

## UNIT - IV

### FORMAL LANGUAGES

#### Languages:-

A Language is a dynamic set of visual, Auditory, or tactile symbols of communication and the elements used to manipulate them. Language can also refer to the use of such systems as a general phenomenon.

#### Symbols and Alphabet:-

A symbol is an abstract entity. It cannot be formerly defined as points in geometry.

EX: Letters, digits, (or) special symbols like \$, @, #, etc.

#### Alphabet:-

A finite collection of symbols denoted by  $\Sigma$ .

EX: English Alphabet  $\Sigma = \{a, b, \dots, z\}$

Binary Alphabet  $\Sigma = \{0, 1\}$

#### String/word:-

A set of symbols from alphabet.

EX: 001, 110, 1111 strings from binary alphabet.

a01 is not a string from binary alphabet.

Word:-

A word over an alphabet can be any finite sequence (or) string, group of letters.

Note:-

- \* The set of all words over an alphabet  $\Sigma$  is usually denoted by  $\Sigma^*$ .
- \* For any alphabet there is only one word of length 0 the empty word is denoted by  $\epsilon$ ,  $\epsilon$  (or)  $\lambda$ .
- \* An empty string can be denoted by  $\epsilon$ .

Note:- \* string: "abc"

Prefix: a, ab

Suffix: c, abc

Substring: a, b, c, bc, ab.

\* Language is a set of words (or) sentences.

Operations on Languages:-

If  $L_1$  and  $L_2$  are two languages then,

(i) Union  $\rightarrow L_1 + L_2$  (or)  $L_1 \cup L_2$

(ii) Concatenation  $\rightarrow L_1 L_2$

(iii) Kleene's closure  $\rightarrow \Sigma^*$

EX:1 If  $\Sigma = \{x\}$  then find  $\Sigma^* = \{\epsilon \cup \Sigma^1 \cup \Sigma^2 \cup \Sigma^3 \dots\}$

$$\Sigma^* = \{\epsilon, x, xx, xxx, \dots\}$$

EX:2 If  $\Sigma = \{a, b\}$  then find  $\Sigma \cup \Sigma^*$

$$\Sigma \cup \Sigma^* = \{\epsilon, a, b, aa, bb, \overline{ab}, \overline{ba}, aaaa, aab, \dots\}$$

Note:- Positive closure  $\Sigma^+ = \Sigma^* - \{\epsilon\}$

$$\Sigma^* = \Sigma^+ \cup \epsilon$$

EX:3 If  $L_1 = \{good, bad\}$ ,  $L_2 = \{boy, girl\}$

then find  $L_1 \cup L_2$  and  $L_1 L_2$ ?

$$L_1 \cup L_2 = \{good, bad, boy, girl\}$$

$$L_1 L_2 = \{goodboy, goodgirl, badboy, badgirl\}$$

## Grammars

Grammar is basically defined as set of 4-tuples  $G = (V, T, P, S)$ , where  $V$  is set of non terminals (variables)

$T$  is set of Terminals (primitive symbols)

$P \rightarrow P$  is the set of Productions (rules) which is relate the non terminals and terminals.

$A$  and  $S \rightarrow$  is start symbol with which strings in grammar are derived.

A production rules has the form  $\alpha \rightarrow \beta$  where  $\alpha$  and  $\beta$  are strings on  $V \cup T$  and atleast one symbol of  $\alpha$  belongs to  $V$ .

EX:  $G = (\{S, A, B\}, \{a, b\}, S, \{S \rightarrow AB, A \rightarrow a, B \rightarrow b\})$

$$V = \{S, A, B\}$$

$$T = \{a, b\}$$

$$S = S$$

$$P = S \rightarrow AB, A \rightarrow a, B \rightarrow b$$

Ans:  $S \rightarrow AB$   
 $\rightarrow aB$   
 $\rightarrow aa$ .

Note:- Terminals  $\rightarrow \{0, 1, b, \dots, z\}, 0-9, (, ) \dots$

Non-terminals  $\rightarrow \{A, \dots, Z\}$ .

$\alpha, \beta, \gamma$  - both terminal, and non-terminals.

## Regular Grammar

Noam Chomsky gave a mathematical model of grammar which is effective for writing computer languages.

### Type-0 Grammar:-

(i) Unrestricted Grammar:- (or) Recursive Grammar:-

A phrase type grammar with no restriction on production.

Production of the form  $\boxed{\alpha \rightarrow \beta}$

$\alpha, \beta$  contains any number of terminals and non-terminals.

This grammar can be modeled using "Turing Machine".

EX:  $AaC \rightarrow bBDE$

$aBD \rightarrow abcDE$ .

### Type-2 Context sensitive Grammar.

\* This grammar generates context sensitive languages.

\* The production is in the form of

$\boxed{\alpha A \beta \rightarrow \alpha \gamma \beta}$

where  $A$  - non terminal.

$\alpha, \beta, \gamma$  - contains any number of terminals, ~~and~~ & non-terminal grammar modeled by "Linear bounded Automata"

Q.

