

Introduction :-

- we have learnt effective use of regular expressions for defining regular languages. But there are certain languages which cannot be defined by regular expressions eg. we cannot write a regular expression for checking well defined formed parenthesis etc. so in order to cover all these things and more complicated features leads us to the idea of a context-free ^{Grammar} languages for the languages.
- while compiling the program, it is checked syntactically first and semantically latter. For checking the correctness of the program syntactically we need CFG.

: Definition :-

→ A context free grammar(CFG) is defined as 4-tuples (V_N, Σ, P, S)

↓
starting
non-terminal

where, V_N = set of non-terminal symbol

Σ = set of terminal symbol

P = Production rule

S = starting non-terminal symbol.

Production rule is in the form,

one non-terminal \rightarrow finite string of terminals and / or non-terminal.

i.e $A \rightarrow \alpha$ where $A \in V_N$ and $\alpha \in (V_N \cup \Sigma)^*$

* Application

- ⇒ The language generated by the context-free grammars are called context-free languages are applied in parser design.
- ⇒ It is also useful for describing block structure in programming languages.
- ⇒ CFG's are useful for describing arithmetic expressions, with arbitrary nesting of balanced parenthesis.

* Derivations

- ⇒ The production rules are used to derive certain strings. The generation of language using specific rules is called derivation.
- ⇒ notations : \Rightarrow and $\xrightarrow{*}$
- ⇒ If $\alpha \Rightarrow \beta$ be a production of P in CFG and a and b are strings in $(V_N \cup \Sigma)^*$ then $\alpha a b \xrightarrow{*} \alpha \beta b$.

Example 1 Construct a context-free grammar G_1 generating all integers (with sign).

Soln Let $G_1 = (V_N, \Sigma, P, S)$ where $V_N = \{S, \langle \text{sign} \rangle, \langle \text{digit} \rangle, \langle \text{integer} \rangle\}$
 $\Sigma = \{0, 1, 2, 3, \dots, 9, +, -\}$

P consists of $S \rightarrow \langle \text{sign} \rangle \langle \text{integer} \rangle$

$\langle \text{sign} \rangle \rightarrow + \text{ or } -$

$\langle \text{integer} \rangle \rightarrow \langle \text{digit} \rangle \langle \text{integer} \rangle \mid \langle \text{digit} \rangle$

$\langle \text{digit} \rangle \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$L(G_1)$ is the set of all integers.

$L(G)$ is the Language of G. It is the set of terminal strings that have derivations from the start symbol.

four components in CFG:-

1. There is the finite set of symbols that form the strings of language being defined. We call this alphabet the terminal or terminal symbol.
2. There is a finite set of variables, also called non-terminal symbol. Each variable represent the language i.e. set of strings.
3. One of the variable represent the language being defined. It is called start symbol. It is denoted by S .
4. There is a finite set of rules or productions that represent the recursive definition of a language. Each production consist of:
 - (a) A variable that is being partially defined by the production. This variable is often called the head of production.
 - (b) The production symbol \rightarrow
 - (c) A strings of zero or more terminals and variables. This string is called the body of production; represents one way to form strings in the language of the variable of the head

note:-

- ↳ The variable symbols are represented by capital letters.
- ↳ The terminal symbols are represented by lowercase letters.
- ↳ One of the variable is designed as a start variable, it usually occurs in the left hand side of the production rule.

example

$$S \rightarrow \epsilon$$

$$S \rightarrow 0SL$$

This is the CFG defining the grammar of all the strings with equal no. of 0's followed by equal no. of 1's.

Here, Two rules define the production P,

$\epsilon, 0, L$ are the terminal defining Σ or T,

S is the non-terminal symbol defining V.

and S is the start symbol from where production starts.

