

Lexical Analysis

Lexical analysis is the first phase of a compiler. It takes modified source code from language preprocessors that are written in the form of sentences. The lexical analyzer breaks these syntaxes into a series of tokens, by removing any whitespace or comments in the source code.

If the lexical analyzer finds a token invalid, it generates an error. The lexical analyzer works closely with the syntax analyzer. It reads character streams from the source code, checks for legal tokens, and passes the data to the syntax analyzer when it demands.

Tokens

Lexemes are said to be a sequence of characters (alphanumeric) in a token. There are some predefined rules for every lexeme to be identified as a valid token. These rules are defined by grammar rules, by means of a pattern. A pattern explains what can be a token, and these patterns are defined by means of regular expressions.

In programming language, keywords, constants, identifiers, strings, numbers, operators and punctuations symbols can be considered as tokens.

For example, in C language, the variable declaration line

```
int value = 100;
```

contains the tokens:

```
int (keyword), value (identifier), = (operator), 100 (constant) and ; (symbol).
```

Specifications of Tokens

Let us understand how the language theory undertakes the following terms:

Alphabets

Any finite set of symbols {0,1} is a set of binary alphabets,

{0,1,2,3,4,5,6,7,8,9,A,B,C,D,E,F} is a set of Hexadecimal alphabets, {a-z, A-Z} is a set of English language alphabets.

Strings

Any finite sequence of alphabets (characters) is called a string. Length of the string is the total number of occurrence of alphabets, e.g., the length of the string java is 4 and is denoted by $|java| = 4$. A string having no alphabets, i.e. a string of zero length is known as an empty string and is denoted by ϵ (epsilon).

Special symbols

A typical high-level language contains the following symbols:-

Arithmetic Symbols	Addition(+), Subtraction(-), Modulo(%), Multiplication(*), Division(/)
Punctuation	Comma(,), Semicolon(;), Dot(.), Arrow(->)
Assignment	=
Special Assignment	+=, /=, *=, -=
Comparison	==, !=, <, <=, >, >=
Preprocessor	#
Location Specifier	&
Logical	&, &&, , , !
Shift Operator	>>, >>>, <<, <<<

Language

A language is considered as a finite set of strings over some finite set of alphabets. Computer languages are considered as finite sets, and mathematically set operations can be performed on them. Finite languages can be described by means of regular expressions.

Regular Expressions

The lexical analyzer needs to scan and identify only a finite set of valid string/token/lexeme that belong to the language in hand. It searches for the pattern defined by the language rules.

Regular expressions have the capability to express finite languages by defining a pattern for finite strings of symbols. The grammar defined by regular expressions is known as **regular grammar**. The language defined by regular grammar is known as **regular language**.

Regular expression is an important notation for specifying patterns. Each pattern matches a set of strings, so regular expressions serve as names for a set of strings. Programming language tokens can be described by regular languages. The specification of regular expressions is an example of a recursive definition. Regular languages are easy to understand and have efficient implementation.

There are a number of algebraic laws that are obeyed by regular expressions, which can be used to manipulate regular expressions into equivalent forms.

Operations

The various operations on languages are:

- Union of two languages L and M is written as
$$L \cup M = \{s \mid s \text{ is in } L \text{ or } s \text{ is in } M\}$$
- Concatenation of two languages L and M is written as
$$LM = \{st \mid s \text{ is in } L \text{ and } t \text{ is in } M\}$$
- The Kleene Closure of a language L is written as
$$L^* = \text{Zero or more occurrence of language } L.$$

Notations

If r and s are regular expressions denoting the languages L(r) and L(s), then

- **Union** : (r) | (s) is a regular expression denoting L(r) \cup L(s)
- **Concatenation** : (r)(s) is a regular expression denoting L(r)L(s)

- **Kleene closure** : $(r)^*$ is a regular expression denoting $(L(r))^*$
- (r) is a regular expression denoting $L(r)$

Precedence and Associativity

- $*$, concatenation $(.)$, and $|$ (pipe sign) are left associative
- $*$ has the highest precedence
- Concatenation $(.)$ has the second highest precedence.
- $|$ (pipe sign) has the lowest precedence of all.

Representing valid tokens of a language in regular expression

If x is a regular expression, then:

x^* means zero or more occurrence of x .

i.e., it can generate $\{ e, x, xx, xxx, xxxx, \dots \}$

x^+ means one or more occurrence of x .

i.e., it can generate $\{ x, xx, xxx, xxxx \dots \}$ or $x.x^*$

$x?$ means at most one occurrence of x

i.e., it can generate either $\{x\}$ or $\{e\}$.

$[a-z]$ is all lower-case alphabets of English language.

$[A-Z]$ is all upper-case alphabets of English language.

$[0-9]$ is all natural digits used in mathematics.

