

## NLP

### Grammar:

It is generally represented as  $(V, T, P, S)$ , where

$V$  = vertices,  $T$  = terminals,  $P$  = production,  $S$  = start symbol.

### Example:

$$\begin{array}{l} S \rightarrow aSB \\ S \rightarrow aB \\ B \rightarrow b \end{array}$$

} grammar,

then ;

$V = \{S, B\}$

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

$$P = \{S \rightarrow aSB, S \rightarrow aB, B \rightarrow b\}$$

$$S = \{S\} \rightarrow \text{start symbol}$$

### Derivation:-

Getting the string from the grammar is

### "Derivation":

Ex:- deriving "aabb"

$$S \Rightarrow aSB$$

left most derivation

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

$$\Rightarrow aAbB$$

$$\Rightarrow aabb$$

In "asB", if you replace 's'

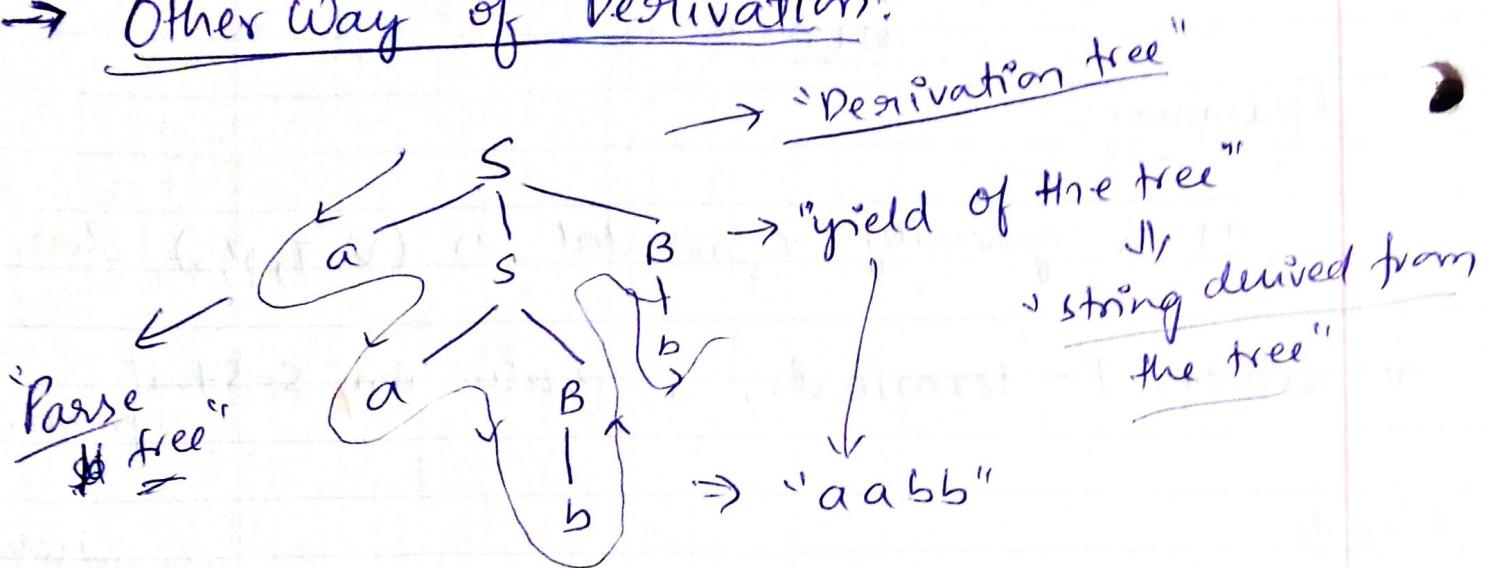
• first as part of derivation, then it is "left most derivation".

else, if you replace "B", first, then it is "right most derivation"

Intermediate forms are called

• sentential forms (or) sequential forms

## → Other Way of Derivation:



## ① How to Construct Grammars?

② Given a Grammar, how to identify that Grammar is representing?

Problem:-

Given, language  $L = \{aa, ab, ba, bb\}$ . Construct Grammar, which accepts set of all strings of length '2'.

$$S \rightarrow aa | ab | ba | bb$$

- Steps: First construct the building blocks, then extend it

Ex:  $(a+b)(a+b)$  is building blk of  $(aa|ab|ba|bb)$

$$\begin{array}{c} A \\ \downarrow \\ a \end{array} \quad \begin{array}{c} A \\ \downarrow \\ b \end{array}$$

$$\Rightarrow S \rightarrow AA$$

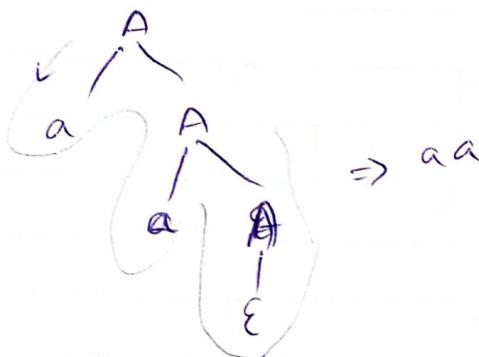
$$A \rightarrow a \mid b$$

Ex 2:

$$a^n \mid n \geq 0$$

a, aa, aaa, aaaa, ...

