

CONTEXT FREE GRAMMARS

* Grammar :

- Grammar is a set of formal rules which check the correctness of the sentences or we can construct the correct sentences using these rules.
- These sentences are known as grammatical sentences and the rules as syntax. The meaning is semantic.
- Machine is a formal system which recognizes as well as generates the class of language.
- These are formal languages and we want to have a relation between the language and the machine.
- Two aspects of language
 - i) Generative Capacity
 - ii) Grammatical Constituents
- Generative Capacity indicates that the grammar should generate all and only the sentences belonging to the language whereas Grammatical constituents gives the stress on primitive structures which compose the sentence.
- The first will consider the consistency and completeness of the formal system whereas second deals with synthesis and analysis.

Grammar Defⁿ:

- Grammar is used to specify syntax for language. The grammar will generate the strings. It is a quadruple and is defined as:

$$G = (V, T, P, S)$$

Where, $V \rightleftharpoons$ A finite set of Variables / Non-terminals

T = Finite set of terminals

P = Set of production / rewriting rules

S = start variable.

Variables → represented by capital letters
Terminals → represented by small letters/operators

e.g: $G = (\{E\}, \{+, *\}, \text{id}, P, E)$

all symbols \downarrow either found in Σ or Δ \hookrightarrow start symbol

Variables Terminals Production

$P : E \rightarrow E+E$ (using grammar of term. & non-term.)
 OR $E \rightarrow E+E / E * E / \text{id}$

OR $E \rightarrow \text{id}$ (in Σ exists in column 2nd row)

now we can read down rows based on the variable.

Derivation of a sentence (word)

into all strings \hookrightarrow (strings of terminals) are said

which are called productions. with associated priorities

e.g: "id + id", derive it. read row Σ as well as Δ

$E \Rightarrow E+E$ (using $E \rightarrow E+E$)

$E+E \Rightarrow \text{id} + E$ (using $E \rightarrow \text{id}$)

$\text{id} + E \Rightarrow \text{id} + \text{id}$ (using $E \rightarrow \text{id}$)

so terminals \hookrightarrow variables appropriate

variables are read down columns until no more

Context Free Grammar (CFG):

→ produces like $\text{cat} \in \Sigma^*$ (language)

- Context free Grammar generates context free language (CFL) that are recognized by push down automata (PDA)

- A grammar is said to be context free grammar if the productions are of the form, as per definition

$A \rightarrow \alpha$ (where $A \in \Delta$ & $\alpha \in \Sigma^*$)

i.e. (one non-terminal) = finite string of terminals and non-terminals (sentential form)

where, α consists to be still a full word

$A \rightarrow$ non-terminal (i.e. $A \in \Delta$)

$\alpha \rightarrow$ some sentential form (Combination of
 variable & terminals)

Here, the left hand side of the production can contain only one non-terminal and right hand side of the production can contain combination of non-terminal and terminals. The start variable can be present on the right hand side of the production.

$$\text{e.g.: } G = (\{S, A, B\}, \{a, b\}, P, S)$$

where P:

$$S \rightarrow aA \mid bB \mid a$$

$$A \rightarrow aA \mid a$$

$$B \rightarrow bB \mid b$$

* Variables:

symbols that are taking part in derivation process but are not part of the derived string

* Terminals:

symbols that are present in the derived sentences.

* Derivation:

The sequence of production used to derive the sentence is known as derivation.

* Derivation of Grammar: page "problem arise"

- A string which is generated from the given product rules and represented in the form of a tree, such tree is known as derivation tree or parse tree or rule tree.
- The derivation tree will have labels to all the nodes. These labels could be non-terminal symbols or terminal symbols or ϵ .
- The root is always labelled start.
- Nodes having labels as non-terminal are the interior nodes.
- All the leaf nodes will be terminal symbols.
- If a particular non terminal is generating a string of terminals and non-terminals of length n then that will have n children. These children will be attached from left to right in sequence.

