

USN 

--	--	--	--	--	--	--	--	--	--



**Solution to Internal Assessment Test I – Mar. 2019**

Sub:	System Software & Compiler Design				Sub Code:	15CS63	Branch:	CSE
Date:	6/03/2019	Duration:	90 min's	Max Marks:	50	Sem/Sec:	6/CSE(A,B,C)	OBE

**1. a) What is Compiler? Explain the various phases of a compiler with a neat diagram. Show the translations for an assignment statement  $B=5+C*D-6$ , clearly indicate the output of each phase.**

**Solution:**

A **compiler** is a program that converts high-level language to assembly language.

**Phases of Compiler**

The compilation process is a sequence of various phases. Each phase takes input from its previous stage, has its own representation of source program, and feeds its output to the next phase of the compiler.

**Lexical Analysis**

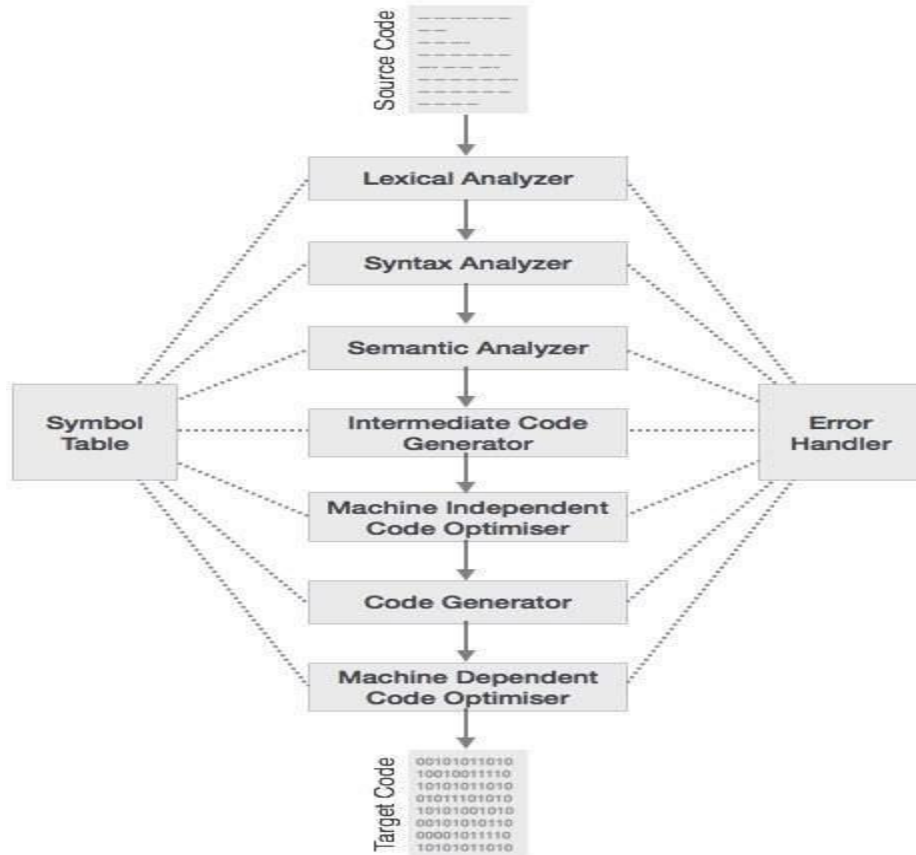
The first phase of scanner works as a text scanner. This phase scans the source code as a stream of characters and converts it into meaningful lexemes. Lexical analyzer represents these lexemes in the form of tokens as:

`<token-name, attribute-value>`

$B=5+C*D-6$

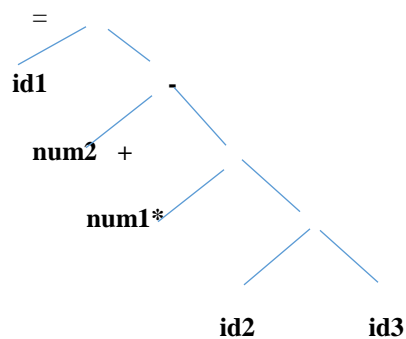
`<id1><= ><num1><+><id2><*><id3><-><num2>`

Divide into 9 tokens -3 ids ,2 nums and 4 operator



### Syntax Analysis

The next phase is called the syntax analysis or **parsing**. It takes the token produced by lexical analysis as input and generates a parse tree (or syntax tree). In this phase, token arrangements are checked against the source code grammar, i.e. the parser checks if the expression made by the tokens is syntactically correct.

