

# UNIT 3 : CONTEXT FREE GRAMMAR AND CONTEXT FREE LANGUAGE

## \* BASIC Elements of grammar:

A grammar is a 4 tuple  $G = \{V, T, P, S\}$  where  $V$  = finite set of non terminal symbol,  $T$  = finite set of terminal symbol,  $P$  = set of production rules and  $S$  = start symbol symbol.

The  $P$  is of the form  $A \rightarrow B$  where  $A$  and  $B \in (VUT)^*$

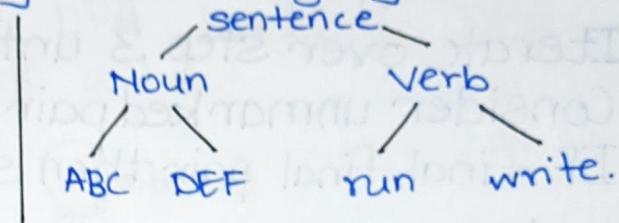
The grammer is normally represented using derivation tree.

The production rules of defining natural language are as follows :

$\text{sentence} \rightarrow \text{Nouns, Verbs}$

$\text{Nouns} \rightarrow 'ABC', 'DEF'$

$\text{Verbs} \rightarrow \text{Run, write.}$



## \* FORMAL DEFINITION OF CONTEXT FREE GRAMMAR

The context free grammar can be formally defined as a set denoted by  $G = \{V, T, P, S\}$  where  $V \nsubseteq T$  are set of non terminals and terminals respectively.

$P$  is a set of production rules, where each production rule is in the form of

Non terminals  $\rightarrow$  Non terminal

Non terminals  $\rightarrow$  Terminals

$S$  is a start symbol.

$$P = \{ S \rightarrow S+S, S \rightarrow 5^*S, S \rightarrow (S), S \rightarrow 4 \}$$

If the language is  $4+4*4$  then we can use the production rules given by  $P$ . The Start symbol is  $S$ . The Number of non terminals in the rules  $P$  is one and the only non-terminal i.e  $S$ . The terminals are  $+$ ,  $*$ ,  $($ ,  $)$ , and  $4$ .

We are using following Conventions.

1) The Capital letters are used to denote the non terminals.

2) The lower case letters are used to denote the terminals.

## SENTENTIAL FORM

Consider  $G = (V, T, P, S)$  be a context-free grammar then we can derive a string  $w$  from it. This  $w$  can be obtained  $(VUT)^*$  where  $V$  denotes the set of non-terminal symbols

The derivation of  $w$  from start symbol  $S$  can be written as  $S \xrightarrow{*} w$  which is called as sentential form

If  $S \xrightarrow{tm} w$  then  $w$  is a left sentential-form

If  $S \xrightarrow{rm} w$  then  $w$  is a right sentential-form

Example : Consider the grammar  $S \rightarrow S+S \mid S^*S \mid 4$

Then to derive  $4+4*4$  we use above grammar as

$$S \xrightarrow{lm} S+S \Rightarrow S+S \xrightarrow{lm} S^*S \Rightarrow S^*S \xrightarrow{lm} 4+S^*S \Rightarrow 4+S^*S \xrightarrow{lm} 4+4^*S$$

Here  $S+S^*S$  is called left sentential-form

$$S \xrightarrow{rm} S+S \Rightarrow S+S \xrightarrow{rm} S^*S \Rightarrow S^*S \xrightarrow{rm} 4 \Rightarrow S+4^*S \xrightarrow{rm} 4+4^*4$$

Here  $S+S^*S$  is called right sentential-form.

Thus we can obtain desired language  $L$  by certain rules

~~Let us~~

Context free Language :

The language of grammar that can be defined by creating the production rules for the given condition

The language generated by context free grammar is called context free language. (CFL)

Let us understand the CFL with the help of various example.

i] Construct the CFG for the language having any number of  $a$ 's over the set.

As we know  $RE = a^*$

lets build production rules for same

$$S \rightarrow aS \text{ rule 1}$$

$$S \rightarrow \epsilon \text{ rule 1}$$

Now if we want string 'aaaa' we can start with start symbols

$$G = \{ \{S\}, \{a\}, \{P\}, S \}$$

$S$

$a a S$

$a a a S$

$a a a a S$

$a a a a \epsilon$

=  $aaaa$ . | we can also do r.e =  $a^* a^*$  |  
|  $S \rightarrow \epsilon$  is possible. |

## DERIVATION

The production rules are used for deriving certain strings. If  $G = (V, T, P, S)$  be some context-free grammar then generation of some language using specific rules is called derivation.

