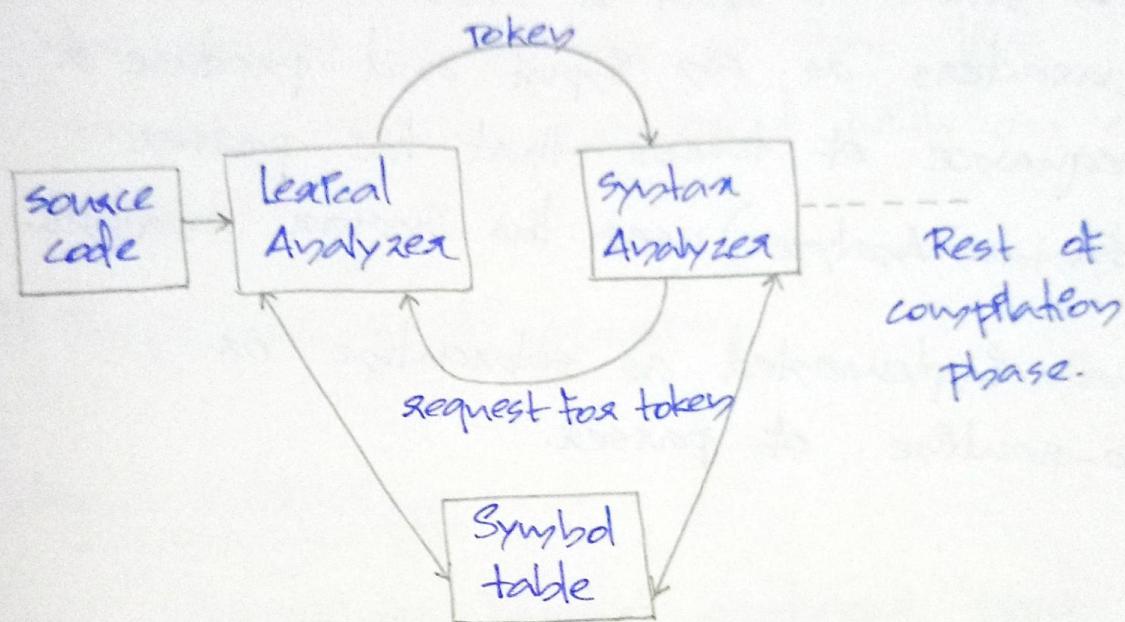


Module : 2

Lexical Analysis

lexical analyzer is the first phase of the compilation also known as scanner. It converts the high level input programs into a sequence of tokens.