(1) Left Most Derivation (LMD)

If at each step in the derivation process, we replace the leftmost variable by its production then the derivation is said to be leftmost derivation.

(2) Right most Derivation (RMD)

If at each step in the derivation process, we replace the rightmost variable by its production then the derivation is said to be rightmost derivation.

X. I. E.

Mahim, Mumbai

Page No. :

Date :

* Generation of CFG:

τ	$L(\tau)$	CFG
a	$\{a\} \cup \{a\}^* = \{a\}$	$S \rightarrow a$
b	$\{b\}$	$S \rightarrow b$
$a b$	$\{a\} \cup \{b\} = \{a, b\}$	$S \rightarrow a \mid b$
a^*	$\{a, aa, aaaa, \dots\}$	$S \rightarrow aS \mid \epsilon$
a^+	$\{a, aa, aaaa, \dots\}$	$S \rightarrow aS \mid a$
$a \cdot b$	$\{ab\}$	$S \mid S \rightarrow ab$
$(ab)^*$	$\{\epsilon, ab, abab, \dots\}$	$S \rightarrow abS \mid \epsilon$
$(afb)^*$	$\{\epsilon, a, b, aa, ab, ba, \dots\}$	$S \rightarrow as \mid bs \mid \epsilon$

Q. Write CFG in form: bab starts with a and ends with b .

$$i) \tau = (a|b)^*$$

$$ii) \tau = ba^*$$

$$iii) \tau = b^* a^*$$

$$iv) \tau = (baa + abb)^*$$

→

$$i) \tau = (a|b)^*$$

$$S \rightarrow as \mid bs \mid \epsilon$$

$$\therefore G = (\{S\}, \{a\}, \{S \rightarrow as \mid bs \mid \epsilon\}, S)$$

$$ii) \tau = ba^*$$

$$S \rightarrow bA$$

$$A \rightarrow aA \mid \epsilon$$

$$\therefore G = (\{S, A\}, \{a, b\}, \{S \rightarrow bA, A \rightarrow aA \mid \epsilon\}, S)$$

$$iii) \tau = b^* a^*$$

$$S \rightarrow BA$$

$$B \rightarrow bB \mid G$$

$$A \rightarrow aA \mid \epsilon$$

$$\therefore G = (\{S, A, B\}, \{a, b\}, \{S \rightarrow BA, B \rightarrow bB \mid \epsilon, A \rightarrow aA \mid \epsilon\}, S)$$

* Ambiguous or Context Free Language : A grammar is said to be Ambiguous if the same string can be generated in two different ways (using more than one LMD/RMD)

→ To determine that the given grammar is ambiguous we need to derive at least one sentence using more than one LMD or RMD.

a. check whether following grammar is ambiguous or not:

$E \rightarrow E+E \mid E * E \mid id$

→ LMD1: Deriving LMD for sentence id + id * id

$E \Rightarrow E+E$ (using $E \rightarrow E+E$)

$\Rightarrow id + E$ (using $E \rightarrow id$)

$\Rightarrow id + E * E$ (using $E \rightarrow E * E$)

$\Rightarrow id + id * E$ (using $E \rightarrow id$)

$\Rightarrow id + id * id$ (using $E \rightarrow id$)

LMD2:

$E \Rightarrow E * E$ (using $E \rightarrow E * E$)

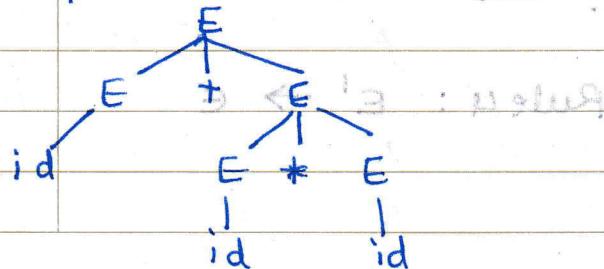
$\Rightarrow E + E * E$ (using $E \rightarrow E+E$)

