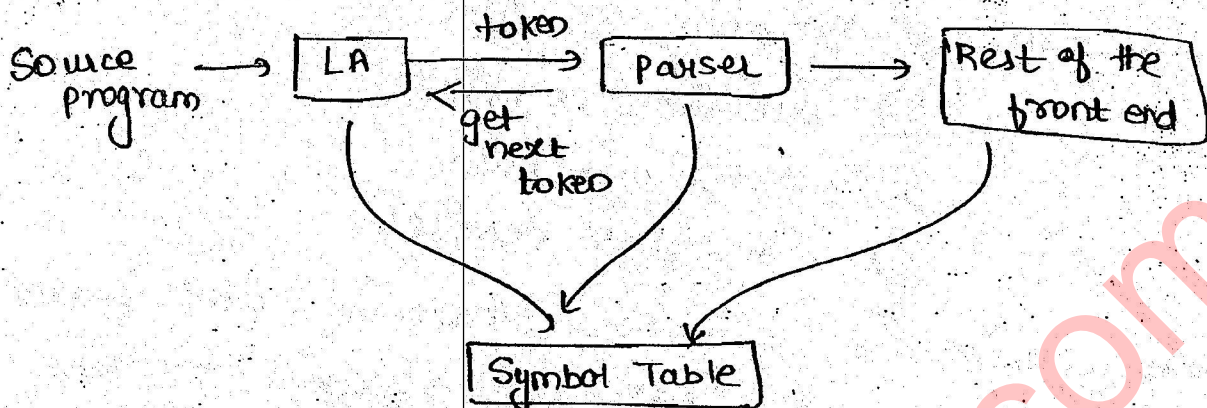


Lexical Analysis



The lexical analyzer scans the source statement character by character, then breaks the statements into stream of tokens. Therefore, lexical analyzer is called as token recognizer.

Other functions -

- ①. maintains separate track of line no.
- ②. skips out white spaces.
- ③. skips out the comment lines.
- ④. To implement a lexical analyzer module we use two concepts -
 - a) finite automata
 - b) Regular expressions.

i.e. the process of lexical analyzer generator starts with token encryption

Each token represented in finite automata from that finite automata lexical analyzer module will be return.

Regular expressions are used to describe tokens of the programming language whereas FA is used to recognize these tokens, i.e.

* Finite automata will be regarded as token recognizer or token scanner.

* Regular expression as token generator or token descriptor.

token (description) \rightarrow RE \rightarrow NFA-E \rightarrow NFA \rightarrow DFA \rightarrow Mini DFA

Code \leftarrow

eg. RE for Integer - \circ .

eg. 123

one or more digit
let $d \rightarrow$ digit.

So, expression = d^+
 $\Rightarrow dd^*$

where $d = 0|1|2|3 \dots |9$

eg. for float - \circ

eg. 12.12

R.E = $d^+ \cdot d^+$
 $= dd^+ \cdot dd^+$

Give a Regular expression for language identifier which defined as letter followed by any no. of letters followed by any no. of digits.

$\Rightarrow ll^+d^+$ where l denotes the letter, d denotes the digit
 $l = a|b|c \dots |z$
 $d = 0|1|2|3 \dots |9$

Q. Give the C language identifier RE.
 let l denotes letters, d denotes digit,
 u denotes underscore.

$$\cancel{(l+u)(l+u)^*d^*(l+u)^*}$$

$$(l+u)(l+u+d)^*$$

← if after digit, letter or - is permitted.

$$(l+u)(l+u+d+\epsilon)^{30}$$

↓ But this will give ∞ which is not possible in C, max limit is 31.