As we know derivation of any string means replacement of non-terminal by its appropriate definition. There may be a situation in which there are many non-terminals then which non-terminal should be replaced by its definition is sometimes confusing to decide.

Here's hence we apply two methods of deriving.

The left most derivation is a derivation in which the leftmost non terminal is replaced first to form the sentential form.

The right most derivation is a derivation in which the rightmost non terminal is replaced first to form the sentential form.

Example.

$$S \rightarrow xyx$$

$$S \rightarrow ayx$$

$$S \rightarrow abx$$

$$S \rightarrow aba.$$

leftmost derivation

$$S \rightarrow xyx$$

$$S \rightarrow xya$$

$$S \rightarrow Xba$$

$$S \rightarrow aba.$$

rightmost derivation

## PARSE TREE :

Parse tree is a graphical representation for the derivation of the given production rules for a given CFG. It is simple way to show how the derivation can be done to obtain some strings from given set of production rules. The parse tree is also called derivation tree.

Properties of parse tree.

- 1) The root node is always node indicating start symbol
- 2) The derivation is read from left to right.
- 3) The leaf nodes are always terminal nodes.
- 4) Interior nodes are always non terminal nodes.

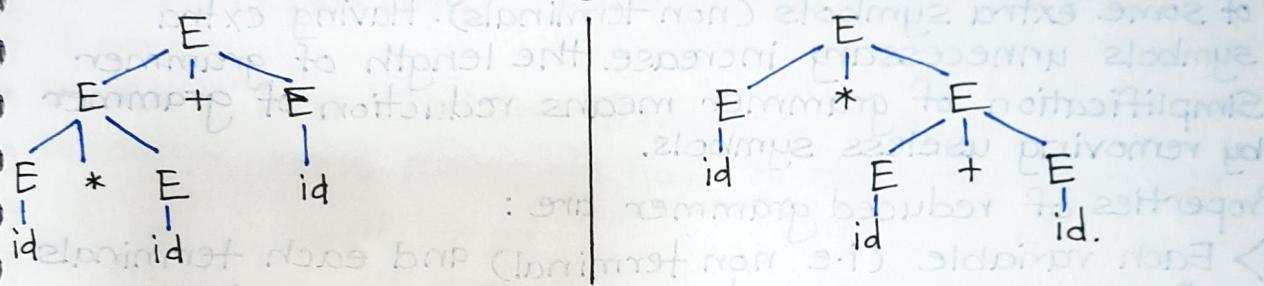
## Ambiguous Grammer.

The grammer can be derived in either leftmost derivation or rightmost derivation one can draw a derivation tree called as parse tree or syntax tree based on these derivations. The parse tree has to be unique even if the derivation is leftmost or rightmost.

If there exists more than one parse trees for a given grammer, that means there can be more than one leftmost and rightmost derivation possible and then that grammer is said to be ambiguous grammer

example :  $G = (V, T, P, S)$  where  $V = \{E\}$ ,  $T = \{\text{id}\}$ ,  $P = \{E \rightarrow E+E, E \rightarrow E^*E, E \rightarrow \text{id}\}$ ,  $S = \{E\}$

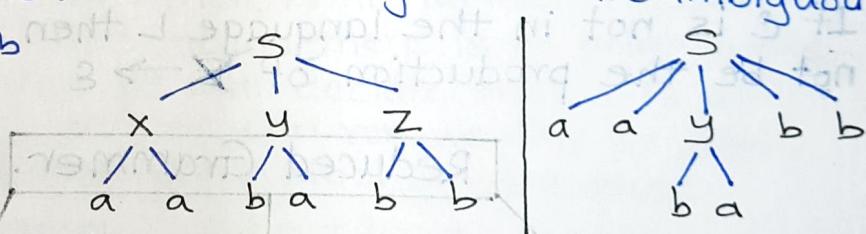
Now, if the string is  $\text{id}^* \text{id} + \text{id}$  then we can draw the two different Parse trees. indicating our  $\text{id}^* \text{id} + \text{id}$ .