$\Rightarrow id + E * E$ (using $E \rightarrow id$)

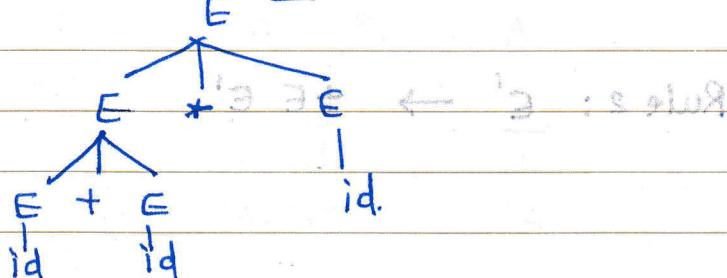
$\Rightarrow id + id * E$ (using $E \rightarrow id$)

$\Rightarrow id + id * id$ (using $E \rightarrow id$)

parse tree for LMD1



LMD2



As there is no unique parse tree for generating a string, the grammar is said to be ambiguous.

- * Unambiguous grammar: A grammar which generates exactly one string for every input string.
- Grammar is unambiguous if it can derive all sentences using exactly one LR(0) item.
- * Ambiguity Resolution: To resolve the ambiguity, no method exists.
- Takema heuristic approach of writing words from right to left.

Rule 1:

Let $A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_n | B_1 | B_2 | B_3 | \dots | B_n$ such that $b_1 + b_2 + \dots + b_n = 3$.

Then, $A \rightarrow B_1 A' | B_2 A' | \dots | B_n A'$ such that $b_1 + b_2 + \dots + b_n = 3$.
 $A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_n A' | \epsilon$ such that $b_1 + b_2 + \dots + b_n = 3$.

Q. Eliminate the ambiguity in the given grammar

$E \rightarrow E + E | E * E | id$

→

$E \rightarrow E + E | E * E | id$
 $\downarrow b_1 + b_2 \leftarrow 3$ $\downarrow b_1 + b_2 \leftarrow 3$ $\downarrow b_1 + b_2 \leftarrow 3$
 $A \rightarrow A + A | A * A | id$ $\alpha_1 \rightarrow \alpha_1 + \alpha_1 | \alpha_1 * \alpha_1 | id$ $\beta_1 \rightarrow \beta_1 + \beta_1 | \beta_1 * \beta_1 | id$

Rule 1: $E \rightarrow id$ $E' \rightarrow E$

\downarrow \downarrow \downarrow

A B_1 A'

Rule 2: $E' \rightarrow + E E'$

\downarrow \downarrow \downarrow

A' α_1 A'

Rule 3: $E' \rightarrow * E E'$

\downarrow \downarrow \downarrow

$A' \rightarrow \alpha_1$ $A' \rightarrow \alpha_2$ $A' \rightarrow \alpha_3$

Rule 4: $E' \rightarrow E$

\downarrow \downarrow \downarrow

$A' \rightarrow E$ $A' \rightarrow E$ $A' \rightarrow E$

* Simplification of a CFG being $\text{2291920} \leftarrow 2$

$AA\alpha \mid \beta \mid d\alpha \leftarrow 2$

* Elimination of Useless Symbols: $\text{2291920} \leftarrow 2$

$\alpha \beta \mid \delta \alpha \leftarrow A$

• Useless Symbol:

A symbol ' x ' is useful if $\Rightarrow \beta$, where

$S \xrightarrow{*} \alpha X \beta \xrightarrow{*} \gamma \delta \alpha \leftarrow A$

where,

$S \rightarrow \text{start symbol} \leftarrow A$

$\alpha, \beta \rightarrow \text{sentential forms}$

for e.g. for $w \Rightarrow \text{sentence}(e, 'q', 'T', 'v') \leftarrow A$

(\Rightarrow) α otherwise ' x ' is useless

$AA\alpha \mid \beta \mid d\alpha \leftarrow 2 \leftarrow q$

- A production where useless symbols are present is known as useless production.