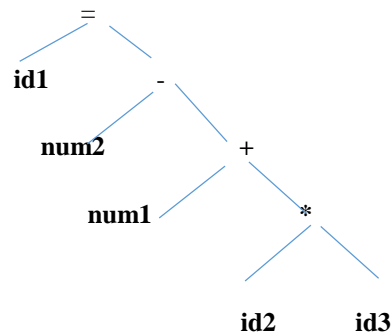


### Semantic Analysis

Semantic analysis checks whether the parse tree constructed follows the rules of language.

For example,

Assignment of values is between compatible data types, and adding string to an integer. Also, the semantic analyzer keeps track of identifiers, their types and expressions; whether identifiers are declared before use or not etc. The semantic analyzer produces an annotated syntax tree as an output.



### Intermediate Code Generation

After semantic analysis the compiler generates an intermediate code of the source code for the target machine. It represents a program for some abstract machine. It is in between the high-level language and the machine language. This intermediate code should be generated in such a way that it makes it easier to be translated into the target machine code.

T1=id2\*id3

T2=num1+T1

T3=T2-num2

Id1=T3

### Code Optimization

The next phase does code optimization of the intermediate code. Optimization can be assumed as something that removes unnecessary code lines, and arranges the sequence of statements in order to speed up the program execution without wasting resources (CPU, memory).

T1=id2\*id3

T2=num1+T1

Id1=T2-num2

## Code Generation

In this phase, the code generator takes the optimized representation of the intermediate code and maps it to the target machine language. The code generator translates the intermediate code into a sequence of (generally) re-locatable machine code. Sequence of instructions of machine code performs the task as the intermediate code would do.

LDA r1,id2

LDA r2,id3

MUL r1,r1,r2

LDA r3,num1

Add r1,r1,r3

LDA r4,num2

SUB r1,r1,r4

STR id1,r1

## Symbol Table

It is a data-structure maintained throughout all the phases of a compiler.

1. All the identifier's names along with their types are stored here.
2. The symbol table makes it easier for the compiler to quickly search the identifier record and retrieve it. The symbol table is also used for scope management.

1	Id1	-----
2	Id2	----
3	Id3	-----

1	Num1	----
2	Num2	-----

2.a) Write the algorithm used for eliminating the left recursion. Eliminate left recursion from the given grammar.  $S \rightarrow A$ ,  $A \rightarrow Ad / Ae / aB / ac$ ,  $B \rightarrow bBc / f$

**Solution:-**

**Left recursion elimination algorithm:**

Arrange the nonterminals in some order  $A_1, A_2, \dots, A_n$ .  
 For (each  $i$  from 1 to  $n$ )  
 {  
     For (each  $j$  from 1 to  $i-1$ )  
     {  
       Replace each production of the form  $A_i \rightarrow A_j \gamma$  by the production  $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$  where  $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$  are all current  $A_j$  productions  
     }  
   Eliminate left recursion among the  $A_i$ -productions  
 }

$S \rightarrow A$ ,  $A \rightarrow Ad / Ae / aB / ac$ ,  $B \rightarrow bBc / f$

**Solution:**

$S \rightarrow A$   
 $A \rightarrow aBA' / acA'$   
 $A' \rightarrow dA' / eA' / \epsilon$   
 $B \rightarrow bBc / f$

2.b) Write the regular definition for unsigned number. Also draw the transition diagram.

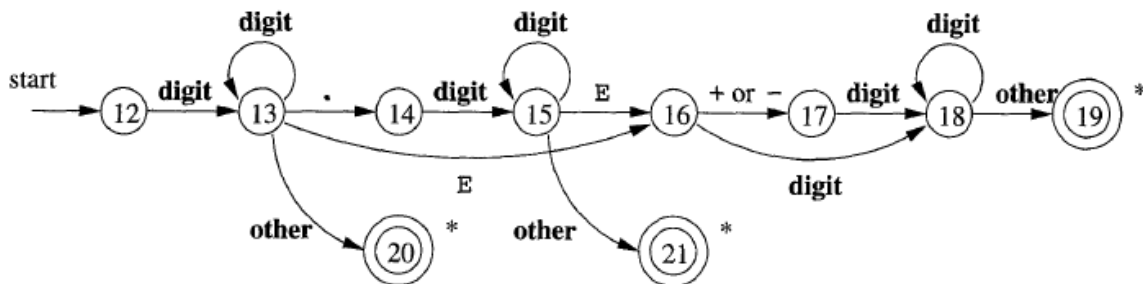
**Solution:**

The regular definition for unsigned number

$letter\_ \rightarrow [A-Za-z\_]$   
 $digit \rightarrow [0-9]$   
 $id \rightarrow letter\_ (letter \mid digit)^*$