## Inherent Ambiguity:

A context free grammer is called inherently ambiguously, if all the production rules in this grammer are ambiguous

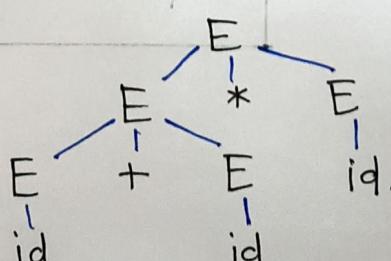
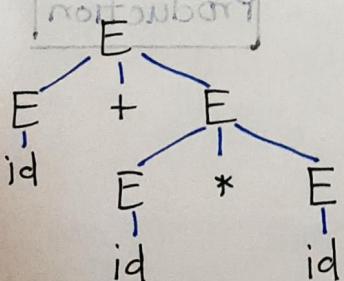
$S \rightarrow xyz \mid aaybb$   
 $x \rightarrow aay \mid aa$   
 $y \rightarrow baz \mid ba$   
 $z \rightarrow bzb \mid bb$



## Removal of ambiguity.

Consider the ambiguous grammer as

$E \rightarrow E+E \mid E^*E \mid \text{id}$



for removing the ambiguity we will apply one rule  
 If the grammar has left associative operator  
 (Such as  $+, -, *, /$ ) then induce the left recursion and  
 if the grammar has right associative operator  
 (Exponential operator) then induce right recursion.

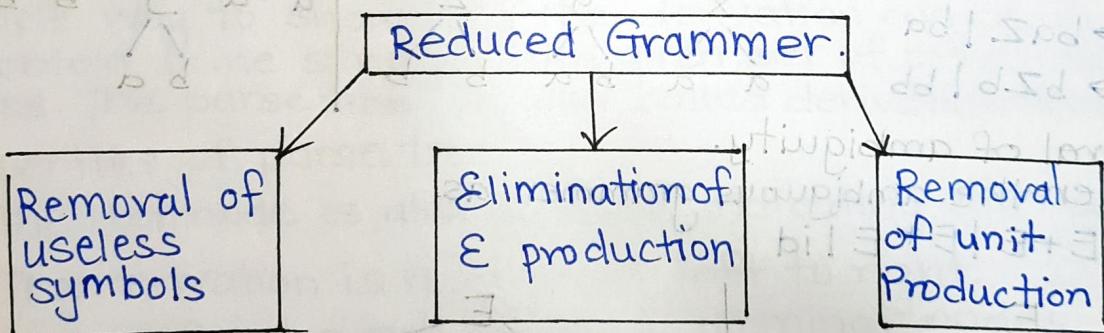
$$\begin{array}{l|l} E \rightarrow E + T & E \rightarrow E + T \mid T \\ E \rightarrow T S & T \rightarrow T * F \mid F \\ T \rightarrow T * F & F \rightarrow id \\ T \rightarrow F & \\ F \rightarrow id & \end{array}$$

### Simplification of CFG

As we have seen various language can effectively be represented by Context free grammar. All the grammars are not always optimized. That means grammar may consist of some extra symbols (non terminals). Having extra symbols unnecessary increase the length of grammar. Simplification of grammar means reduction of grammar by removing useless symbols.

Properties of reduced grammar are :

- 1> Each variable (i.e non terminal) and each terminals of G appears in the derivation of some word in L
- 2> There should not be only production as  $x \rightarrow y$  where x and y are non terminals.
- 3> If  $\epsilon$  is not in the language L then there need not be the production of  $x \rightarrow \epsilon$



## REMOVAL OF USELESS SYMBOL

Any symbol is useful when it appears when it appears on the right hand side in the production rule and generates some terminal string. If no such derivation exists then it is supposed to be an useless symbol.

A symbol  $P$  is useful if there exists some derivation. In the following form.

$$S \xrightarrow{*} \alpha PB \quad \& \quad \alpha PB \xrightarrow{*} w$$

Then  $P$  is said to be useful symbol.

Where  $\alpha$  &  $B$  may be some terminal or non terminal symbol and will help us to derive certain string  $w$ .

example :

$$G = (V, T, P, S) \text{ where } V = \{S, T, X\}, T = \{0, 1\}$$

$$S \rightarrow OT \mid IT \mid X \mid 0 \mid 1 \quad \text{rule 1}$$

$$T \rightarrow 00 \quad \text{rule 2}$$

In above CFG, nonterminals are  $S, T, X$

To derive some string we have to start from symbols

$$\begin{array}{ll} S \\ \overline{OT} \\ 000 \end{array} \quad \begin{array}{l} S \rightarrow OT \\ T \rightarrow 00 \end{array}$$

Thus we can reach to certain string after following these rules.