• Elimination procedure:

Given CFG $G = (V, T, P, S) \leftarrow 2 \leftarrow q$

Define $G' = (V', T', P', S)$ be the CFG with no useless production (symbol) such that

$L(G') = L(G)$

Step 1: Initialize P' to P

Step 2: find useless symbols

* Symbol is useless if it is not reachable from start symbol.

* Symbol is useless if it cannot derive any sentence (\Leftrightarrow any combination of terminals)

Step 3: Eliminate all occurrence of useless symbols

* Elimination of Unit Production

- Unit production :

A production of the form $A \rightarrow B$ where A & B are variables is called as Unit production.

Note:

Chain of unit production is :-

$$A \rightarrow X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_n \rightarrow B$$

- Elimination procedure :

Given $G = (V, T, P, S)$

Define $G' = (V', T', P', S)$ be the CFG with no unit productions such that $L(G') = L(G)$.

Step 1 : Initialize P' to P .

Step 2 : Find unit productions.

Say $A \rightarrow B$ is a unit production and $B \rightarrow \beta_1 | \beta_2 | \dots | \beta_n$ are non-unit productions of B, then add non-unit productions of B to A (if not present)
i.e. $A \rightarrow \beta_1 | \beta_2 | \dots | \beta_n$.

Step 3 : Delete all unit productions.

- Non - Unit productions :

A production which is not of the form $A \rightarrow B$ where A & B are variables is called as non-unit productions.

* Elimination of Null production

* NULL production:

A production of the form
 $x \rightarrow \epsilon$ or $x \rightarrow \lambda$, where x is a variable
 is called NULL production

* Nullable Variable:

A variable ' x ' is said to be nullable if

$\exists x \rightarrow \epsilon \text{ or } x \xrightarrow{*} \epsilon$

* Elimination Procedure:

Given CFG, $G = (V, T, P, S)$

Define $G' = (V', T, P', S)$ be the CFG with no
 null productions such that $L(G') = L(G) - \{\epsilon\}$

Step 1: Initialize P' to P .

Step 2: Find the nullable variables $\leftarrow A$

Say $A \rightarrow \alpha$ is a production where α
 contains some nullable variables, then add
 to A , the productions obtained by deleting
 all possible subsets of nullable variables
 from α (if not present)

Step 3: Delete null productions.

NORMAL FORMS

→ The normal forms are used to rewrite the CFG as per the specified conditions. When certain condition is imposed on the RHS, the grammar is said to be in the Normal Form.

* Chomsky Normal Form (CNF) - older work

- Any ϵ -free CFL can be generated by a CFG whose productions are of the form

$$A \rightarrow BC \quad \text{or} \quad A \rightarrow a$$

where, A, B & C are variables and a is a terminal

$$A \leftarrow a, \quad AA \leftarrow a, \quad AB \leftarrow a$$

$$ABA \leftarrow a$$

- In the Chomsky's Normal Form there will be restriction on the RHS.

* Conversion procedure from CFG to CNF:

- 1) Simplify the given CFG by eliminating Null, Unit and Useless productions.
- 2) Add to the solution those productions that are already in CNF.
- 3) For the non-CNF productions:
 - a) Replace the terminals by some variables.
 - b) Limit the no. of variables on RHS.