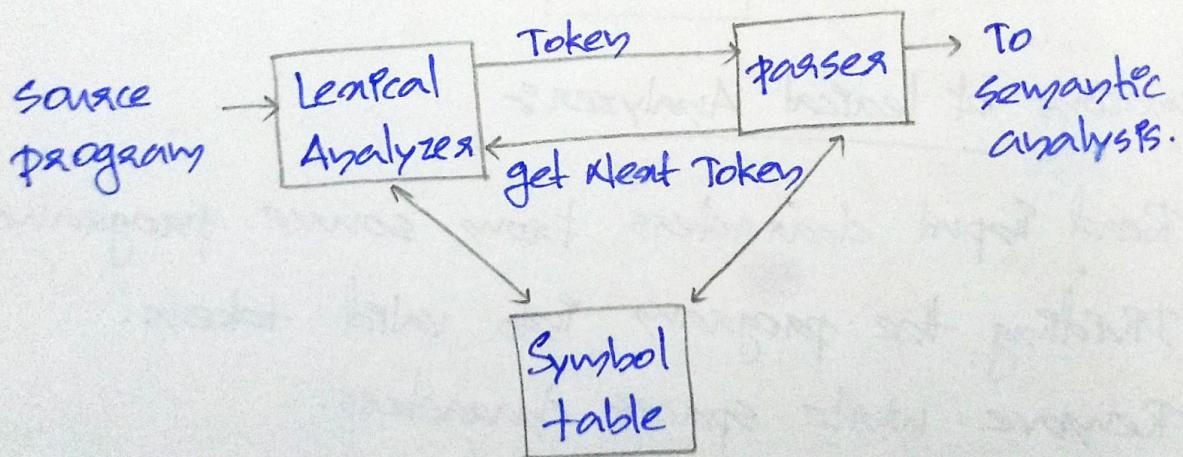
Functions of lexical Analyzer

- ① Read input characters from source programs.
- ② Dividing the program into valid tokens.
- ③ Remove white spaces characters.
- ④ Remove comments (//, #).
- ⑤ Helps to identify tokens into the symbol table.

⑥ It generates lexical errors.

* Role of lexical Analyzer:

- ① Lexical analyzer is the first phase of a compiler.
- ② Main task is to read a stream of characters as an input and produce a sequence of tokens that the parser (Syntax Analyzer) uses for syntax analysis.
- ③ Usually implemented as subroutine or co-routine of parser.



- ④ Separation of the input source code into tokens.
 - ⑤ stripping out the comments and unnecessary white spaces by the form of blank, tab and newline from the source code.
 - ⑥ keeping track of line numbers while scanning the newline characters so that line numbers can be associated with an error message.
- ⑦ Preprocessing of Macros.

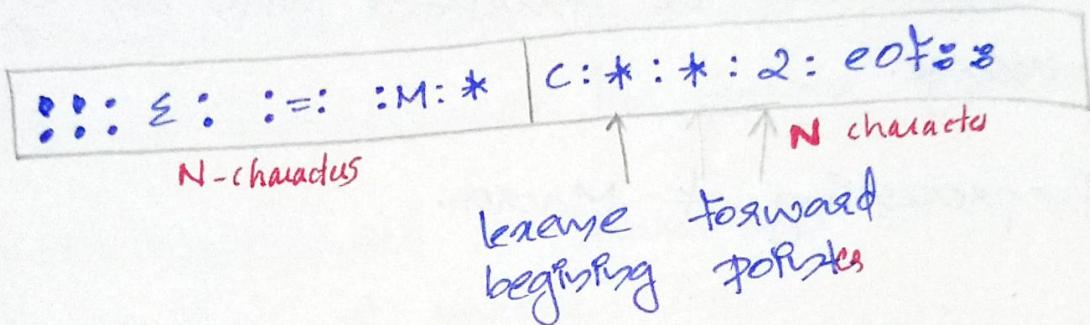
Input Buffering:-

Sometimes lexical analyzer need to look ahead several characters beyond the lexeme for a pattern before a match can be announced.

* As large amount of time can be consumed moving characters we need two buffers scheme to handle large look ahead

safety.

- * Buffer is divided into two n -characters halves. Where n is the number of characters on one disk block.
- * Read n input characters into each half of the buffer with one system read command.



- * If fewer than n characters remain in the input than a special character 'eot' makes the end of source file.
- * Two pointers to the input buffer are maintained.
- * Storing of characters between the two pointers is the current leave.

- * Initially both pointers point to the first character of the next lexeme to be found.
- * One called the forward pointer scans ahead until a match for a pattern is found.
- * Once the next lexeme is determined, the forward pointer is set to the characters at its right end.
- * If the forward pointer is about to move past the halfway mark, right half is filled with a new input character.
- * If the forward pointer is about to move past the right end of the buffer, the left half is filled with a new character and the forward pointer wraps around to the beginning of the buffer.
- * After the lexeme is processed both pointers are set to the characters

immediately past the lexeme.