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



Ex 3:

$$(a+b)^*$$

$$S \rightarrow aS \mid bS \mid \epsilon$$

"abab"



Ex 4: Set of all strings of length, atleast, 2.

①

$$(a+b)(a+b)(a+b)^*$$

$$S \rightarrow AAB$$

$$A \rightarrow a/b$$

$$B \rightarrow aB \mid bB \mid \epsilon$$

②

"at most 2"

$$(a+b+\epsilon)(a+b+\epsilon)$$

A            A

$$S \rightarrow AA$$

$$A \rightarrow a \mid b \mid \epsilon$$

③ starting & ending with same symbol.

$$S \rightarrow aAa \mid bAb \mid a \mid b \mid \epsilon$$

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

③

"starting with 'a'"  
and ending with "b"

$$S \rightarrow AAb$$

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

④

starting and ending with diff symbols

$$a(a+b)^*b + b(a+b)^*a$$

$$S \rightarrow aAb \mid bAa$$

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

⑥ Construct grammar such that,

$$a^n b^n / n \geq 1$$
$$\Rightarrow S \rightarrow a S b / a b$$

⑦ Generate set of palindromes.

$$\rightarrow w w^R \cup w a w^R \cup w b w^R$$

Sol:-  $S \rightarrow a S a / b S b / a / b / \epsilon$

## CLASSIFICATION OF GRAMMARS

• According to Chomsky,

- Type 3
- Type 2
- Type 1
- Type 0.

Type 3:- Regular Grammar / Finite Automata machine

If all the productions are of the form,

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

A, B  $\in V$   
 $\alpha, \beta \in T^*$

→ right linear grammar  
as 'B' is towards right side

$$A \rightarrow B\alpha | \beta$$

$A, B \in V$

$\alpha, \beta \in T^*$

left-linear grammar (LLG)

Ex:-

$$\begin{array}{l} A \rightarrow a\underline{B} | a \\ B \rightarrow a\underline{B} | b\underline{B} | a | b \end{array} \quad \left. \begin{array}{c} \\ \end{array} \right\} RLG$$

$$\begin{array}{l} A \rightarrow Ba | a \\ B \rightarrow Ba | Bb | a | b \end{array} \quad \left. \begin{array}{c} \\ \end{array} \right\} LLG$$

Ex:-

$$A \rightarrow Ba | a \rightarrow LLG$$

$$B \rightarrow aB | a \rightarrow RLG$$

X Type 3. It cannot be combination of LLG, RLG.

Type 2 :- /context-Free-Grammar.