But if  $S \rightarrow X$  then there is no further rule as a definition to  $X$ . That means there is no point in the rule  $S \rightarrow X$ . Hence we can declare that  $X$  is a useless symbol. and we can remove this so after removal of this useless symbol CFG becomes.

$$G = (V, T, P, S) \text{ where } V = \{S, T\}, T = \{0, 1\}$$

and  $P \{ S \rightarrow OT \mid IT \mid 0 \mid 1 \}$ ,

$$T \rightarrow 00 \}$$

$S$  is a start symbol

## Elimination of $\epsilon$ Productions from grammar.

As we have seen in finite Automata and regular expression that  $\epsilon$  or a null string indicates a string with no value. we have also seen that even though some NFA contains  $\epsilon$  moves we can convert that NFA without  $\epsilon$  moves. Even in Context Free grammar, If at all there is  $\epsilon$  production. we can remove it, without changing the meaning of the grammar.

Thus  $\epsilon$  productions are not necessary in grammar.

$$S \rightarrow 0S \mid 1S \mid 0 \mid 1, \text{ for } S \rightarrow 0S \mid 1S \mid \epsilon.$$

example:

$$S \rightarrow XYX, X \rightarrow 0X \mid \epsilon, Y \rightarrow 1Y \mid \epsilon$$

$$\therefore S \rightarrow XYX$$

$$\text{if 1st } X \text{ in rhs} = \epsilon \quad \therefore S \rightarrow XY$$

$$\text{if last } X \text{ in Lhs} = \epsilon \quad \therefore S \rightarrow XYX$$

$$\text{if } Y = \epsilon \quad \therefore S \rightarrow XX$$

$$\text{if } X \& Y \text{ are } \epsilon \text{ then } \therefore S \rightarrow X$$

$$\text{when both } X \text{ replaced by } \epsilon \therefore S \rightarrow Y$$

$$S \rightarrow XY \mid YX \mid XX \mid Y \mid X$$

$$\therefore X \rightarrow 0X$$

if we place  $\epsilon$  at rhs of  $X$  then

$$X \rightarrow 0$$

$$X \rightarrow 0X \mid 0$$

$$Y \rightarrow 1Y \mid 1$$

collectively we can write it as

$$S \rightarrow XY \mid YX \mid XX \mid X \mid Y$$

$$X \rightarrow 0X \mid 0$$

$$Y \rightarrow 1Y \mid 1$$

## Removing Unit Productions

The unit productions are the productions in which one non-terminal gives another non-terminals

example: if  $x \rightarrow y$ ,  $y \rightarrow z$ ,  $z \rightarrow x$  then  $x, y, z$  are unit productions

Example: If  $S \rightarrow 0A|1B|c$ ,  $A \rightarrow 0S|00$ ,  $B \rightarrow 1|A$   
 $c \rightarrow 01$

$c \rightarrow 01$  is a unit production so write a rule for what  $c$  gives.

$\therefore S \rightarrow 0A|1B|01$

$B \rightarrow A$  is also unit production so write a rule for what  $A$  gives

$S \rightarrow 0A|1B|01$ ,  $A \rightarrow 0S|00$ ,  $B \rightarrow 1|0S|00$ ,  $c \rightarrow 01$

## NORMAL FORMS :

There are two important normal forms

1. Chomsky's Normal Form

2. Griebach Normal form

1. Chomsky's Normal form:

It can be defined as :-

Non terminal  $\rightarrow$  Non terminal . Non-terminal  
Non terminal  $\rightarrow$  terminal

The given CFG must be converted into above format

then we can say that the grammar is in C<sub>FG</sub> CNF

Before converting the grammar into CNF it must be in reduced form.

That means remove all the useless symbols,  $\epsilon$ -productions and unit productions from it

Thus reduced grammar can be then converted to CNF

Example: Find CNF equivalent of.

$$S \rightarrow aA bB$$

$$A \rightarrow aA | a$$

$$B \rightarrow bB | b$$

Step 1: Simplification of grammar

The grammar is already in a simple form without

1)  $\Sigma$ -productions

2) Unit productions

3) Useless symbols

Step 2: Every symbol in  $\alpha$ , in a production of the form  $A \rightarrow \alpha$  where  $|\alpha| \geq 2$  should be a variable.