$digit \rightarrow [0-9]$   
 $digits \rightarrow digit^+$   
 $number \rightarrow digits (. digits)? (E [+ -]? digits)?$

The transition diagram for unsigned number



3 a) What are the key problems with top down parsing? Write a recursive descent parser for the grammar:  $cAd$ ,  $A \rightarrow ab \mid a$  for the input string  $w=cad$ .

**Solution:**

**Problem with top down parsing:**

1. The key problem is that of determining the production to be applied for a nonterminal, say A. Once an A-production is chosen, the rest of the parsing process consists of "matching" the terminal symbols in the production body with the input string.
2. This may require backtracking to find the correct A-production to be applied.
3. Elimination of left recursion of grammar is required
4. Left factoring is required

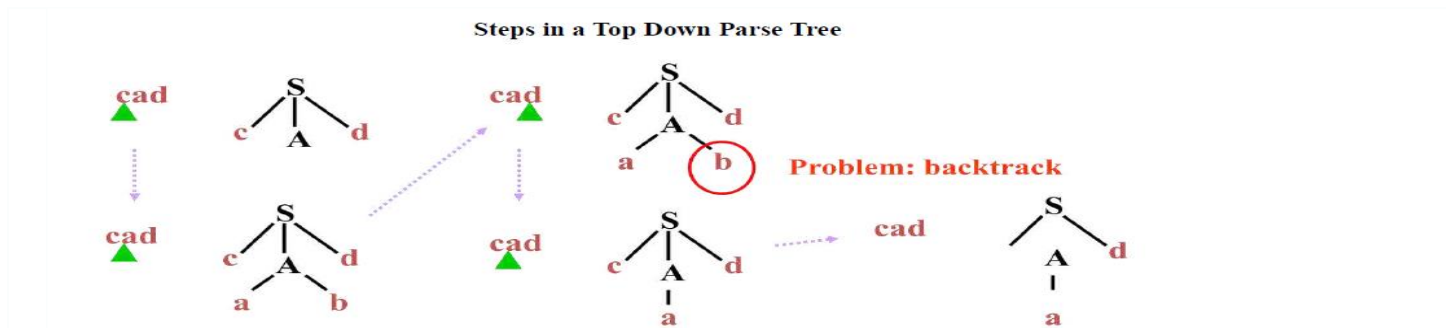
**A recursive descent parser for the grammar:**  $S \rightarrow cAd$ ,  $A \rightarrow ab \mid a$  for the input string  $w=cad$ .

- General category of Parsing Top-Down
- A recursive-descent parsing program consists of a set of procedures, one for each nonterminal.
- Execution begins with the procedure for the start symbol, which halts and announces success if its procedure body scans the entire input string.
- Choose production rule based on input symbol
- May require backtracking to correct a wrong choice.

Example:  $S \rightarrow cAd$

$A \rightarrow ab \mid a$

**input: cad**



3. (b) What is Token, Lexeme and pattern? Explain with an example.

**Solution:**

A **token** is a pair consisting of a token name and an optional attribute value. The token name is an abstract symbol representing a kind of lexical unit, e.g., a particular keyword, or a sequence of input characters denoting an identifier. The token names are the input symbols that the parser processes. We will often refer to a token by its token name.

A **pattern** is a description of the form that the lexemes of a token may take. In the case of a keyword as a token, the pattern is just the sequence of characters that form the keyword. For identifiers and some other tokens, the pattern is a more complex structure that is *matched* by many strings.

A **lexeme** is a sequence of characters in the source program that matches the pattern for a token and is identified by the lexical analyzer as an instance of that token.