If the If all the productions are of the form,

$$A \rightarrow \alpha, \text{ where } A \in V$$

$$\alpha \in (V \cup T)^*$$

Ex:-

$$A \rightarrow aAb | ab$$

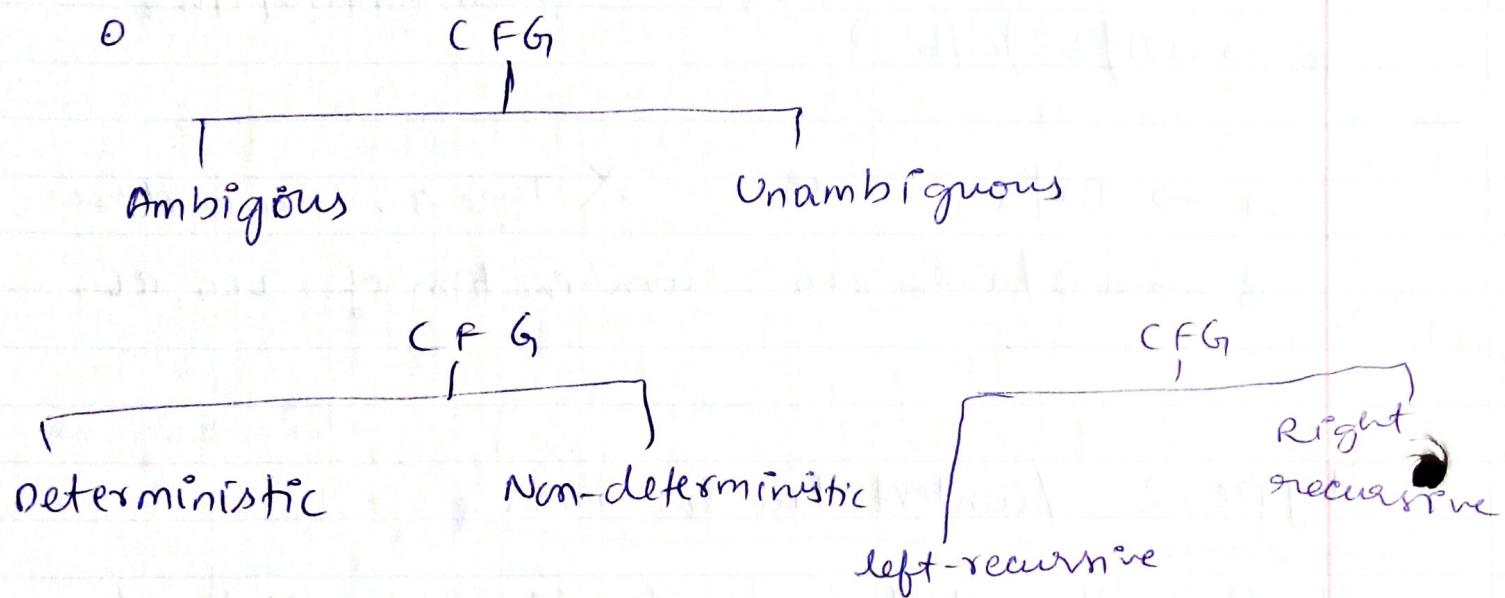
$$\begin{array}{l} A \Rightarrow aAb \\ \Rightarrow a\underline{ab}b \end{array}$$

We call it context-free, because we can replace 'A' without looking at the context.

→ But in context-sensitive Grammar (Type-i), we will have  $\textcircled{A} A \rightarrow ab | aAb$   
we have to look, whether before 'A', 'a' is present (or) not.

- Language generated is "Context-free language"
- machine is "Push Down Automata" (PDA)

### CFG - Dichotomy :-



③ If we want to find whether a string,  $w$ , exists in the language generated by CFG. i.e.,  $w \in L(G)$  (or not.)

sol:- we have to increase the length of the string,  $w$  too by using the grammar to  $|w|$ .

but, if we have  $A \rightarrow \epsilon \rightarrow$  epsilon productions  
 $\quad \quad \quad A \xrightarrow{\text{or}} B \rightarrow$  unit productions,

then again the length of the string reduces

and we may not be able to check it correctly,  
 So, we have to make sure that our grammar production doesn't contain "epsilon"  $\epsilon$   
"Unit productions". then, we can increase the length by '1' at each time, and check whether string is present in grammar or not

- ① One such algorithm to find the membership of a string in grammar is "CYK algorithm"

### Elimination of $\epsilon$ -Productions:

- If the language contains ' $\epsilon$ ', we will not able to eliminate  $\epsilon$  to remove ' $\epsilon$ ' productions . . .
- To eliminate  $\epsilon$ -productions, you should remove the "nullable" variables:
- A variable is nullable if it generates ' $\epsilon$ ', or immediately (or) after some sentential forms

$$\text{Ex:- } A \rightarrow \epsilon -$$

$s \rightarrow a \cancel{A} b \Rightarrow s \rightarrow a \cancel{(A)} b$   $\rightarrow$  so, 's' is <sup>not</sup> nullable as it doesn't generate 'null'

Ex:-

To eliminate ~~left~~ 'ε' productions in below ex

$$S \rightarrow aSb \mid aAb \cancel{Bab}$$

$$A \rightarrow \epsilon$$

Sol: ①

$A \rightarrow \epsilon$  is a nullable variable

② Next, write productions with & without nullable variable

$$S \rightarrow aSb \mid aAb \mid ab$$

$$A \rightarrow \epsilon$$

③ finally eliminate , nullable <sup>variable</sup> productions, and productions which contain nullable variable.

Ex:-

$$S \rightarrow AB ; A \rightarrow aAA/\epsilon ; B \rightarrow bBB/\epsilon$$

1. Find nullable variables.

$$A \rightarrow \epsilon ; B \rightarrow \epsilon \Rightarrow \{A, B, S\}$$

2. Write productions with ~~left~~

- when 'S' generates 'ε', then 'ε' belongs to the language, so, we'll not be able to remove <sup>all "null" productions</sup> 'ε' from the Grammar.