(प्र० १) एक सिर्फ टाइप नहीं हो सकता।

$$AD \leftarrow a, \quad d \leftarrow a$$

$$\begin{aligned} d &\leftarrow a \\ ad &\leftarrow a \\ b &\leftarrow a \\ AD &\leftarrow a \\ A &\leftarrow a \\ a &\leftarrow a \end{aligned}$$

$$a \leftarrow a$$

X. I. E.

Mahim, Mumbai

Aa | as | a Page No. 2

Date: _____

18d ← 8

* Greibach Normal Form (GNF):

- *Harnef baifl qmfa ni ei commopj nqfj sft*
- Any ϵ -free CFL can be generated by a CFG whose productions are of the form

$$A \rightarrow a\tau$$

where,

$A \rightarrow$ a variable

$a \rightarrow$ a terminal

$\tau \rightarrow$ string of τ variables which can be empty

$$Bd | B \leftarrow d$$

$$Aa \leftarrow a$$

$$Aad | ad \leftarrow a$$

$$Ad | d \leftarrow A$$

$$aa \leftarrow a$$

$$Bd | Bd \leftarrow a$$

$$ABA \leftarrow a$$

* Conversion procedure for CFG to GNF:

▷ Simplify the given CFG by eliminating Null, unit and useless productions.

2) Use any combination of Rule1 and Rule2 to get the CFG in GNF.

Rule1: Let $A \rightarrow B\alpha$ be some A-productions and $B \rightarrow \beta_1 | \beta_2 | \dots | \beta_n$ be B-productions

then we can write A-productions as

$$A \rightarrow \beta_1\alpha | \beta_2\alpha | \dots | \beta_n\alpha | A\beta_1 | A\beta_2 | \dots | A\beta_n \leftarrow a$$

Rule2: Let $A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_r$ be A-productions and $A \rightarrow \beta_1 | \beta_2 | \dots | \beta_s$ be the remaining A-productions then introduce a new variable B, then

B-productions : $B \rightarrow \alpha_i | \alpha_i B$, $1 \leq i \leq r$

A-productions : $A \rightarrow \beta_i | \beta_i B$, $1 \leq i \leq s$

Note: Let $A \rightarrow \alpha E F \# a$, then we can rewrite A as:

$A \rightarrow \alpha E F X_1 X_2$ ($X_1, X_2 \rightarrow$ auxiliary variables),

$$X_1 = b, X_2 = a$$

CHOMSKY'S HIERARCHY

- Regular expression language, Context free language are few example of languages.
- These languages along with some new classes of languages were classified by Noam Chomsky.
- According to Noam Chomsky, there are four types of grammars as follows:

Type 0 → Unrestricted Grammar / Semi-Thue Grammar / phrase structure grammar

Type 1 → Context Sensitive Grammar

Type 2 → Context Free Grammar

Type 3 → Regular Grammar

- 1) Type 0 → Unrestricted Grammar
- The grammar is also known as Phrase Structure grammar / Semi-Thue Grammar.
- The largest class defined in CHOMSKY's HIERARCHY is TYPE 0 grammar.
- The productions in type 0 grammars are of the form $\alpha \rightarrow \beta$ where α & β are strings formed by the symbols in the grammar and also α is not NULL.
- There are no restrictions on the production rules in this case.
- It generates Recursively Enumerable Language (REL).
- The language generated (strings) will be recognized by the Turing Machine. Hence a TM can be designed such that its algorithmic way can check the acceptance.

of any strings belonging to any 'Type 0' grammar.

e.g.: Let, $V = \{A, B, S\}$

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

$$G = (V, T, P, S)$$

$$P: S \rightarrow SA$$

$$SA \rightarrow AB$$

$$A \rightarrow a$$

2) Type 1: Context Sensitive Grammars (CSG)

→ It is a grammar in which the productions are of the form $\alpha \rightarrow \beta$ such that the length of the string β , will be at least as long as the length of α .

→ It means, the left hand side of each production is never longer than the right hand side.
 $|\beta| \leq |\alpha|$

→ It generates context sensitive language (CSL).

→ The generated languages are recognized by linear bounded automata (LBA).