This can be done by adding two productions

$Ca \rightarrow a$  and  $Cb \rightarrow b$

The set of productions after the above changes is

$$P_2 = \left\{ \begin{array}{l} S \rightarrow Ca A \epsilon b B \\ A \rightarrow Ca A | a \\ B \rightarrow Cb B | b \\ Ca \rightarrow a \\ Cb \rightarrow b \end{array} \right\}$$

Step 3: Convert to CNF

$$S \rightarrow Ca C_1$$

$$C_1 \rightarrow AC_2$$

$$C_2 \rightarrow Cb B$$

$$A \rightarrow Ca A | a$$

$$B \rightarrow Cb B | b \quad \text{Corresponding to}$$

$$\left\{ \begin{array}{l} Ca \rightarrow a \\ Cb \rightarrow b \end{array} \right\} \quad S \rightarrow Ca A Cb B$$

$$P_3 = \left\{ \begin{array}{l} S \rightarrow Ca C_1, C_1 \rightarrow AC_2, C_2 \rightarrow Cb B \\ A \rightarrow Ca A | a, B \rightarrow Cb B | b, \\ Ca \rightarrow a, Cb \rightarrow b \end{array} \right\}$$

### GREIBACH NORMAL FORM

A context free grammar  $G = (V, T, P, S)$  is said to be in GNF if every production is of the form

$$A \rightarrow a \alpha$$

where  $a \in T$  is a terminal and  $\alpha$  is a string of zero or more variables.

Non Terminal  $\rightarrow$  one terminal, any number of non terminal.

The Language should be without  $\epsilon$ .  
Right hand side of production should start with a terminal followed by a string of non-terminals of length zero or more.

Algorithm for conversion from CFG to GNF

1. Eliminate  $\epsilon$ -productions, unit productions and useless symbols from the grammar.
2. In production of the form  $A \rightarrow x_1 x_2 \dots x_i \dots x_n$  other than  $x_i$ , every other symbol should be a variable  $x_i$  could be a terminal.
3. Rename variables as  $A_1, A_2, A_3 \dots A_n$  to be a production of  $A$ .
4. Modify the productions to ensure that if there is a production  $A_i > A$  then  $i$  should be  $\leq j$ .  
If there is a production  $A_i > A_j \alpha$  with  $i > j$  then we have must generate productions substituting for  $A_j$ .
5. Repeating step 4 : Several times will guarantee that for every production  $A_i \rightarrow A_j \alpha \quad \nexists i \leq j$ .
6. Removing left recursion from every production of the form  $A_k \rightarrow A_k \alpha \dots \beta$ -productions should be added to remove left recursion.
7. Modify  $A_i$ -productions to the form  $A_i \rightarrow a \alpha$  where  $a$  is a terminal and  $\alpha$  is a string of non-terminals.
8. Modify  $B_i$ -productions to the form  $B_i \rightarrow a \alpha$  where  $a$  is a terminal and  $\alpha$  is a string of non-terminals.

CONSTRUCT A GRAMMAR IN CFG GNF in which is equivalent to the grammar

$$S \rightarrow AA | a \quad A \rightarrow SS | b.$$

Step 1: grammar is already in simple form.

We can proceed to renaming of variables

Variable S and A are renamed as A<sub>1</sub> and A<sub>2</sub> respectively

The set of productions after renaming becomes,

$$A_1 \rightarrow A_2 A_2$$

$$A_2 \rightarrow A_1 A_1$$

$$A_1 \rightarrow a$$

$$A_2 \rightarrow b.$$

Step 2: Every production of the form A<sub>i</sub> → A<sub>j</sub> with i > j must be modified to make i ≤ j

A<sub>2</sub>-production A<sub>2</sub> → A<sub>1</sub>A<sub>1</sub> should be modified to make i < j modified.

we must substitute A<sub>2</sub>A<sub>2</sub>|a for A<sub>1</sub>

we shouldn't touch the 2nd A<sub>1</sub> of A<sub>1</sub>A<sub>1</sub>

$$A_2 \rightarrow A_1 A_1 \rightarrow A_2 \rightarrow A_2 A_2 A_1$$

$$A_2 \rightarrow a A_1$$

The resulting set of productions is :

$$A_1 \rightarrow A_2 A_2 | a \quad A_2 \rightarrow A_2 A_2 A_1 | a A_1 | b$$

Step 3 : renaming left recursions.