* code to advance forward pointers

* code to advance forward pointers

if forward at end of first half
then begin reload second half:

forward = forward + 1

end

* else if forward at the end of second
half then begin reload first half

Move forward to beginning of first
half

end

* else

forward = forward + 1

lexer

Eg: abc = pqr * xyz

↓ forward point

leave beg

↓ forward

:a:b:c:=:p: :q:r:x:*:y:z:eof

leave beg

↓ forward

leave beg

abc = Identifier. (first lexeme)

= assign operator (second lexeme)

else [↓ forward] :a:b:c:=:p: :q:r:x:y:z:eof

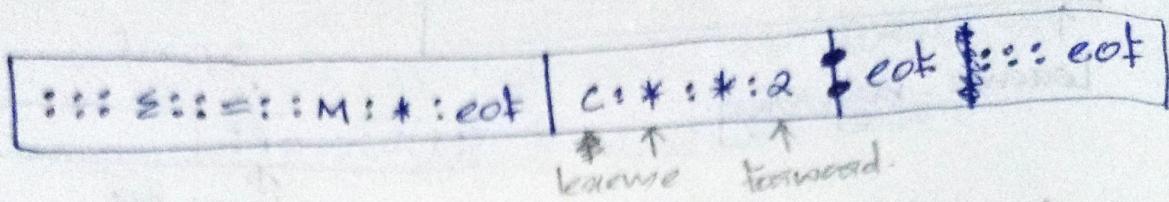
Sentinels:

Except at the end of the buffer halves,
the above code requires two tests for
each advance of the forward pointer.

- * In first model, we want to test after every incrementation of fp
- * We can reduce the two tests to one if we extend each buffer half to hold a

sentinel characters at the end.

- * Sentinel is a special character that cannot be part of the source programs (eof)



- * Look ahead code with sentinels:

forward = forward + 1;

① If forward = eof

then begin
 ~~begin~~

If forward at end of ~~both first half~~

then

begin second second half

forward = forward + 1

end

else if forward at end of second half

then

begin

second part back:

wave toward to beginning of first part

wave toward to beginning of str

Specifications of tokens:-

These are three specifications of tokens:-

- ① strings.
- ② languages.
- ③ Regular expressions.

① Strings:-

An alphabet or character class is a finite set of symbols. A string over an alphabet is a finite sequence of symbols drawn from that alphabet.

a, b → alphabet

aa
ab
ba
bb } → strings

② Languages:-

A Language is any countable set of strings over some ~~finit~~ alphabet. fixed alphabet.

* In language theory, the terms sentence and word are often used as synonyms for string. The length of a string s , usually written $|s|$, is the number of occurrences of symbols in s .

Eg:- banana is a string of length 6.

* The empty string denoted as ϵ , is the string of length zero.

* Symbols = $\{0, 1\}$

String = $0, 1, 01, 10, 11, 00, 111, 100, 101$

Language = $\{1, 11, 111, 1111\}$

* Operations on strings:

The following string related terms are commonly used:

① A prefix of string s is any string obtained

by removing zero or more symbols from the end of string s .

Eg: ban is a prefix of banana.

② A suffix of string s is any string obtained by removing zero or more symbols from the beginning of s .

Eg: nana is a suffix of banana.

③ A substring of s is obtained by deleting any prefix and any suffix from s .

Eg: an is a ~~not~~ substring of banana.

④ The proper prefixes, suffixes and substrings of a string s are those prefixes, suffixes and substrings of s that are not ϵ or not equal to s itself.
(empty)

⑤ A subsequence of s is any string formed by deleting zero or more not necessarily

consecutive positions of 's'

Eg: baaz is a subsequence of banana.

* Operations on languages

The following are the operations that can be applied to languages.

- ① Union.
- ② Concatenation.
- ③ Kleene closure.
- ④ Positive closure

* the following examples shows the operations on strings

Let, $L = \{0, 1\}$ and $S = \{a, b, c\}$

① Union :

$$L \cup S = \{0, 1, a, b, c\}$$

② Concatenation :

$$L \cdot S = \{0a, 1a, 0b, 1b, 0c, 1c\}$$

③ Kleene closure: (0 or More Occurrences)

$$L^* = \{ \epsilon, 0, 1, 00, 11, 10, 01, \dots \}$$

④ Positive closure:

$$L^+ = \{ 0, 1, 00, 11, 10, \dots \} \text{ rather}$$

It is similar to Kleene closure except ϵ .

