:** The program detects if a token is prematurely interrupted. For example:

- In identifiers: If an invalid character appears (e.g., ‘invalid_*id’), it reports ”Invalid identifier format.”
 - In numbers: If an invalid letter follows an exponential notation (e.g., ‘12.3eX’), it reports ”Invalid number format.”
4. **Number-Identifier Handling:** The scanner treats ‘12ab’ as two tokens: ‘NUM’ and ‘ID’. Such sequences are not flagged by the lexer but will be handled by the parser.

Instructions to Run the Program

1. Save the code in a file named `scanner.l`.
2. Generate the scanner using the command:

```
flex scanner.l
```

3. Compile the generated C code:

```
gcc lex.yy.c -o scanner -lfl
```

4. Create a test input file named `input.txt` with the content shown below.
5. Run the scanner:

```
./scanner input.txt
```

Test Input

```
/* Test identifiers */
valid_id
another@id
third.id
invalid_*id
x_1
y.2
z@3
```

```
/* Test numbers */
123
45.67
```

2.3E-4
12.3eX
89.1e+2

```
/* Test keywords */  
if then else  
IF THEN ELSE  
If Then Else  
  
/* Test invalid cases */  
abc_*  
123abc
```

Program Output

```
saliert@saliert:~/Downloads/Compiler Design/scanner$ flex scanner.l  
saliert@saliert:~/Downloads/Compiler Design/scanner$ gcc lex.yy.c -o scanner -ll  
saliert@saliert:~/Downloads/Compiler Design/scanner$ ./scanner input.txt  
TOKEN: ID (valid_id) at line 2, position 1  
TOKEN: ID (another@id) at line 3, position 1  
TOKEN: ID (third.id) at line 4, position 1  
Error: Invalid identifier format at line 5, position 1: 'invalid_*id'  
TOKEN: ID (x_1) at line 6, position 1  
TOKEN: ID (y.2) at line 7, position 1  
TOKEN: ID (z@3) at line 8, position 1  
TOKEN: NUM (123) at line 11, position 1  
TOKEN: NUM (45.67) at line 12, position 1  
TOKEN: NUM (2.3E-4) at line 13, position 1  
Error: Invalid number format at line 14, position 1: '12.3eX'  
TOKEN: NUM (89.1e+2) at line 15, position 1  
TOKEN: IF at line 18, position 1  
TOKEN: THEN at line 18, position 4  
TOKEN: ELSE at line 18, position 9  
TOKEN: IF at line 19, position 1  
TOKEN: THEN at line 19, position 4  
TOKEN: ELSE at line 19, position 9  
TOKEN: IF at line 20, position 1  
TOKEN: THEN at line 20, position 4  
TOKEN: ELSE at line 20, position 9  
Error: Unclosed comment at end of file at line 25, position 1: ''
```