Example following is a CFG for the language $L = \{w c w^R | w \in (a, b)^*\}^2$

Solution Let G_1 be the CFG for the Language $L = \{w c w^R | w \in (a, b)^*\}$

$$G_1 = (V_N, \Sigma, P, S)$$

$$\text{Here, } V_N = \{S\}$$

$$\Sigma = \{a, b, c\}$$

and P is given by

$$S \rightarrow aSa$$

$$S \rightarrow bSb$$

$$S \rightarrow c$$

let us check that abcbba can be derived from given CFG

$$S \Rightarrow aSa \quad (\text{use the } S \rightarrow aSa)$$

$$\Rightarrow abSba \quad (\text{use the } S \rightarrow bSb)$$

$$\Rightarrow abbSbbba \quad (\text{use the } S \rightarrow bSb)$$

$$\Rightarrow abcbba \quad (\text{use the } S \rightarrow c)$$

Example 2 Write a CFG, which generates string of balanced parenthesis.

Solution This grammar will accept the balanced left and right parenthesis. For example $() ()$ is acceptable. $(((())))$ this is also acceptable.

Let us design the CFG for this,

Let CFG be,

$$G_1 = (V_N, \Sigma, P, S) \text{ where}$$

$$V_N = \text{set of non-terminal} = \{S\}$$

$$\Sigma = \{C,)\}$$

and set of production P is given by

$$S \rightarrow SS$$

$$S \rightarrow (S)$$

$$S \rightarrow \epsilon$$

Now we see what this grammar generates

$$S \rightarrow SS$$

$$\rightarrow (S)S$$

$$\rightarrow (S)SS$$

$$\rightarrow (S)(S)S$$

$$\rightarrow (S)(S)(S)$$

$$\Rightarrow () () ()$$

Thus above grammar will always generate balanced pair of parenthesis.

Example write a CFG, which generates palindrome for binary numbers.

Solution Grammar will generate palindrome for binary numbers, that is 00, 010, 11, 101...

let CFG be $G = (V_N, V_T, \Sigma, P, S)$

$V_N = \text{set of non-terminals} = \{S\}$

$\Sigma = \text{set of terminals} = \{0, 1\}$

and P is the Production Rule defined as

$$S \rightarrow 0S0 \mid 1S1$$

$$S \rightarrow 0 \mid 1$$

Obviously this grammar generates palindrome for binary numbers it can be seen by the following derivation.

$$S \Rightarrow 0S0$$

$$\Rightarrow 01S01$$

$$\Rightarrow 010S010$$

$\Rightarrow 0101010$ which is palindrome.

Example write a CFG for a regular expressions.

$$\tau = 0^* L (0+1)^*$$

Solution let us analyse regular expression

$$\tau = 0^* L (0+1)^*$$

Clearly regular expression is the set of strings which starts with any number of 0's followed by one and end with any combination of 0's and 1's.

Let CFG be $G_1 = (V_N, \Sigma, P, S)$

where,

$$V_N = \text{set of non-terminal symbol}$$
$$= \{S, A, B\}$$

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

and productions P are defined as

$$S \rightarrow A \sqcup B$$

$$A \rightarrow 0A \sqcup \epsilon$$

$$B \rightarrow 0B \sqcup 1B \sqcup \epsilon$$

Let us see the derivation of the string 00101

$$S \Rightarrow A \sqcup B$$

$$\Rightarrow 0A \sqcup B$$

$$\Rightarrow 00A \sqcup 01B$$

$$\Rightarrow 00101$$

Clearly, G_1 is the CFG for regular expression σ .

Example Write a CFG which generates strings having equal number of a's and b's.