Token	INFORMAL DESCRIPTION	SAMPLE LEXEMES
if	characters i, f	if
else	characters e, l, s, e	else
Comparison	comparison < or > or <= or >= or == or !=	<=
Id	letter followed by letters and digits	Total

Numeral	any numeric constant	3.4
Literal	anything but ", surrounded by "	"hello world"

In many programming languages, the following classes cover most or all of the tokens:

1. One token for each keyword. The pattern for a keyword is the same as the keyword itself.
2. Tokens for the operators, either individually or in classes such as the token
3. One token representing all identifiers.
4. One or more tokens representing constants, such as numbers and literal strings.
5. Tokens for each punctuation symbol, such as left and right parentheses, comma, and semicolon

Example: Consider the following C statement

```
printf ("Total = %d\n", score) ;
```

both `printf` and `score` are lexemes matching the pattern for token **id**, and

`"Total = %d\n"` is a lexeme matching literal.

**4.a) Given the grammar:  $E \rightarrow E - T \mid T$ ,  $T \rightarrow T / F \mid F$ ,  $F \rightarrow (E) \mid \text{num}$**

Construct the predictive parsing table and show the moves made by predictive parser on the input string **w=6/3-1**

**Solution:**

**Apply left recursion**

$E \rightarrow TE'$   
 $E' \rightarrow -TE' \mid \epsilon$   
 $T \rightarrow FT'$   
 $T' \rightarrow /FT' \mid \epsilon$   
 $F \rightarrow (E) \mid \text{num}$

**First and follow set**

$\text{First}(E) = \{ (, \text{num} \}$        $\text{follow}(E) = \{ \$, ) \}$   
 $\text{First}(E') = \{ -, \epsilon \}$        $\text{follow}(E') = \{ -, \$, ) \}$   
 $\text{First}(T) = \{ (, \text{num} \}$        $\text{follow}(T) = \{ -, \$, ) \}$   
 $\text{First}(T') = \{ /, \epsilon \}$        $\text{follow}(T') = \{ -, \$, ) \}$   
 $\text{First}(F) = \{ (, \text{num} \}$        $\text{follow}(F) = \{ -, /, \$, ) \}$



### Parsing table

NT	ACTION					
	(	)	/	-	num	\$
<b>E</b>	E -> TE'				E -> TE'	
<b>E'</b>		E' -> ε		E' -> -TE'		E' -> ε
<b>T</b>	T -> FT'				T -> FT'	
<b>T'</b>		T' -> ε	T' -> /FT'	T' -> ε		T' -> ε
<b>F</b>	F -> (E)				F -> num	

### Tracing the string

Stack	Input	action
<b>\$E</b>	Num/num-num\$	<b>E-&gt;TE'</b>
<b>\$E'T</b>	Num/num-num\$	<b>T -&gt; FT'</b>
<b>\$E'T'F</b>	Num/num-num\$	<b>F -&gt;num</b>
<b>\$E'T'num</b>	Num/num-num\$	<b>match</b>
<b>\$E'T'</b>	/num-num\$	<b>T' -&gt; /FT'</b>
<b>\$E'T'F/</b>	/num-num\$	<b>match</b>
<b>\$E'T'F</b>	num-num\$	<b>F -&gt;num</b>
<b>\$E'T'num</b>	num-num\$	<b>Match</b>
<b>\$E'T'</b>	-num\$	<b>T' -&gt; ε</b>
<b>\$E'</b>	-num\$	<b>E' -&gt; -TE'</b>
<b>\$E'T-</b>	-num\$	<b>match</b>
<b>\$E'T</b>	num\$	<b>T -&gt; FT'</b>
<b>\$E'T'F</b>	num\$	<b>F-&gt;num</b>
<b>\$E'T'num</b>	num\$	<b>match</b>
<b>\$E'T'</b>	\$	<b>T' -&gt; ε</b>
<b>\$E'</b>	\$	<b>E' -&gt; ε</b>
<b>\$</b>	<b>\$</b>	<b>accept</b>

5 (a) Write the algorithm to find FIRST and FOLLOW for the given grammar

#### Solution:

To compute First(X) for all grammar symbols X, apply following rules until no more terminals or ε can be added to any First set:

## FIRST and FOLLOW sets

- If  $X$  is terminal,  $\text{FIRST}(X) = \{X\}$ .
- If  $X \rightarrow \epsilon$  is a production, then add  $\epsilon$  to  $\text{FIRST}(X)$ .
- If  $X$  is a non-terminal, and  $X \rightarrow Y_1 Y_2 \dots Y_k$  is a production, and  $\epsilon$  is in all of  $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_k)$ , then add  $\epsilon$  to  $\text{FIRST}(X)$ .

