

Fr. Conceicao Rodrigues College of Engineering
Department of Computer Engineering

Academic Term : Jan-May 2024 - 25

Class : T.E. (Computer - A)

Subject Name : System Programming and Compiler Construction Subject

Code : (CPC601)

Practical No:	02
Title:	Lexical Analyzer
Date of Performance:	11/02/2025
Date of Submission:	19/02/2025
Roll No:	9913
Name of the Student:	Mark Lopes

Evaluation:

Sr. No	Rubric	Grade
1	Time Line (2)	
2	Output(3)	
3	Code optimization (2)	
4	Postlab (3)	

Signature of the Teacher :

Experiment No 2

Aim: Write a program to implement Lexical analyzer

Learning Objective: Converting a sequence of characters into a sequence of tokens.

Theory:

THE ROLE OF LEXICAL ANALYZER

The lexical analyzer is the first phase of a compiler. Its main task is to read the input characters and produce as output a sequence of tokens that the parser uses for syntax analysis. Upon receiving a “get next token” command from the parser, the lexical analyzer reads input characters until it can identify the next token.

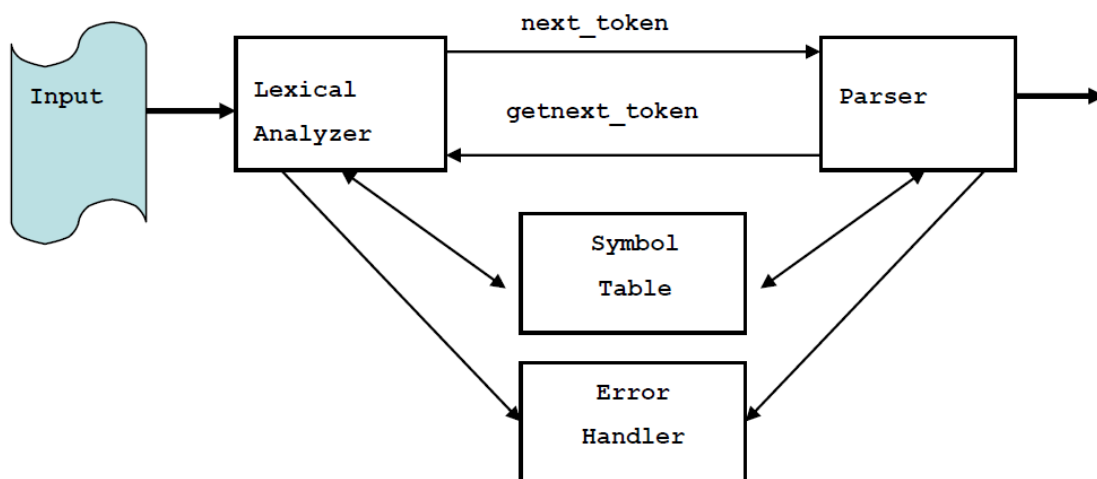


Figure 4.1 Interaction of Lexical Analyzer with Parser

Since the lexical analyzer is the part of the compiler that reads the source text, it may also perform certain secondary tasks at the user interface. One such task is stripping out from the source program comments and white spaces in the form of blank, tab, and new line characters. Another is correlating error messages from the compiler with the source program. Sometimes lexical analyzers are divided into a cascade of two phases first called “scanning” and the second “lexical analysis”. The scanner is responsible for doing simple tasks, while the lexical analyzer proper does the more complex operations.

Implementation Details

1. Read the high level language as source program
2. Convert source program in to categories of tokens such as Identifiers, Keywords, Constants, Literals and Operators.

Test cases:

1. Input undefined token

Code:

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <string.h>
4  #include <ctype.h>
5
6  #define MAX_TOKEN_LEN 100
7  #define MAX_TOKENS 1000
8
9  // Token type arrays
10 const char* keywords[] = {
11     "auto", "break", "case", "char", "const", "continue", "default", "do",
12     "double", "else", "enum", "extern", "float", "for", "goto", "if",
13     "int", "long", "register", "return", "short", "signed", "sizeof", "static",
14     "struct", "switch", "typedef", "union", "unsigned", "void", "volatile", "while"
15 };
16
17 const char operators[] = {
18     '+', '-', '*', '/', '%', // arithmetic
19     '=', '<', '>', '!', // relational
20     '&', '|', '^', '~', // bitwise
21     ':', '?' // others
22 };
23
24 const char delimiters[] = {
25     '(', ')', '{', '}', '[', ']',
26     ';', ',', '.', ' '
27 };
28
29 const char* data_types[] = {
30     "int", "char", "float", "double", "void",
31     "long", "short", "signed", "unsigned"
32 };
33
34 // Arrays to store unique tokens
35 char identifiers[MAX_TOKENS][MAX_TOKEN_LEN];
36 char numbers[MAX_TOKENS][MAX_TOKEN_LEN];
37 char operators_found[MAX_TOKENS];
38 char delimiters_found[MAX_TOKENS];
39
40 int id_count = 0, num_count = 0, op_count = 0, del_count = 0;
41
42 // Helper function to check if a string is in an array
43 int exists_in_array(char arr[MAX_TOKENS][MAX_TOKEN_LEN], int count, char* str) {
44     for (int i = 0; i < count; i++) {
45         if (strcmp(arr[i], str) == 0) return 1;
46     }
47     return 0;
48 }
```