Appendix A: Code

Listing 1: Lexical Analyzer Code: `scanner.1`

```
1 %{
2 #include <stdio.h>
3 #include <string.h>
4 #include <ctype.h>
5
6 /* Token definitions */
7 #define ID_TOKEN 1
8 #define NUM_TOKEN 2
9 #define IF_TOKEN 3
10 #define THEN_TOKEN 4
11 #define ELSE_TOKEN 5
12
13 int line_number = 1;
14 int char_position = 1;
15 int token_start_pos = 1;
16 void error_message(char* msg);
17 %}
18
19 %option noyywrap
20 %option caseless
21
22 /* States */
23 %x COMMENT
24
25 /* Regular Definitions */
26 LETTER      [A-Za-z]
27 DIGIT       [0-9]
28 WHITESPACE  [ \t\r]
29 SPECIAL     [.\@_]
30 EXPONENT    [Ee] [+ -]? {DIGIT}+
31
32 ID          {LETTER}({LETTER}|{DIGIT})*({SPECIAL}({LETTER}|{
    DIGIT}))?)?
33 INVALID_ID  {LETTER}({LETTER}|{DIGIT})*{SPECIAL}?[^A-Za-z0-9.
    @_ \n \t]+({LETTER}|{DIGIT})*
34 NUM         {DIGIT}+(\.{DIGIT}+)?({EXPONENT})?
35 INVALID_NUM {DIGIT}+(\.{DIGIT}+)?({EXPONENT})?[A-Za-z]+
36
```

```
37 %%
38
39 "/*"          {
40     token_start_pos = char_position;
41     char_position += yyleng;
42     BEGIN(COMMENT);
43 }
44
45 <COMMENT>"*/"  {
46     char_position += yyleng;
47     BEGIN(INITIAL);
48 }
49
50 <COMMENT>\n    {
51     line_number++;
52     char_position = 1;
53 }
54
55 <COMMENT>.      {
56     char_position += yyleng;
57 }
58
59 <COMMENT><<EOF>> {
60     error_message("Error: Unclosed comment at end of file");
61     return 0;
62 }
63
64 {WHITESPACE}    {
65     char_position += yyleng;
66 }
67
68 \n              {
69     line_number++;
70     char_position = 1;
71 }
72
73 if              {
74     token_start_pos = char_position;
75     char_position += yyleng;
76     printf("TOKEN: IF at line %d, position %d\n", line_number,
77           token_start_pos);
```

```
77     return IF_TOKEN;
78 }
79
80 then      {
81     token_start_pos = char_position;
82     char_position += yyleng;
83     printf("TOKEN: THEN at line %d, position %d\n", line_number,
84           token_start_pos);
85     return THEN_TOKEN;
86 }
87
88 else      {
89     token_start_pos = char_position;
90     char_position += yyleng;
91     printf("TOKEN: ELSE at line %d, position %d\n", line_number,
92           token_start_pos);
93     return ELSE_TOKEN;
94 }
95
96 {ID}      {
97     token_start_pos = char_position;
98     char_position += yyleng;
99     printf("TOKEN: ID (%s) at line %d, position %d\n", yytext,
100           line_number, token_start_pos);
101     return ID_TOKEN;
102 }
103
104 {INVALIDID_ID}  {
105     token_start_pos = char_position;
106     error_message("Error: Invalid identifier format");
107     char_position += yyleng;
108     /* Continue scanning after the invalid identifier */
109 }
110
111 {NUM}      {
112     token_start_pos = char_position;
113     char_position += yyleng;
114     printf("TOKEN: NUM (%s) at line %d, position %d\n", yytext,
115           line_number, token_start_pos);
116     return NUM_TOKEN;
117 }
```

```
114
115 {INVALID_NUM}    {
116     token_start_pos = char_position;
117     error_message("Error: Invalid number format");
118     char_position += yyleng;
119     /* Continue scanning after the invalid number */
120 }
121
122 [+\\-*/=<>!]      {
123     token_start_pos = char_position;
124     char_position += yyleng;
125     printf("TOKEN: OPERATOR (%s) at line %d, position %d\\n",
126           yytext, line_number, token_start_pos);
127     return yytext[0];
128 }
129 .                  {
130     token_start_pos = char_position;
131     error_message("Error: Invalid character");
132     char_position += yyleng;
133     /* Consume the invalid character */
134 }
135
136 %%
137
138 void error_message(char* msg) {
139     fprintf(stderr, "%s at line %d, position %d: '%s'\\n",
140           msg, line_number, token_start_pos, yytext);
141 }
142
143 int main(int argc, char** argv) {
144     if (argc < 2) {
145         printf("Usage: %s input_file\\n", argv[0]);
146         return 1;
147     }
148
149     FILE* input = fopen(argv[1], "r");
150     if (!input) {
151         fprintf(stderr, "Error: Cannot open input file '%s'\\n",
152               argv[1]);
153         return 1;
154     }
155 }
```

```
153     }
154
155     yyin = input;
156
157     while (yylex()) {
158         // Token processing is handled in the rules above
159     }
160
161     fclose(input);
162     return 0;
163 }
```