The A<sub>2</sub> productions A<sub>2</sub> → A<sub>2</sub>A<sub>2</sub>A<sub>1</sub> | aA<sub>1</sub> | b contains left recursions from A<sub>2</sub> productions can be removed through introduction of B<sub>2</sub>-production

$$A_2 \rightarrow a A_1 B_2 | b B_2$$

$$B_2 \rightarrow A_2 A_1 B_2 | A_2 A_1.$$

The resulting set of productions is.

$$A_1 \rightarrow A_2 A_2 | a.$$

$$A_2 \rightarrow a A_1 B_2 | a B_2 | a A_1 | b.$$

$$B_2 \rightarrow A_2 A_1 B_2 | A_2 A_1$$

Step 4:  $A_2$  production are in GNF

$A_1 \& B_2$  productions can be converted to GNF  
with the help of  $A_2$ -productions

$A_2 \rightarrow aA_1B_2 \mid bB_2 \mid aA_1 \mid b \dots$  in GNF

$A_1 \rightarrow A_2A_2$

Substitute

$A_1 \rightarrow aA_1B_2A_2 \mid bB_2 \mid aA_1 \mid b \quad - \textcircled{1}$

$A_1 \rightarrow aA_1B_2A_2 \mid bB_2A_2 \mid aA_1A_2 \mid bA_2$

$A_1 \rightarrow a \dots$  in GNF

Now for  $B_2$  productions

$B_2 \rightarrow A_2A_1B_2 \Rightarrow$  substitute  $\textcircled{1}$  for 1st  $A_2$

$B_2 \rightarrow aA_1B_2A_1B_2 \mid bB_2A_1B_2 \mid aA_1A_1B_2 \mid bA_1B_2$

$B_2 \rightarrow A_2 \neq A_1 \Rightarrow$  substitute  $\textcircled{1}$  for 1st  $A_2$

$B_2 \rightarrow aA_1B_2A_1 \mid bB_2A_1 \mid aA_1B_1 \mid bA_1$ .

The final set of productions is.

$A \rightarrow aA_1B_2 \mid bB_2 \mid aA_1 \mid b$ .

$A_1 \rightarrow aA_1B_2A_2 \mid bB_2A_2 \mid aA_1A_2 \mid bA_2 \mid a$ .

$B_2 \rightarrow aA_1B_2A_1B_2 \mid bB_2A_1B_2 \mid aA_1A_1B_2 \mid bA_1B_2 \mid$   
 $aA_1B_2A_1 \mid bB_2A_1 \mid aA_1A_1 \mid bA_1$

$V = (A_1, A_2, B_2)$

$T = (a, b)$

$S = A_1$

PUMPING LEMMA OF CFG  
Break down a CFG into  $uvxyz$ .  $|w| \geq n$   
such that  $v, y \neq \epsilon$ ,  $|vxy| \leq n$   
now pump up  $v, v \notin y$ .  $i > 0$

$uv^ixy^iz$

check if pumped string is a CFG or not

$v \notin y$  can be pumped multiple times.

$\frac{0000011111}{u \quad v \quad x \quad y \quad z}$ .  $|w| = 10$   $n = 10$  CFL: ~~On<sup>n</sup>1<sup>n</sup>~~

$$\therefore |vxy| = 6 \leq n.$$

$$u = 00$$

$$v = 00$$

$$x = 01$$

$$y = 11$$

$$z = 11$$

$\therefore$  pumping  $v \notin y$ .

$$i = 2$$

$uv^2x y^2 z$ .

00 0000 01 1111 11

new string: 000000011111. is in CFL.

$\therefore 0000011111$  is in CFL.

## CHOMSKY HIERARCHY

The chomsky hierarchy represents the class of language that are accepted by different machine. The category of languages in Chomsky Hierarchy is as given below

language class	language	grammar	Machine	example
Type 3	Regular	Regular grammar	FSM i.e NFA to DFA	$a^*b^*$
Type 2	Context free	Context free	PDA	$a^n b^n$
Type 1	Decidable language	Context sensitive	Linear bounded automata	$a^n b^n c^n$
Type 0	Computable language	Unrestricted grammar	Turing machine	

## Applications of CFG:

- 1) parsing technique
- 2) Markup Language
- 3) XML and Document type definitions.