Solution

Let CFG be $G_1 = (V_N, \Sigma, P, S)$

$$V_N = \{S\}$$

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

where P is defined as,

$$S \rightarrow aSbS \sqcup bSaS \sqcup \epsilon$$

Let us derive a string $co = bbabaa bbba$

$$\begin{aligned} S &\Rightarrow b \underline{S} a S \\ &\Rightarrow b b \underline{S} a S a S \\ &\Rightarrow b b a S b \underline{S} a S a S \\ &\Rightarrow b b a S b a a S b S b S a S a S \\ &\Rightarrow b b a b a a S b S b S a S a S \\ &\Rightarrow b b a b a a b a a S b S a S a S \\ &\quad b b a b a a b b S a S a S \\ &\Rightarrow b b a b a a b a b S a S a S \\ &\Rightarrow b b a b a a b b a a S \\ &\Rightarrow b b a b a a b b a a \end{aligned}$$

Example Design a CFG for the language $L = \{a^n b^m : n \neq m\}$

solution If $n \neq m$ then there are only two cases are possible

Case 1.

$n > m$

let us say language L on condition $n > m$ is

$$L_1 \text{ and } L_2 = \{a^n b^m : n > m\}$$

let G_L be the CFG for the language L_1 ,

$$G_L = (V_n^L, \Sigma^L, P^L, S^L)$$

$$V_n^L = \{S_1, A, S\}$$

$$\Sigma^L = \{a, b\}$$

Then production P_L are defined as,

$$S^L \rightarrow AS_L$$

$$S_L \rightarrow aS_L b | \epsilon$$

$$A \rightarrow aAa$$

Case 2

$$n < m$$

let us say L on condition $n < m$ is L_2

$$L_2 = \{0^n 1^m : n < m\}$$

and let CFG for G_2 be $G_2 = (V_N^2, \Sigma^2, P^2, S)$

$$V_N^2 = \{S^2, S_2, B\}$$

~~$$\Sigma^2 = \{a, b\}$$~~

P^2 are defined as,

$$S^2 \rightarrow S_2 B$$

$$S_2 \rightarrow aS_2 b | \epsilon$$

$$B \rightarrow bBb$$

By combining G_L and G_2 we can write the CFG for L as

$$S \rightarrow S^L | S^2$$

* Derivation

⇒ CFG generates string according to the following process called derivation.

- ① we start by writing down the start symbol of the CFG.
 - ② At each step of the derivation, we may replace any non-terminal symbol of the string generated so far by the right hand side of any production that has the symbol on the left.
 - ③ The process ends when it is impossible to apply a step of type described in ②.
- ⇒ If the string at the end of this process consists entirely of terminals, it is a string in the language generated by the grammar called yield.

⇒ There are two approaches of derivation

1. Body to head approach (Bottom up)
2. Head to Body (Top down)

1. Body to head :-

⇒ Here, we take strings known to be in the language of each of the variables of the body, concatenate them in the proper order with any terminals appearing in the body, the resulting string is the language of the variable in the head.

Consider grammar,

$$S \rightarrow S+S$$

$$S \rightarrow S|S$$

$$S \rightarrow (S)$$

$$S \rightarrow S-S$$

$$S \rightarrow S*S$$

$$S \rightarrow a$$

Here given $a + (axa)la-a$
now, applying body to head approach,

S.N	String inferred	Variable Production	String used.
1.	a	S	$S \rightarrow a$
2.	axa	S	$S \rightarrow S+S$
3.	(axa)	S	$S \rightarrow (S)$
4.	$(axa)la$	S	$S \rightarrow SIS$
5.	$(axa)la-a$	S	$S \rightarrow S-S$
6.	$a + (axa)la-a$	S	$S \rightarrow S+S$

Thus in this approach we start with any terminal appearing in the body and use the available rules from body to head.

2. Head to Body :-

Here, we used production from head to body. we expand the start symbol using a production, whose head is the start symbol. Here we expand the resulting string until all strings of terminal are obtained. Here we have two approaches

(i) Left most Derivation:-

Here, leftmost symbol(variable) is replaced first

(ii) Right most Derivation :-

Here, rightmost symbol is replaced first.