$\text{FOLLOW}(A)$  for a nonterminal  $A$  is defined as to be the set of terminals  $a$  that can appear immediately to the right of  $A$  in some sentential form

That is, the set of terminals  $a$  such that there exists a derivation of the form  $S \Rightarrow \alpha A a \beta$ , for some  $\alpha$  and  $\beta$

**To compute  $\text{FOLLOW}(A)$  for all nonterminals  $A$ , apply the following rules until nothing can be added to any  $\text{FOLLOW}$  set**

- If  $\$$  is the input end-marker, and  $S$  is the start symbol,  $\$ \in \text{FOLLOW}(S)$ .
- If there is a production,  $A \rightarrow \alpha B \beta$ , then  $(\text{FIRST}(\beta) - \epsilon) \subseteq \text{FOLLOW}(B)$ .
- If there is a production,  $A \rightarrow \alpha B$ , or a production  $A \rightarrow \alpha B \beta$ , where  $\epsilon \in \text{FIRST}(\beta)$ , then  $\text{FOLLOW}(A) \subseteq \text{FOLLOW}(B)$ .

5. (b) Compute FIRST and FOLLOW for the grammar

$S \rightarrow aBDh \quad B \rightarrow cC \quad C \rightarrow bC / \epsilon \quad D \rightarrow EF \quad E \rightarrow g / \epsilon \quad F \rightarrow f / \epsilon$

**Solution:**

$\text{First}(S) = \{a\}$

$\text{First}(B) = \{c\}$

$\text{First}(C) = \{b\}$

$\text{First}(D) = \{g, f, \epsilon\}$

$\text{First}(E) = \{g, \epsilon\}$

$\text{First}(F) = \{f, \epsilon\}$

$\text{Follow}(S) = \{\$\}$

$\text{Follow}(B) = \{g, f, h\}$

$\text{Follow}(C) = \{g, f, h\}$

$\text{Follow}(D) = \{h\}$

$\text{Follow}(E) = \{f, h\}$

$\text{Follow}(F) = \{h\}$

6 (a) Show that the following grammar is ambiguous. Write an equivalent unambiguous grammar for the same.  $E \rightarrow E+E \mid E * E \mid E-E \mid (E) \mid id$

**Solution:**

Two LMD for same string so grammar is ambiguous

Example:

**String  $3*2+5$  can be derived in 2 ways-LMD**

I) First leftmost derivation

$E \Rightarrow E * E$   
 $\Rightarrow id * E$   
 $\Rightarrow 3 * E + E$   
 $\Rightarrow 3 * id + E$   
 $\Rightarrow 3 * 2 + E$   
 $\Rightarrow 3 * 2 + id$   
 $\Rightarrow 3 * 2 + 5$

II) Second leftmost derivation

$E \Rightarrow E + E$   
 $\Rightarrow E * E + E$   
 $\Rightarrow id * E + E$   
 $\Rightarrow 3 * E + E$   
 $\Rightarrow 3 * id + E$   
 $\Rightarrow 3 * 2 + id$   
 $\Rightarrow 3 * 2 + 5$

6. b) Eliminate left factoring from the following grammar:

$S \rightarrow bSSa a S \mid bSSa S b \mid bSSa \mid a$   
 $A \rightarrow aAb \mid bc \mid aAc$

1)  $S \rightarrow bSSa a S \mid bSSa S b \mid bSSa \mid a$

**Solution:**

$S \rightarrow bSSa S' \mid a$   
 $S' \rightarrow aS \mid S b \mid \epsilon$

2)  $A \rightarrow aAb \mid bc \mid aAc$

**Solution:**

$A \rightarrow aAA' \mid bc$   
 $A' \rightarrow b \mid c$

### 7.(a) Explain the input buffering strategy used in lexical analysis phase.

#### Solution:

#### Input Buffering:

- Some efficiency issues concerned with the buffering of input.
- A two-buffer input scheme that is useful when lookahead on the input is necessary to identify tokens.
- Techniques for speeding up the lexical analyser, such as the use of sentinels to mark the buffer end.
- There are three general approaches to the implementation of a lexical analyser:
  1. Use a lexical-analyser generator, such as Lex compiler to produce the lexical analyser from a regular expression based specification. In this, the generator provides routines for reading and buffering the input.
  2. Write the lexical analyser in a conventional systems-programming language, using I/O facilities of that language to read the input.
  3. Write the lexical analyser in assembly language and explicitly manage the reading of input.