Ex:  $S \rightarrow aAB$

$$AB \rightarrow aAB$$

$$B \rightarrow b.$$

Type-2 Context free Grammar

The production is of the form  $A \rightarrow \gamma$

left side  $A$  should be non-terminal.

right side  $\gamma$  contain any no. of terminal and non terminal.

grammar modeled by pushdown Automata [PDA].

Ex:  $S \rightarrow aA$

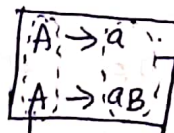
$$S \rightarrow b$$

$$S \rightarrow aA.$$

Q.

Type-3 Regular Grammar:-

Production is in the form of



should be non-terminal

should contain any no. of terminal and non-terminal grammar.

This grammar is modelled using "Finite Automata".

Right Linear Grammar:-

A grammar is said to be right linear if all productions are of the form.

$$A \rightarrow xB$$

$$A \rightarrow x.$$

where  $A, B \in V$  and  $x \in T$ .

Eg:  $S \rightarrow abS / b$  ,  ~~$S \rightarrow abS / b$~~

## Left Linear Grammar:-

A grammar is said to be left linear grammar. if all productions are of the form.

$$A \rightarrow Bx$$

$$A \rightarrow x.$$

where  $A, B \in V$  and  $x \in T$ .

EG:-  $S \rightarrow sbblb$

~~Recap~~

## Pumping Lemma for Regular Language:-

A Language is said to be regular. if it is accepted either by a finite Automaton (or) it has a regular grammar generating it. In order to prove that a Language is not regular the most commonly used technique is "Pumping lemma".

The lemma gives a pumping property that a sufficiently long word has a subword (non-empty) that can be pumped.

## Theorem: (Pumping lemma).

### statement:-

Let  $L$  be a regular language over  $T$ . Then there exists a constant  $k$  depending on  $L$  such that for each  $w \in L$  with  $|w| \geq k$  there exists  $x, y, z \in T^*$  such that  $w = xyz$  and

1.  $|xy| \leq k$

2.  $|y| > 1$

3.  $xy^i z \in L \forall i \geq 0$

where  $k$  is no more than the number of states in the minimum state Automaton accepting  $L$ .

### Proof:-

Let  $M = (K, \Sigma, \delta, q_0, F)$  be a deterministic finite state Automaton (DFSA) accepting  $L$ . Let  $K = \{q_1, \dots, q_k\}$ .

Let  $w = a_1 \dots a_m \in L$ , where  $a_i \in \Sigma, 1 \leq i \leq m, m \geq k$ .

Let the transitions on  $w$  be as shown below

$$q_1 a_1 \dots a_m \vdash a_1 q_2 a_2 \dots a_m \vdash \dots \vdash a_1 \dots a_m a_{m+1}$$

where  $q_j \in K, 1 \leq j \leq m+1$ . Here  $a_1 \dots a_{j-1} q_j a_j \dots a_m$  means the FSA is in  $q_j$  state after reading  $a_1 \dots a_{j-1}$  and the input head is pointing to  $a_j$ . Clearly in the above transitions,  $m+1$  states are visited but  $M$  has only  $k$  states. Hence there exists  $q_i, q_j$  such that  $q_i = q_j$ . Hence for

$$q_1 a_1 \dots a_m \vdash a_1 q_2 a_2 \dots a_m \vdash (a_1 \dots a_{j-1} q_j a_j \dots a_m \dots \vdash a_1 \dots a_{j-1} q_j a_j \dots a_m) \vdash \dots \vdash a_1 \dots a_m q_{m+1}$$



where the transitions between the brackets start and end at  $q_i$ . Processing a string  $a^t$  for  $t \geq 0$ . Hence if  ~~$x = a$~~   
 $x = a_1 \dots a_{j-1}$ ,  $y = a_j \dots a_j$ ,  $z = a_{j+1} \dots a_m$ ,  $xy^t z \in L \quad \forall t \geq 0$

where  $|xy| \leq k$  since  $q_i$  is the first state identified to repeat in the transition and  $|y| \geq 1$  hence the lemma.

The Potential application of this lemma is that it can be used to show that some languages are non-regular.

EX:1 show that  $L = \{a^n b^n \mid n \in \mathbb{N}\}$  is not regular.

let  $L = \{a^n b^n \mid n \in \mathbb{N}\}$

If  $L$  is regular, then by Pumping lemma

$\exists$  a constant 'k' which is satisfying the pumping lemma conditions.

Now choose  $w = a^k b^k$

Clearly  $|w| > k$ , then  $w = xyz$ ,  $|xy| \leq k$  and  $|y| \geq 1$ .

If  $|x| = p$ ,  $|y| = q$ ,  $|z| = r$ ,  $p+q+r = 2k$  and  $p+q \leq k$ .

Hence  $xy$  consists of only a's and since  $|y| \geq 1$ ,  $xz \notin L$  as number of a's in  $xz$  is less than  $k$  and  $|z|_b = k$

~~$\rightarrow$~~

Hence the pumping lemma is not true, for  $i=0$  as  $xy^i z$  must be in  $L$  for  $i \geq 0$ .

$\text{H.W.}$  Hence  $L$  is not regular.