Example Consider the grammar G1 with production
 $S \rightarrow ass \mid b$

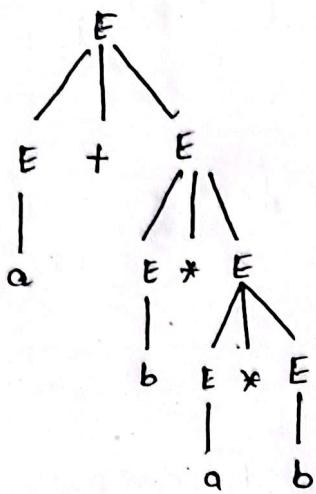
Now, we have

$$\begin{aligned} S &\rightarrow ass \\ &\Rightarrow aasss \\ &\Rightarrow aabss \\ &\Rightarrow aabasss \\ &\Rightarrow aabbss \\ &\Rightarrow aababss \\ &\Rightarrow aababbb \end{aligned}$$

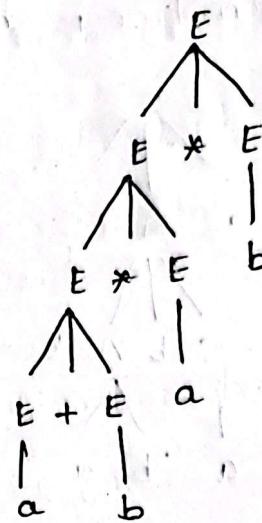
The sequence followed is left-most derivation

Now right-most derivation:

$$\begin{aligned} S &\rightarrow ass \\ &\Rightarrow asb \\ &\Rightarrow aassb \\ &\Rightarrow aasassb \\ &\Rightarrow aaSasbb \\ &\Rightarrow aaSabb \\ &\Rightarrow aababbb \end{aligned}$$



LMD parse tree



RMD parse tree

Example Consider the grammar G.

$$S \rightarrow A1B$$

$$A \rightarrow 0A1B$$

$$B \rightarrow 0B1B1B$$

1. Construct the parse tree for 00101.

LMD

$$S \rightarrow \underline{A}1B$$

$$\Rightarrow 0\underline{A}1B$$

$$\Rightarrow 00\underline{A}1B$$

$$\Rightarrow 00\underline{B}1B$$

$$\Rightarrow 001\underline{B}$$

$$\Rightarrow 0010\underline{B}$$

$$\Rightarrow 00100$$

$$\Rightarrow 00101$$

RMD

$$S \rightarrow \underline{A}1B$$

$$\Rightarrow A\underline{1}B$$

$$\Rightarrow A1\underline{B}$$

$$\Rightarrow A10\underline{B}$$

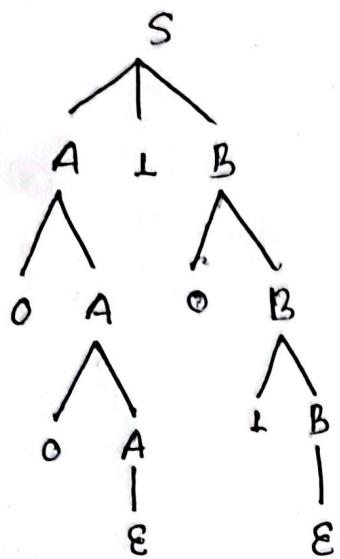
$$\Rightarrow A101\underline{B}$$

$$\Rightarrow A1010\underline{B}$$

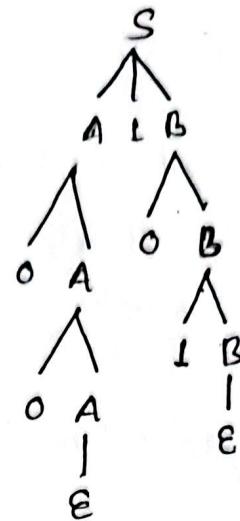
$$\Rightarrow 0A101$$

$$\Rightarrow 00A101$$

$$\Rightarrow 00101$$



LMD parse tree



RMD parse tree

* Ambiguity in Grammar:-

⇒ A grammar $G_1 = (V, \Sigma, P, S)$ is said to be ambiguous if there is a string $w \in L(G)$ for which we can derive two or more distinct derivation trees at S and yielding w . In other words, a grammar is ambiguous if there exist more than one leftmost derivation or more than one rightmost derivation for the same string.