#### Buffer pairs:

- Because of a large amount of time can be consumed moving characters, specialized buffering techniques have been developed to reduce the amount of overhead required to process an input character.
- The scheme to be discussed:
- Consists a buffer divided into two N-character halves.

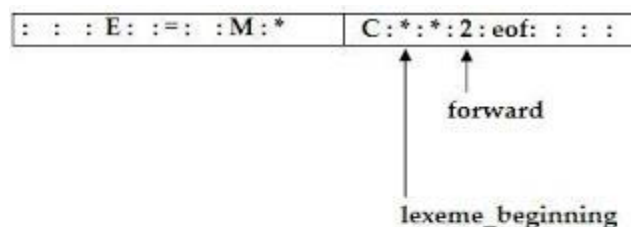


Fig 1.8 An input buffer in two halves

N – Number of characters on one disk block, e.g., 1024 or 4096.

1. Read N characters into each half of the buffer with one system read command.
2. If fewer than N characters remain in the input, then eof is read into the buffer after the input characters.
3. Two pointers to the input buffer are maintained.
4. The string of characters between two pointers is the current lexeme.
5. Initially both pointers point to the first character of the next lexeme to be found.
6. Forward pointer, scans ahead until a match for a pattern is found.
7. Once the next lexeme is determined, the forward pointer is set to the character at its right end.

8. If the forward pointer is about to move past the halfway mark, the right half is filled with N new input characters.
9. If the forward pointer is about to move past the right end of the buffer, the left half is filled with N new characters and the forward pointer wraps around to the beginning of the buffer.
10. *Disadvantage of this scheme:*
11. This scheme works well most of the time, but the amount of lookahead is limited.
12. This limited lookahead may make it impossible to recognize tokens in situations where the distance that the forward pointer must travel is more than the length of the buffer.
13. For example: DECLARE ( ARG1, ARG2, ... , ARGn ) in PL/1 program;
14. Cannot determine whether the DECLARE is a keyword or an array name until the character that follows the right parenthesis.

### Sentinels:

1. In the previous scheme, must check each time the move forward pointer that have not moved off one half of the buffer. If it is done, then must reload the other half.
2. Therefore the ends of the buffer halves require two tests for each advance of the forward pointer.
3. This can reduce the two tests to one if it is extend each buffer half to hold a sentinel character at the end.
4. The sentinel is a special character that cannot be part of the source program. (eof character is used as sentinel).

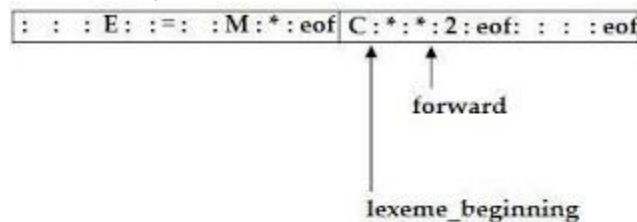


Fig 1.9 Sentinels at the end of each buffer half

1. In this, most of the time it performs only one test to see whether forward points to an eof.
2. Only when it reach the end of the buffer half or eof, it performs more tests.
3. Since N input characters are encountered between eof's, the average number of tests per input character is very close to 1.

**7.(b1) Explain the structure of LEX program. Write a LEX program to count number of words, characters, spaces and lines present in a given input file.**

**Solution:**

**ALEXprogramconsistsofthreeparts:**

```

declarations
%%
translation rules
%%

```

auxiliary procedures

The declarations section includes declarations of variables, constants, and regular definitions.

The translation rules of a lex program are statements of the form

R1 {action1}

R2 {action2}

.....

Rn {actionn} where each R<sub>i</sub> is a regular expression and each action<sub>i</sub> is a program fragment describing what action the lexical analyzer should take when pattern R<sub>i</sub> matches a lexeme. Typically, action<sub>i</sub> will return control to the parser. In Lex actions are written in C; in general, however, they can be in any implementation language.

The third section holds whatever auxiliary procedures are needed by the actions.

### **7.b2) Program to count the number of characters, words, spaces and lines in a given input file.**

```
% {
#include<stdio.h>
int line=0,cha=0,word=0,space=0;
% }
%%
[\n] {line++;}
[ \t" "] {space++;}
[a-zA-z0-9]* {word++;cha+=yyleng;}
%%
int main(int argc, char *argv[])
{
yyin=fopen(argv[1], "r");
yylex();
fclose(yyin);
printf("No. of lines=%d\n, characters=%d\n, words=%d\n, spaces=%d\n", line, cha, word, space);
return 0;
}
```