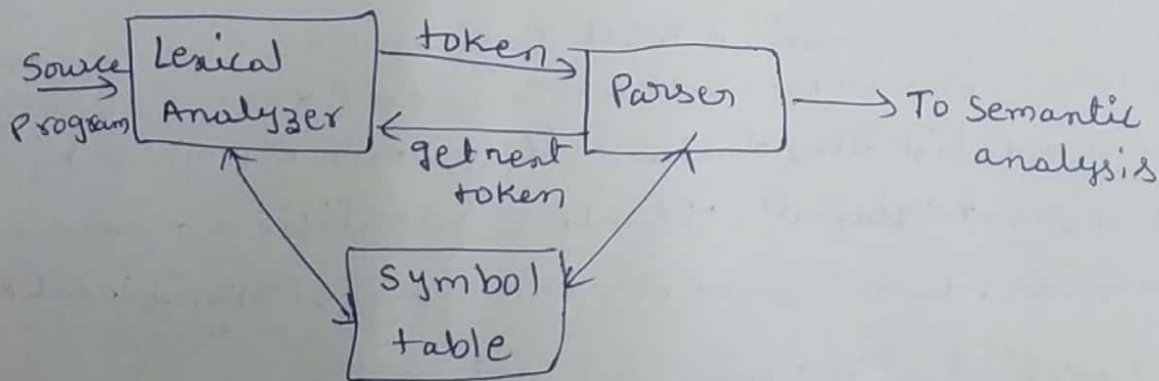


Lexical Analysis:- ^{1st Phase} LA (Scanner) reads the source program one character at a time, turning the source program into a sequence of automatic units called tokens. (2) It reads the source program character by character and returns the tokens of the SP.

Role of the LA:-

1. As the 1st Phase of a compiler, the main task of the lexical analyzer is to read the i/p char's of the SP, group them into lexemes & produce as o/p a sequence of tokens for each lexeme in the SP.
2. The stream of tokens is sent to the Parser for syntactic analysis.
3. When the lexical analyzer discovers a lexeme constituting an identifier, it needs to enter that lexeme into the symbol table. In some cases, info regarding the kind of identifier may be read from the symbol table by the LA to assist it in determining the proper token it must pass to the Parser.
5. These interactions are suggested in fig (below).



Interaction b/w the Lexical analyzer & the Parser.

6. Commonly, the interaction is implemented by making Parser call the L.A.

7. The call, suggested by the getNextToken command, causes the lexical analyzer to read char's from its i/p until it can identify the next lexeme and produce for the next token.

8. Since the LA is part of the compiler that reads the source text, it may perform ~~contain~~ other tasks besides identification of lexemes.

(i) stripping out comments & whitespace (blank, \n, tab, ... etc).

(ii) correlating error messages generated by the compiler with the SP.

(iii) the LA may keep track of the no's of new line chars seen, so it can associate with a line no with error message.

(b) In some compilers, the lexical analyzer makes a copy of the SP with the error messages inserted at the appropriate positions

(c) some times LA are divided into a cascade of 2 processors.

(i) scanning consists of the simple processes that do not require tokenization of the i/p, such as deletion of comments & ^{white spaces} ~~compaction~~

(ii) lexical analysis proper is the more complex portion, which produces tokens from the o/p of the scanner. ^{why to separate into two phases}

Token, patterns & lexemes:- We use 3 related but distinct terms in LA.

Lexeme:- A lexeme is a seq of characters in the SP that matches Pattern for a token & is identified by lexical analyzer as an instance of that token. (Lexeme is matched against Pattern to generate Token)

Pattern:- The rule associated with each set of strings is called P

Token:- A Pattern explains what can be a token these P are defined by means of REG.

Token:- Token is a valid seq of characters which are given by lexeme in SP. tokens are keywords, operators, identifiers, numbers, punctuation symbols identified as possible tokens (Token is invalid it generates error).

Ex:- `int sum(int a, int b);`

<u>lexeme</u>	<u>token</u>	<u>Pattern</u>
int	keyword	datatype
sum	identifier	$L = L(L D)^*$
(open PS	open brace
int	keyword	datatype
a	identifiers	$L = L(L D)^*$
,	PS	comma
int		
b		
;		

int	a=10;	
int	keyword	DT
a	identified	$L = L(L D)^*$
=	operator	Assignment
10	Number/digit	Num/dig
,	delimiter	delimiter
	SP	