Representation occurrence of symbols using regular expressions

letter = $[a - z]$ or $[A - Z]$

digit = $0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$ or $[0-9]$

sign = $[+ \mid -]$

Representation of language tokens using regular expressions

Decimal = $(\text{sign})^?(\text{digit})^+$

Identifier = $(\text{letter})(\text{letter} \mid \text{digit})^*$

The only problem left with the lexical analyzer is how to verify the validity of a regular expression used in specifying the patterns of keywords of a language. A well-accepted solution is to use finite automata for verification.

Finite Automata

Finite automata is a state machine that takes a string of symbols as input and changes its state accordingly. Finite automata is a recognizer for regular expressions. When a regular expression string is fed into finite automata, it changes its state for each literal. If the input string is successfully processed and the automata reaches its final state, it is accepted, i.e., the string just fed was said to be a valid token of the language in hand.

The mathematical model of finite automata consists of:

- Finite set of states (Q)
- Finite set of input symbols (Σ)
- One Start state (q_0)
- Set of final states (q_f)
- Transition function (δ)

The transition function (δ) maps the finite set of state (Q) to a finite set of input symbols (Σ), $Q \times \Sigma \rightarrow Q$

Finite Automata Construction

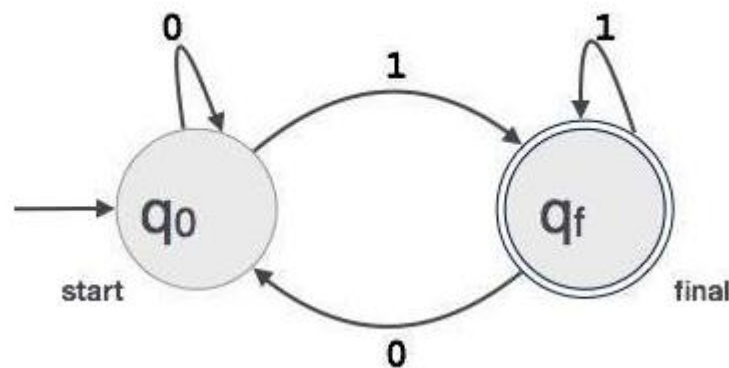
Let $L(r)$ be a regular language recognized by some finite automata (FA).

- **States** : States of FA are represented by circles. State names are of the state is written inside the circle.
- **Start state** : The state from where the automata starts, is known as start state. Start state has an arrow pointed towards it.
- **Intermediate states** : All intermediate states has at least two arrows; one pointing to and another pointing out from them.
- **Final state** : If the input string is successfully parsed, the automata is expected to be in this state. Final state is represented by double circles. It may have any odd

number of arrows pointing to it and even number of arrows pointing out from it. The number of odd arrows are one greater than even, i.e. **odd = even+1**.

- **Transition** : The transition from one state to another state happens when a desired symbol in the input is found. Upon transition, automata can either move to next state or stay in the same state. Movement from one state to another is shown as a directed arrow, where the arrows points to the destination state. If automata stays on the same state, an arrow pointing from a state to itself is drawn.

Example : We assume FA accepts any three digit binary value ending in digit 1. FA = $\{Q(q_0, q_f), \Sigma(0,1), q_0, q_f, \delta\}$



Longest Match Rule

When the lexical analyzer read the source-code, it scans the code letter by letter; and when a whitespace, operator symbol, or special symbols occurs, it decides that a word is completed.

For example:

```
int intvalue;
```

While scanning both lexemes till 'int', the lexical analyzer cannot determine whether it is a keyword *int* or the initials of identifier int value.

The Longest Match Rule states that the lexeme scanned should be determined based on the longest match among all the tokens available.

The lexical analyzer also follows **rule priority** where a reserved word, e.g., a keyword, of a language is given priority over user input. That is, if the lexical analyzer finds a lexeme that matches with any existing reserved word, it should generate an error.

Applications of Lexical Analysis

- **Compilers:** Lexical analysis is an important part of the compilation process, as it converts the source code of a program into a stream of tokens that can be more easily processed by the compiler.
- **Interpreters:** Lexical analysis is also used in interpreters, which execute a program directly from its source code without the need for compilation.
- **Text editors:** Many text editors use lexical analysis to highlight keywords and other elements of the source code in different colors, making it easier for programmers to read and understand the code.
- **Code analysis tools:** Lexical analysis is used by tools that analyze the source code of a program for errors, security vulnerabilities, and other issues.
- **Natural language processing:** Lexical analysis is also used in natural language processing (NLP) to break down natural language text into individual words and phrases that can be more easily processed by NLP algorithms.
- **Information retrieval:** Lexical analysis is used in information retrieval systems, such as search engines, to index and search for documents based on the words they contain.