Q. In a programming language P, an identifier defined as letter followed by any no. of letters and digits combinations.

$$\Rightarrow l(l+d)^*$$

Q. In some programming language, identifier defined over the alphabet digit, letters, sp (space symbol)

beginning with letter ending with digit.

let l denote the letter, d denote the digit and S denote the special symbol.

$$\Rightarrow l(l+d+s)^*d$$

RE \rightarrow C - Keyword

C \rightarrow int | float | char | main | . . . | extern

Q. for operators ?

Relational = < | <= | = | != | > | >=

for special symbol -

SP : ; | . | # | → | ↑ | (|) . . .

Tokens Types - :

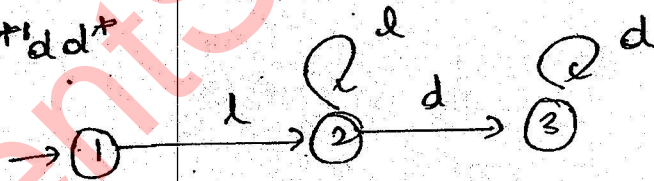
- (a) Identifier
- (b) operator
- (c) Constant
- (d) keywords
- (e) special symbol / punctuation symbol.

Q. In some programming language an identifier define as all names formed using the alphabet & letters, digits, beginning with digits ending with 2 letters - followed by dot of length atmost 32 characters.

let l denotes letters, d denote digits, s denote dot; P denotes special symbol.

$$d(d+l+\epsilon)^{28}lls$$

Let RE = ll^+dd^+
 FA =



state 1 - :

```

c = getcharr()
if (c == l) goto state 2
else
fail()
    
```

state 2 - :

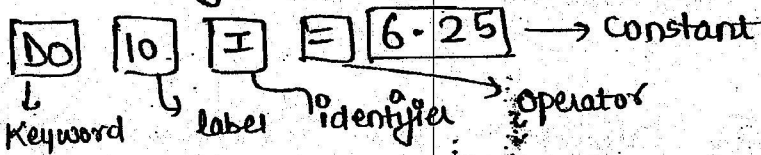
```

c = getcharr()
if (c == l) go to state 2
elseif (c == d) go to state 3
else
return();
    
```


But this is very complex process, instead of this nowadays lex, Flex is used.

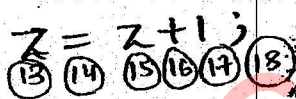
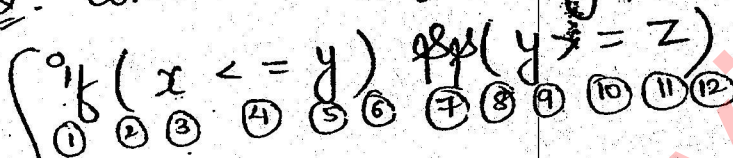
These are open source softwares.

Q. Consider the following FORTRAN statement identify the no. of tokens in the statement.

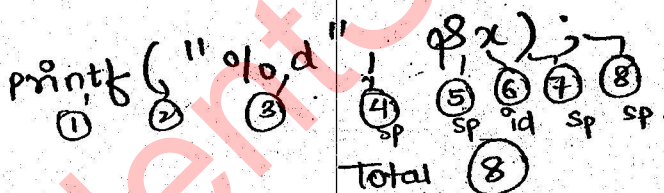


So total = 5.

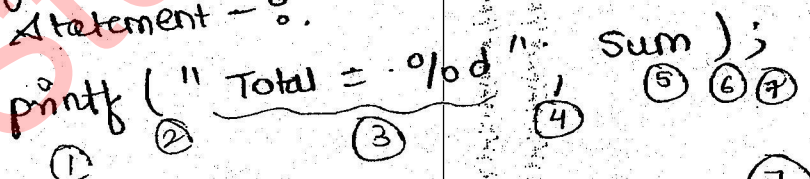
Q. Consider the following C statement



Total = 18.



Q. Find the no. of tokens in the following statement -



Total = 7

Q Which one of the following is not a token of C program?

- (a) 1.02e + 2
- (b) ~~max~~ MAX
- (c) 123.33
- (d) None of these.

Q Which one of the following is not a valid token in C?

