

Name: Digvijay Jondhale

Roll No: PC-32

PRN: 1032201770

SSCD LCA 05

Aim: Generate lexical analyzer for C language using Lex tool.

Theory:

- Token lexeme and pattern:
- > Token: In the context of the programming language and lexical analysis, a token is a fundamental unit of a language's syntax. Tokens represent specific elements or symbols in the source code, such as keywords (e.g., "if", "while"), identifiers (e.g., variable names), literals (e.g., numbers or strings), and operators (e.g., "+", "="). Tokens are used by the parser to understand the structure of the code and perform syntactic analysis.
- > lexeme: A lexeme is the actual sequence of characters in the source code that corresponds to a specific token. It is the concrete representation of a token in the source code. For example, in the statement "int x=42";, the lexemes for the tokens are "int", "x", "=", and "42".
- > Pattern: A pattern is a description or regular expression that defines the structure or format of lexemes. It specifies the rules that lexemes must follow to be recognized as a particular token. For example, a pattern for recognizing integer literals in a programming language might be "\d+".

- Use of Regular Expression (RE) in specifying lexical structure of a language:
 - > Define Patterns: Programmers use regular expressions to define patterns for various tokens in the language. For example, a pattern for identifying identifiers may be "[a-zA-Z_][a-zA-Z0-9_]*", which matches valid variable names.
 - > Tokenization: The lexer or lexical analyzer processes the source code character by character and tries to match the input against the defined regular expressions. When a pattern matches a portion of the input, it generates a token with the associated lexeme.
 - > Generating Tokens: Tokens are created based on the patterns that matched. Each token contains the lexeme (the matched text) and the token type (e.g., identifier, keyword, number). These tokens are then passed on to the parser for further syntactic analysis.
 - > Handling Ambiguity: Regular expressions can help handle ambiguity by allowing you to specify rules for resolving situations where the input could match multiple patterns. This is often done by giving priority to the first matching pattern encountered or using other disambiguation rules.
- Format of lex Specification and Execution Steps of a lex file (4, 2):

Lexical analyzers for programming languages are often generated using tools like lex (or its open-source counterpart, flex). These tools use a specification file with a specific format, typically with the ".l" extension. Here's an overview of the

◦ Lex specification file format (*.L):

- > Definitions: You can define regular expressions and macros at the beginning of the file. These definitions are enclosed in "%{" and "%}" delimiters.
- > Rules: The main part of the file consists of rules that specify regular expressions and associated actions. Rules are written in the format:

regex action;

where "regex" is a regular expression pattern, and "action" is the code or action to be executed when the pattern is matched.

- > User Code: You can include C code within the specification file, usually enclosed in "%{" and "%}" delimiters, for custom actions or additional code.

◦ Execution Steps:

- 1) Write the Lex Specification: Create a ".L" file that defines the lexical structure of the language using regular expressions and associated actions.
- 2) Generate Lexer Code: Use a tool like Lex or Flex to generate C code for the lexical analyzer based on the ".L" file. This code includes functions for tokenizing the input source code.
- 3) Compile and Link: Compile the generated C code and link it with your parser or the rest of your compiler or interpreter.
- 4) Execute: Run the resulting executable, which will tokenize the input source code according to the defined lexical rules and generate a sequence of tokens that can be used for further parsing and analysis.