Regular expressions:

Each regular expression α denotes a language $L(\alpha)$. Here the rules that define the of regular expressions over some alphabet Σ and the languages that those expressions denotes :-

① ϵ is a regular expression and $L(\epsilon)$ is $\{\epsilon\}$, that is the language whose set number is the empty string.

② If 'a' is a symbol in Σ then 'a' is a regular expression and

$L(a) = \{a\}$. That is the language with one string of length one with 'a' in its one position.

③ Suppose x and s are regular expressions denoting the language " $L(x)$ and $L(s)$ ". Then,

(i) $(x) \# (s)$ is a regular expression denoting the language $\underline{L(x) \cup L(s)}$.

(ii) $(x) \cdot (s)$ is a regular expression denoting the language $L(x) \cdot L(s)$.

(iii) $(x)^*$ is a regular expression denoting $(L(x))^*$.

(iv) (x) is a regular expression denoting $L(x)$.

④ May operator * has highest precedence and is left associative concatenation has second precedence.

④ Regular Definitions

Giving names to regular expressions is referred to as a regular definition.

- * If Σ is an alphabet of basic symbols, then a regular definition is a sequence of definitions of the form:

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

$$d_3 \rightarrow r_3$$

- * d_1, d_2, d_3 are the names given to the regular expression (definition).

① Each d_i is a distinct name.

② Each r_i is a regular expression over the alphabet $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$

examples

Eg: identifiers (used in pascal)

regular definition:

letters $\rightarrow A|B| \dots |z|$ $|a|b| \dots |z|$

in this case
 r_2 can use d_1
 r_3 " d_2
 r_4 " d_3

digit \rightarrow 0|1|2 ... 19

$d_1 \rightarrow \text{letter}(r_1)$
 $d_2 \rightarrow \text{digit}(r_2)$
 $d_3 \rightarrow \text{id}(r_3)$

id \rightarrow letter (letter | digit)*

$\text{id} \rightarrow r_1(r_1 | r_2)^*$

Eg: Unsigned number (5280) 39.37, 6.335E4

digit \rightarrow 0|1|2 ... 19

$d_1 \rightarrow r_1$

digits \rightarrow (digit)* digit digit*

$d_2 \rightarrow r_2$

Optional fraction \rightarrow digits / E

$d_3 \rightarrow r_3$

Optional exponent \rightarrow (E (+|-|-) digit) / E

$d_4 \rightarrow r_4$

number \rightarrow digits optional fraction

$d_5 \rightarrow r_5$

optional exponent.

Regular sets

A language that can be defined by a regular expression is called a regular set.

If two regular expressions x and s

denote the same regular set, we say they are equivalent and write $x = s$.

- * There are a number of algebraic laws or regular expressions that can be used to manipulate R to equivalent forms.
- * For instance, $s|t = t|s$ is commutative.
 $r((s|t)) = (r(s)|t)$
 $r((s|t))$
 $(r(s)|t)$

short hand

certain constants occur so frequently in regular expressions that it is convenient to introduce notational short hand for them.

① One or more instances (+):

The unary postfix operator + means *

one or more instances.

- * If R is a regular expression that denotes the language $L(x)$, then ~~$L(x)$~~ $(x)^+$ is a regular expression that denotes the language $(L(x))^+$.

- * Thus the regular expression a^* denotes the set of all strings of one or more a 's.
- * The operator $+$ has the same precedence and associativity as the operator $*$.

② zero or one instance (?)

The unary position operator $?$ means zero or one instance.

- * The notation $r?$ is a shorthand for $r \cup \epsilon$.
- * If r is a regular expression, then $L(r?)$ is a regular expression that denotes the language.

Character classes:-

The notation $[abc]$ where a, b , and c are alphabet symbols, denotes the regular expression $ab|bc$.