- (a) ab123
- (b) 123
- (c) #include → as it is a group of tokens, not a single token.
- (d) 12.2

* CONTEXT FREE GRAMMARS

Context free grammars are useful to define syntaxes of high level languages.

Hence, all high level languages are called CFL's.

eg. C, C++, Java.

A context free grammar is a 4 tuple system (V, T, P, S) and every production is of the form

$A \rightarrow \alpha$

where, $A \in V$

$\alpha \in (V \cup T)^*$

i.e., left side symbol must be non-terminal, right side one can be combination of variables and terminals.

HLL \rightarrow $a = b + c + d$ ^{represented through} CFG.

eg. Integer : 123
 $\langle \text{Integer} \rangle \rightarrow \langle \text{Integer} \rangle \langle \text{Digit} \rangle \mid \langle \text{Digit} \rangle$
 $\langle \text{Digit} \rangle \rightarrow 0 \mid 1 \mid 2 \mid 3 \dots \mid 9$
 Here assume that start symbol is $\langle \text{Integer} \rangle$.

eg. floating : 123.123.
 $\langle \text{floatno} \rangle \rightarrow \langle \text{Int} \rangle \langle \text{dot} \rangle \langle \text{Int} \rangle$
 $\langle \text{Int} \rangle \rightarrow \langle \text{Int} \rangle \langle \text{digit} \rangle \mid \langle \text{digit} \rangle$
 $\langle \text{digit} \rangle \rightarrow 0 \mid 1 \mid 2 \mid 3 \dots \mid 9$
 $\langle \text{dot} \rangle \rightarrow \cdot$

eg. $a * b + c - (d)$
 $\rightarrow \langle \text{EXP} \rangle \rightarrow \langle \text{EXP} \rangle * \langle \text{EXP} \rangle \mid \langle \text{EXP} \rangle + \langle \text{EXP} \rangle \mid \langle \text{EXP} \rangle - \langle \text{EXP} \rangle \mid (\langle \text{EXP} \rangle) \mid a \mid b \mid c \mid d$

eg. In TOC, we write as -

$E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid a$

Here
 $\text{EXP} \rightarrow \text{EXP} + \text{Term} \mid \text{Term}$
 $\text{Term} \rightarrow \text{Term} * \text{Factor} \mid \text{Factor}$
 $\text{Factor} \rightarrow (\text{EXP}) \mid a$

Let

$$w = \cancel{E * E} : a * b + c$$

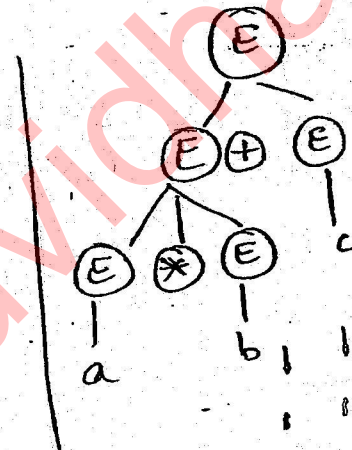
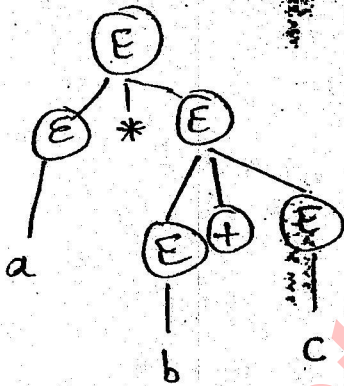
LMD

$E \rightarrow E * E$
 $\rightarrow a * E$
 $\rightarrow a * E + E$
 $\rightarrow a * b + E$
 $\rightarrow a * b + c$

RMD

$E \rightarrow E * E$
 $E \rightarrow E * E + E$
 $E \rightarrow E + E + c$
 $E \rightarrow E * b + c$
 $E \rightarrow a * b + c$

Parse tree -



Since 2 parse tree possible \therefore ambiguous