EX:2 show that  $L = \{a^i b^j \mid i, j \geq 1, i \neq j\}$  is not regular.

Definition:- A Grammar  $G = (V, T, P, S)$  is said to be Context free if all Production in  $P$  have the form  $A \rightarrow \alpha$ , where  $A \in V$  and  $\alpha \in (V \cup T)^*$ .

$V \rightarrow$  non-terminals.

$T \rightarrow$  Terminals.

$P \rightarrow$  Productions

$S \rightarrow$  The start symbol.

~~Context free language.~~

Context free Language:- (CFL)

A Language generated by a CFG is called a Context Free Language. (CFL).

EX:1 terminal:  $a$

nonterminal:  $S$

Productions:  $S \rightarrow aS$

$S \rightarrow \epsilon$

Is a simple CFG that defines  $L(G) = a^*$

where  $V = \{S\}$   $T = \{a\}$

EX:2 The CFG for defining Palindrome over  $\{a, b\}$ .

The Productions  $P$  are  $S \rightarrow \epsilon | a | b$

$S \rightarrow aSa$

$S \rightarrow bSb$ .

and the Grammar is  $G = (\{S\}, \{a, b\}, P, S)$

EX: 3 Derive 'a<sup>n</sup>' from a Grammar.

Terminal : a

Nonterminal : S

Productions :  $S \rightarrow aS$   
 $S \rightarrow \epsilon$

Solution:-

$$S \Rightarrow aS$$

$$\Rightarrow aaS$$

$$\Rightarrow aaaS$$

$$\Rightarrow aaaaS$$

$$\Rightarrow aaaa\epsilon$$

$$\Rightarrow aaaa.$$

The Language has strings  $\{ \epsilon, a, aa, aaa, \dots \}$

EX: 4 Find Language and derive 'abbaaba' from the following grammar.

Terminals : a, b

Non-terminals : S, X

Productions :  $S \rightarrow Xaax$

$X \rightarrow aX \mid bX \mid \epsilon$

Solution:-

$$S \Rightarrow Xaax$$

$$\Rightarrow aXaax$$

$$\Rightarrow abXaax$$

$$\Rightarrow abbXaax$$

$$\Rightarrow abb\epsilon aax$$

$$\Rightarrow abbaabx$$

$$\Rightarrow abbaabax$$

$$\Rightarrow abbaaba\epsilon$$

$$\Rightarrow abbaaba.$$

EX: 5 Give the language defined by grammar  $G$ .

$G = \{ \{S, C\}, \{a, b\}, P, S \}$ , where  $P$  is given by  
 $S \rightarrow aCa$   $C \rightarrow aCa \mid b$

Solution:-

$$\begin{aligned} S &\Rightarrow aCa \\ &\Rightarrow aaCa \\ &\Rightarrow aaaa \end{aligned}$$

$$L(G) = a^n b a^n \text{ for } n \geq 1.$$

EX: 6 Give the language defined by grammar  $G$ .

$G = \{ \{S\}, \{0, 1\}, P, S \}$ , where  $P$  is given by  $S \rightarrow 0S1 \mid \epsilon$

Solution:-

$$\begin{aligned} S &\Rightarrow 0S1 \\ &\Rightarrow 00S11 \\ &\Rightarrow 000S111 \\ &\Rightarrow 000\epsilon111 \end{aligned}$$

$$L(G) = 0^n 1^n \quad \forall n \geq 0$$

Leftmost and Rightmost Derivations:-

EX: 1 Consider the CFG  $G = (\{S, X\}, \{a, b\}, P, S)$ , where

Productions are  $S \rightarrow baxaS \mid ab$

$X \rightarrow Xab \mid aa$ .

Find LMD and RMD for string  $w = baaqababaab$ .

The following is a LMD:

$S \Rightarrow baxas$   
 $\Rightarrow baxabas$   
 $\Rightarrow baxababas$   
 $\Rightarrow baaaaababas$   
 $\Rightarrow baaaababab.$

The following is a RMD:

$S \Rightarrow baxas$   
 $\Rightarrow baxaab$   
 $\Rightarrow baxabaaab$   
 $\Rightarrow baxababaaab.$   
 $\Rightarrow baaaababaaab.$

Note:- Any word that can be generated by a given CFG can have LMD/RMD.

H.W

Ex: 2 Consider the CFG

$S \rightarrow aB \mid bA$

$A \rightarrow a \mid aS \mid bAA$

$B \rightarrow b \mid bS \mid aBB$

Find LMD and RMD for string  $w = aabbabba.$

Problems:-

Problem:1 Find LMD and RMD for string 00101 in

grammar given below  $S \rightarrow B|A$ ,  $A \rightarrow 0A|\epsilon$ ,  $B \rightarrow 1B|0B|\epsilon$

Solution:- LMD, RMD are same.

- $S \Rightarrow B$
- $\Rightarrow 0B$
- $\Rightarrow 00B$
- $\Rightarrow 001B$
- $\Rightarrow 0010B$
- $\Rightarrow 00101B$
- $\Rightarrow 00101\epsilon$
- $\Rightarrow 00101.$