Example

$$S \rightarrow ABaaB$$

$$A \rightarrow aAa$$

$$B \rightarrow b$$

for any string aab ;

we have two leftmost derivation as.

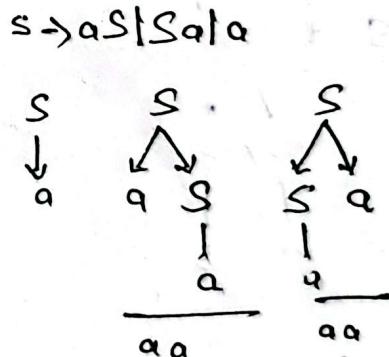
$$S \rightarrow \underline{AB}$$

$$\rightarrow AaB$$

$$\rightarrow aaB$$

$$\rightarrow aab$$

$$\text{Also, } S \rightarrow aab \\ \rightarrow aab$$



Here, we see there are two ways to generate string aab . Thus ambiguous.

for example

$G_1 = (\{S\}, \{a, b, +, *\}, P, S)$ where P consist
 $S \rightarrow S+S \mid S* \mid a \mid b$, there are two derivation trees.

for $a+ab$.

Soln ~~we have two left most derivation tree~~

$$S \rightarrow S+S$$

$$\rightarrow a+S$$

$$\rightarrow a+S*S$$

$$\rightarrow a+a*S$$

$$\rightarrow a+a+b$$

Also,

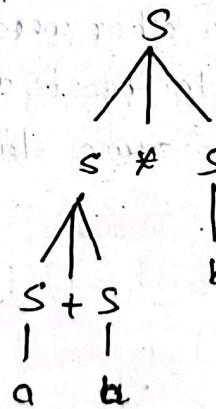
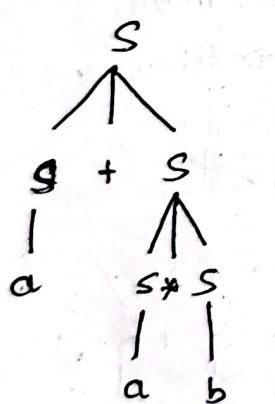
$$S \rightarrow S*S$$

$$\rightarrow S*b$$

$$\rightarrow S+S*b$$

$$\rightarrow S+a*b$$

$$\rightarrow a+a*b$$



Thus, the given grammar is ambiguous.

* Removal of Ambiguity

- ⇒ In programming Languages, where there should be only one interpretation of each statement, ambiguity must be removed when possible.
- ⇒ often we can achieve this by rewriting the grammar in an equivalent unambiguous form.

Example. Grammar

$$\begin{aligned} E &\rightarrow I \\ E &\rightarrow E+E \\ E &\rightarrow E \times E \\ E &\rightarrow (E) \\ I &\rightarrow a+b+c \end{aligned}$$

is ambiguous and the ambiguity is due to the precedence of operator.
If we correct the precedence then ambiguity may be removed.

Here, two causes of ambiguity in the grammar.

1. The precedence of operators is not respected.
2. The sequence of identical operators can group either from left or from the right. We would see two different parse trees for $E+E+E$. Since addition is associative, it does not matter whether we group from the left or the right, but to remove ambiguity we must take one.

The unambiguous grammar is

$$E \rightarrow T$$

$$T \rightarrow F$$

$$F \rightarrow I$$

$$E \rightarrow E + T$$

$$T \rightarrow T * F$$

$$F \rightarrow (E)$$

$$I \rightarrow a b | c$$

The sole parse tree for $a+b*c$ is