```

50 // Helper function to check if a character is in an array
51 int exists_in_char_array(char arr[MAX_TOKENS], int count, char ch) {
52     for (int i = 0; i < count; i++) {
53         if (arr[i] == ch) return 1;
54     }
55     return 0;
56 }
57
58 // Check if token is a keyword
59 int is_keyword(char* str) {
60     for (int i = 0; i < sizeof(keywords) / sizeof(char*); i++) {
61         if (strcmp(str, keywords[i]) == 0) return 1;
62     }
63     return 0;
64 }
65
66 // Check if token is a data type
67 int is_data_type(char* str) {
68     for (int i = 0; i < sizeof(data_types) / sizeof(char*); i++) {
69         if (strcmp(str, data_types[i]) == 0) return 1;
70     }
71     return 0;
72 }

```

```

74 // Check if character is an operator
75 int is_operator(char ch) {
76     for (int i = 0; i < sizeof(operators); i++) {
77         if (ch == operators[i]) return 1;
78     }
79     return 0;
80 }
81
82 // Check if character is a delimiter
83 int is_delimiter(char ch) {
84     for (int i = 0; i < sizeof(delimiters); i++) {
85         if (ch == delimiters[i]) return 1;
86     }
87     return 0;
88 }
89
90 // Helper function to print a line separator
91 void print_separator(int width, char style) {
92     printf("+");
93     for (int i = 0; i < width-1; i++) {
94         printf("%c", style);
95     }
96     printf("+");
97     for (int i = 0; i < width-1; i++) {
98         printf("%c", style);
99     }

```

```

100     printf("+\n");
101 }
102
103 // Helper function to print wrapped text in a column
104 void print_wrapped_text(const char* text, int width) {
105     int len = strlen(text);
106     int pos = 0;
107     int chars_remaining = len;
108     int first_line = 1;
109
110     while (chars_remaining > 0) {
111         int chars_to_print = chars_remaining;
112         if (chars_to_print > width - 2) {
113             chars_to_print = width - 2;
114             // Find last space before width limit
115             while (chars_to_print > 0 && text[pos + chars_to_print] != ' ' && text[pos + chars_to_print] != ',') {
116                 chars_to_print--;
117             }
118             if (chars_to_print == 0) chars_to_print = width - 2; // If no space found, force break
119         }
120

```

```

121         if (first_line) {
122             first_line = 0;
123         } else {
124             printf("| %*s | ", width-2, "");
125         }
126
127         printf("%-*.*s%s\n",
128             chars_to_print, chars_to_print,
129             text + pos,
130             width - chars_to_print - 2, "");
131
132         pos += chars_to_print;
133         while (pos < len && (text[pos] == ' ' || text[pos] == ',')) pos++;
134         chars_remaining = len - pos;
135     }
136 }
137
138 int main() {
139     printf("#####\n");
140     printf("##### LEXICAL ANALYSIS #####\n");
141     printf("#####\n");
142     FILE *file_ptr;
143     char ch, token[MAX_TOKEN_LEN];
144     int i = 0;

```

```

145
146     file_ptr = fopen("fibonacci.c", "r");
147     if (NULL == file_ptr) {
148         printf("File can't be opened\n");
149         return EXIT_FAILURE;
150     }
151
152     while ((ch = fgetc(file_ptr)) != EOF) {
153         if (isspace(ch)) continue;
154
155         if (isalpha(ch) || ch == '_') {
156             i = 0;
157             token[i++] = ch;
158             while ((ch = fgetc(file_ptr)) != EOF && (isalnum(ch) || ch == '_')) {
159                 token[i++] = ch;
160             }
161             token[i] = '\0';
162             ungetc(ch, file_ptr);
163
164             if (!is_keyword(token) && !is_data_type(token) && !exists_in_array(identifiers, id_count, token)) {
165                 strcpy(identifiers[id_count++], token);
166             }
167         }
168         else if (isdigit(ch)) {
169             i = 0;

```