→ The productions of the form $\alpha \rightarrow \beta$ where, $\alpha \in \text{non-terminal}$ and $\beta \in \text{string of terminals } \{a, b, \dots\}^*$ ($a, b, \dots \in T$)

→ The replacement of α by β is possible only if it is preceded by α and is followed by β . Hence the context in which α is appearing will be important for replacement. If the context is not appropriate the rule will not be applied.

→ e.g. Let $G = (V, T, P, S)$ where, $V = \{S, A, B\}$, $T = \{a, b\}$, P, S (where, $P: S \rightarrow aAb$, $A \rightarrow aB$, $A \rightarrow aA$, $A \rightarrow b$, $A \rightarrow a$)

X.I.E.

Mahim, Mumbai

* Context Type 2 \rightarrow Context Free Grammar (CFG):

- CFG generates context free languages (CFL).
- The languages generated are recognized by push down automata (PDA).
- A grammar is said to be CFG, if the productions are of the form, i.e. (one non-terminal) \rightarrow finite string of sentential form.

$$[A \rightarrow \alpha]$$

A \rightarrow non-terminal

d \leftarrow 9

α \rightarrow some sentential form (V + T)

- Here, the LHS. of the production can contain only one non-terminal and RHS. of the production can contain combination of non-terminal and terminal. The start variable can be present on the RHS of the production.

→ e.g.:

$$G = (S, A, B, \{a, b\}, P, S)$$

where P:

$$S \rightarrow aA \mid bB \mid a$$

$$A \rightarrow aa \mid ab$$

$$B \rightarrow bb \mid ba$$

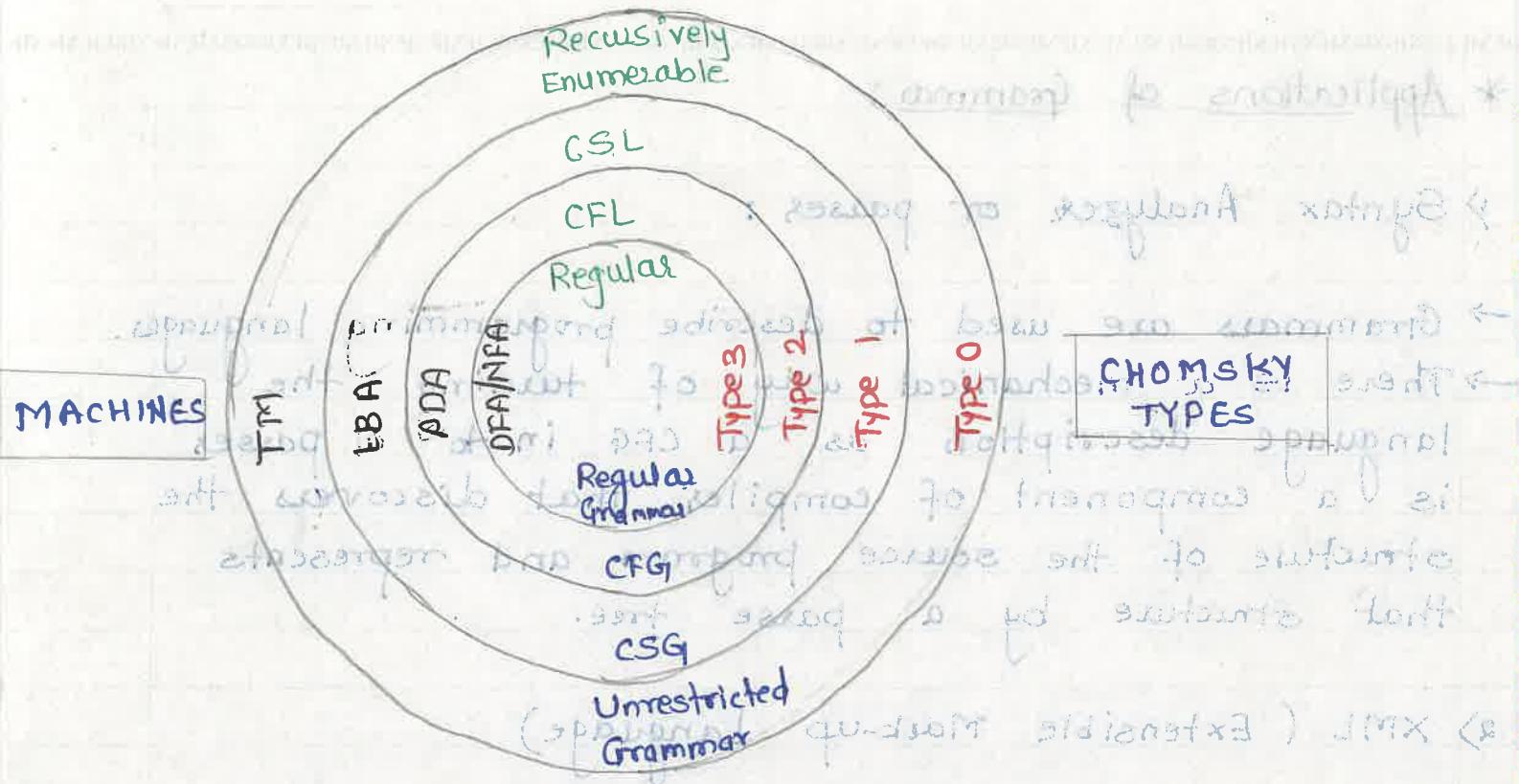
* Type 3 \rightarrow Regular Grammar (RG)