- * The character class such $[a-z]$ denotes the regular expression $a|b|c|d|\dots|x$.
- * we can describe patters as being strings generated by the regular expression,
 $[A-Z \ A-Z] \ [A-Za-zA-Z0-9]^*$

Eg @Identifiers

letter $\rightarrow [A-Za-zA-Z]$

digit $\rightarrow [0-9]$

id \rightarrow letter - (letter - | digit)*

② unsigned numbers.

digit $\rightarrow [0-9]$

digit \rightarrow digit+

number \rightarrow digits (. digits)? (E [+ -] ? digits)?

Recognition of tokens using finite automata:

Indicate an accepting state by a double circles, and if there is an action to be taken typically returning a token end an

attribute value to the parser.

- * It is necessary to retract the forward pointer one position (that is, the lexeme does not including include the symbol that got us to the accepting state). Then we shall additionally place a mark at the accepting state.
- * One state is designated the start state, or initial state, it is indicated by an edge, labeled 'start' entering from nowhere.
- * The transition diagram always begins in the start state before any input symbols have been read.
- * Tokens can be recognized with the help of transition diagrams.

Eg: Recognition of tokens (identifiers).

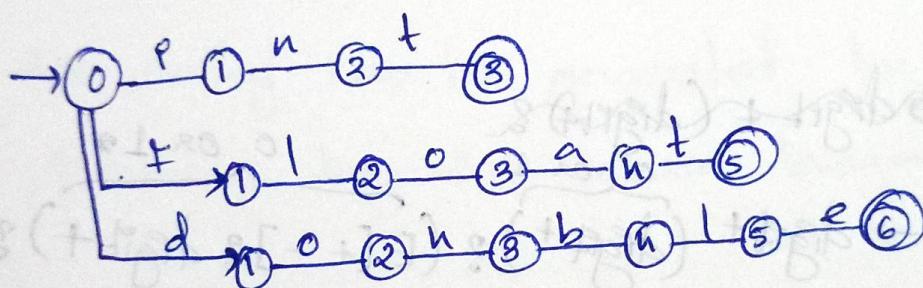
letter \rightarrow $a|b| \dots |z|$
 digit \rightarrow $0|1|2| \dots |9|$
 Id \rightarrow letter (letter/digit)*

⇒ Transition diagram,



② Recognition of keywords:

We can use transition diagrams to recognizing keywords like int, float, double, for, if etc.



③ Recognition of variables:

We can use - at beginning , a digit at the later part. So,

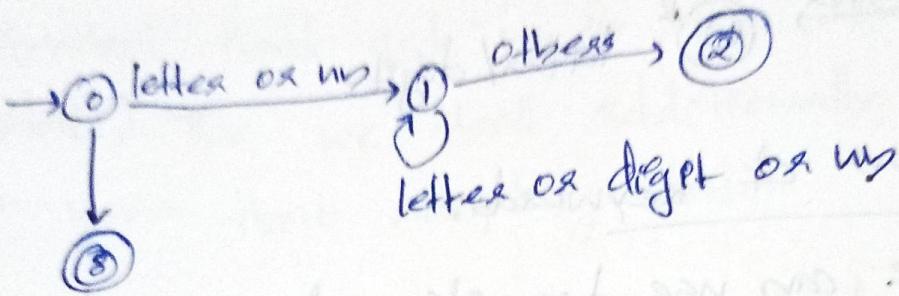
$$w \rightarrow -$$

letter $\rightarrow a | b | c | \dots | z | A | B | \dots | Z |$

digit $\rightarrow 0 | 1 | 2 | \dots | 9 |$

word $\rightarrow (\text{letter} | \text{num})^*$ ($\text{letter} | \text{digit} | \text{num}$) *

for var,



④ Recognition of constants:

Eg: 123, 123.45

number \rightarrow digit + (digit)* &

number \rightarrow digit + $(\underbrace{\text{digit}^+}_0)_? (\underbrace{E [+-]}_0 \underbrace{\text{digit}^+}_1)_?$

digit [$(\cdot \text{digit})^0{}_1$] ? $(E [+-] {}^0{}_1 \text{digit})^0{}_1$