```

170         token[i++] = ch;
171         while ((ch = fgetc(file_ptr)) != EOF && (isdigit(ch) || ch == '.')) {
172             token[i++] = ch;
173         }
174         token[i] = '\0';
175         ungetc(ch, file_ptr);
176
177         if (!exists_in_array(numbers, num_count, token)) {
178             strcpy(numbers[num_count++], token);
179         }
180     }
181     else if (is_operator(ch)) {
182         if (!exists_in_char_array(operators_found, op_count, ch)) {
183             operators_found[op_count++] = ch;
184         }
185     }
186     else if (is_delimiter(ch)) {
187         if (!exists_in_char_array(delimiters_found, del_count, ch)) {
188             delimiters_found[del_count++] = ch;
189         }
190     }
191 }
192

```

```

193     fclose(file_ptr);
194
195     int column_width = 50;
196
197     // Print table header with dashes
198     print_separator(column_width, '-');
199     printf("| %-*s | %-*s |\n", column_width-2, "Category", column_width-2, "Value");
200     print_separator(column_width, '_'); // Using underscore after Category
201
202     // Print IDENTIFIER
203     printf("| %-*s | ", column_width-2, "IDENTIFIER");
204     char id_buffer[MAX_TOKENS * MAX_TOKEN_LEN] = {0};
205     for (int j = 0; j < id_count; j++) {
206         strcat(id_buffer, identifiers[j]);
207         if (j < id_count - 1) strcat(id_buffer, ", ");
208     }
209     print_wrapped_text(id_buffer, column_width);
210     print_separator(column_width, '_'); // Add separator after each category
211
212     // Print OPERATOR
213     printf("| %-*s | ", column_width-2, "OPERATOR");
214     char op_buffer[MAX_TOKENS * 2] = {0};
215     for (int j = 0; j < op_count; j++) {
216         char temp[3] = {operators_found[j], 0};
217         strcat(op_buffer, temp);

```



```

228     if (j < num_count - 1) strcat(num_buffer, ", ");
229 }
230 print_wrapped_text(num_buffer, column_width);
231 print_separator(column_width, '_'); // Add separator after each category
232
233 // Print DELIMITER
234 printf("| %-s |", column_width-2, "DELIMITER");
235 char del_buffer[MAX_TOKENS * 2] = {0};
236 for (int j = 0; j < del_count; j++) {
237     char temp[3] = {delimiters_found[j], 0};
238     strcat(del_buffer, temp);
239     if (j < del_count - 1) strcat(del_buffer, ", ");
240 }
241 print_wrapped_text(del_buffer, column_width);
242 print_separator(column_width, '-'); // Using dash for final separator
243
244 // Calculate and print total tokens
245 int keyword_count = sizeof(keywords) / sizeof(keywords[0]);
246 int datatype_count = sizeof(data_types) / sizeof(data_types[0]);
247 int total_tokens = id_count + num_count + op_count + del_count + keyword_count + datatype_count;
248 printf("\nTotal number of tokens: %d\n", total_tokens);
249
250 return 0;
251 }

```

Output:

```

#####
#####          LEXICAL ANALYSIS          #####
#####
+-----+-----+
| Category | Value |
+-----+-----+
| IDENTIFIER | include, stdio, h, printf, n, If, the, number, |
| | of, terms, is, smaller, than, printf, Invalid, |
| | Number, First, two, series, prev1, prev2, loop, |
| | that, prints, fibonacci, i, Print, current, term |
| | and, update, previous, curr, d, main, Printing, |
| | first |
+-----+-----+
| OPERATOR | <, >, /, =, +, % |
+-----+-----+
| NUMBER | 1, 0, 2, 9 |
+-----+-----+
| DELIMITER | ., (, ), {, :, }, , |
+-----+-----+

Total number of tokens: 95

```

Post Lab Questions:

1. Explain the role of automata theory in compiler design.
2. What are the errors that are handled by Lexical analysis phase?

DATE	
PAGE NO.	

19/02/25

Q.1 Role of automata theory in compiler design.
→ Automata theory is fundamental in compiler design, enabling efficient parsing and processing of program languages.

1. Lexical analysis:- The compiler first phase uses finite Automata to recognize tokens, keywords, identifiers, generator. Regular expression define token pattern which are converted into DFA for fast scanning.

2. Syntax analysis:- The syntax analyzer ensures the program follows the language. Grammar is context free grammar. Parser uses PDA for nested structures.

3. Semantic analysis and code generation:- Automata theory helps build abstract syntax tree ensure type correctness and generate intermediate code.

4. Optimization of code generation:- Automata helps in control flow analysis register allocation and instruction scheduling for efficient execution.

5. Runtime execution:- ~~the~~ virtual machines use automata for interpreting and executing compiled code.

Q.2 What are the errors that are handled by lexical analysis phase?

-
1. Invalid character :- Unrecognized symbols or illegal characters that do not match any defined token pattern.
 2. Unterminated string or comments :- strings that start but don't end.
 3. Improper number format :- Numerical literals with invalid formats like floating point number with multiple decimal points or characters.
 4. Identifier length :- Identifiers that exceed max. allocated allowed target length set by the language specification.
 5. Illegal escape sequence :- Improper use of escape sequence in string.