- If the production rules are of the form in which the RHS contains at the most one non-terminal symbol, the grammar is known as regular grammar.
- Regular grammar generates regular language (RL).

- The generated languages are recognized by FSM.
- These regular languages can also be expressed by simpler expressions called as Regular Expression (RE).
- The productions are of the form, ~~containing at least two non-terminals~~ if there is only one non-terminal
 $S \rightarrow ST^* / T^*$ or $S \rightarrow T^* S / T^*$
 where, S is non-terminal, T is non-terminal and T^* is any combination of terminals.
- The grammar is said to be left linear if the productions contain many one non-terminal symbol which occurs at the leftmost.
 e.g.: $S \rightarrow a$ same strict \leftarrow (leftmost-rightmost) i.e.
 $S \rightarrow Pb$
 $P \rightarrow b$ leftmost-rightmost $\leftarrow A$
- The grammar is said to be right linear if the productions contain many one non-terminal symbol which occurs at the rightmost.
 e.g.: $S \rightarrow ab$ no non-terminal after b so it is right linear $S \rightarrow aPb$ since a is leftmost $P \rightarrow b$ no non-terminal after b
- The grammar is said to be linear if any non-terminal generates a single non-terminal on the RHS.
- The grammar are said to be equivalent if the strings generated by them are identical.

LANGUAGES

source: Theory
Computat" - Kar



GRAMMARS

Fig: Chomsky Hierarchy of formal languages

Sr. No.	Grammar Type	Grammar accepted	Language accepted	Automation
1.	Type 0	Unrestricted grammar	Recursively Enumerable Language (REL)	Turing Machine (TM)
2.	Type 1	Context Sensitive Grammar (CSG)	Context Sensitive Language (CSL)	Linear Bounded Automata (LBA)
3.	Type 2	Context free grammar (CFG)	Context Free Language (CFL)	Pushdown Automata (PDA)
4.	Type 3	Regular grammar	Regular Language	Finite Automaton (FA)

* Applications of Grammars:

1) Syntax Analyzers or parsers:

- Grammars are used to describe programming languages.
- There is a mechanical way of turning the language description as a CFG into a parser. A parser is a component of compiler that discovers the structure of the source program and represents that structure by a parse tree.

2) XML (Extensible Markup Language)

- An essential part of XML is DTD (Document Type Definition) which is essentially a CFG that describes the allowable tags that and the ways in which these tags may be nested.