### 3. Write with and without nullable variables

$S \rightarrow AB / B / A / \epsilon$

$A \rightarrow AAA / AA / a$

$B \rightarrow BBB / BB / b$

all removed ' $\epsilon$ '  
 productions except  
 $S \rightarrow \epsilon$ , as it ' $\epsilon$ '  
 contains in the language

Ex:

$S \rightarrow AbaC \cancel{Bac}$

$A \rightarrow BC$

$\Rightarrow$  nullable variables = {B, C, A}

$B \rightarrow b / \epsilon$

$\Rightarrow S \rightarrow AbaC | bac | Aba | ba$

$C \rightarrow D / \epsilon$

$A \rightarrow BC | B | C$

$D \rightarrow d$

$B \rightarrow b$   
 $C \rightarrow D$   
 $D \rightarrow d$

BUT, removing ' $\epsilon$ ' productions, adding

UNIT - productions

• but they are also useless. For example,

$A \rightarrow B ; B \rightarrow C ; C \rightarrow D ; D \rightarrow d$

• To generate 'd',  $A \Rightarrow B \Rightarrow C \Rightarrow D \Rightarrow d.$

neither the length  
 nor terminals  
 increasing, therefore  
 we don't use unit

prod in CFG

## Elimination of Unit Productions:

$$S \rightarrow Aa|B$$

$$B \rightarrow A|bb$$

$$A \rightarrow a|bc|B$$

- we should make sure that even while eliminating unit productions, it should not affect the grammar of the language.
- So, you write all the productions, after removing unit productions. and then add productions for each variable, which will be formed by removing unit productions.

$$S \rightarrow Aa|bb|a|bc$$

$$B \rightarrow bb|a|bc$$

$$A \rightarrow a|bc|bb$$

If we remove,  $B \rightarrow$   
we miss

$$\{ S \rightarrow B \rightarrow bb$$

$$\} S \rightarrow B \rightarrow A \rightarrow a \rightarrow bc$$

$$\{ B \rightarrow A \rightarrow a \rightarrow bc$$

$$\} A \rightarrow B \rightarrow bb$$

Ex 2:

$$S \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow C/b$$

$$C \rightarrow D$$

$$D \rightarrow E$$

$$E \rightarrow a$$

Remove  
unit prod

$$S \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow b/a$$

$$C \rightarrow a$$

$$D \rightarrow a$$

$$E \rightarrow a$$

$$\rightarrow B \rightarrow C \rightarrow D \rightarrow b \rightarrow a$$

$$\rightarrow C \rightarrow D \rightarrow a$$

$$D \rightarrow b \rightarrow a$$



From 'S', you can reach 'a' & 'B', but not 'C',

'D', 'E'. So, these are useless symbols. You can delete those useless symbols.

### Useful Symbol:

• A symbol is said to be useful, if it is able to derive a part of the string (or) a terminal.

• It should be reachable from "start" symbol.

Steps:

Example:  
$$S \rightarrow AB/a$$
  
$$A \rightarrow BC/b$$
  
$$B \rightarrow aB/c$$

$$C \rightarrow aC/B$$

i) Add all terminals to the set of Useful Symbols, U

ii)

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

2. Add "symbols", which are generating ~~wrong~~  
symbols some terminals  
 $\Rightarrow U = \{a, b, S, A\}$
3. Add symbols, which  
 Add production symbols, whose right hand side of  
 production is made up of only useful symbols

$B \rightarrow aB | c$   $\times \Rightarrow$  'a' is there but it has 'B' with it, so  
 it cannot be included

- 'B' & 'c' are not useful as it'll not generate string.

Ex:  $B \rightarrow aB | c \Rightarrow B \rightarrow aB$   
 $c \rightarrow ac | B \Rightarrow aac$   
 $\Rightarrow aaaB$   
 $\Rightarrow aaa \cdot aB \Rightarrow$  we'll not be able  
 to get rid of 'B'  
 $\Rightarrow aaa \cdot a$  at any time  
 (or even 'c')

- Now delete 'B' & 'c' productions, and in  
 the productions, where 'B' & 'c' are present.

$\Rightarrow$   
So,  $S \rightarrow aB | a$   
 $S \rightarrow a$   $\times$   
 $A \rightarrow b$   $\emptyset$

- Now see if every symbol is reachable from ' $S$ ' or not.

Final Answer

$$S \rightarrow a$$

## Chomsky Normal Form:

A context-free grammar,  $G$  is said to be in CNF, if all its production rules are of the form

$$A \rightarrow BC$$

$$A \rightarrow a$$

$$A \rightarrow \epsilon$$

where  $A, B, C \rightarrow$  Non-terminal symbols

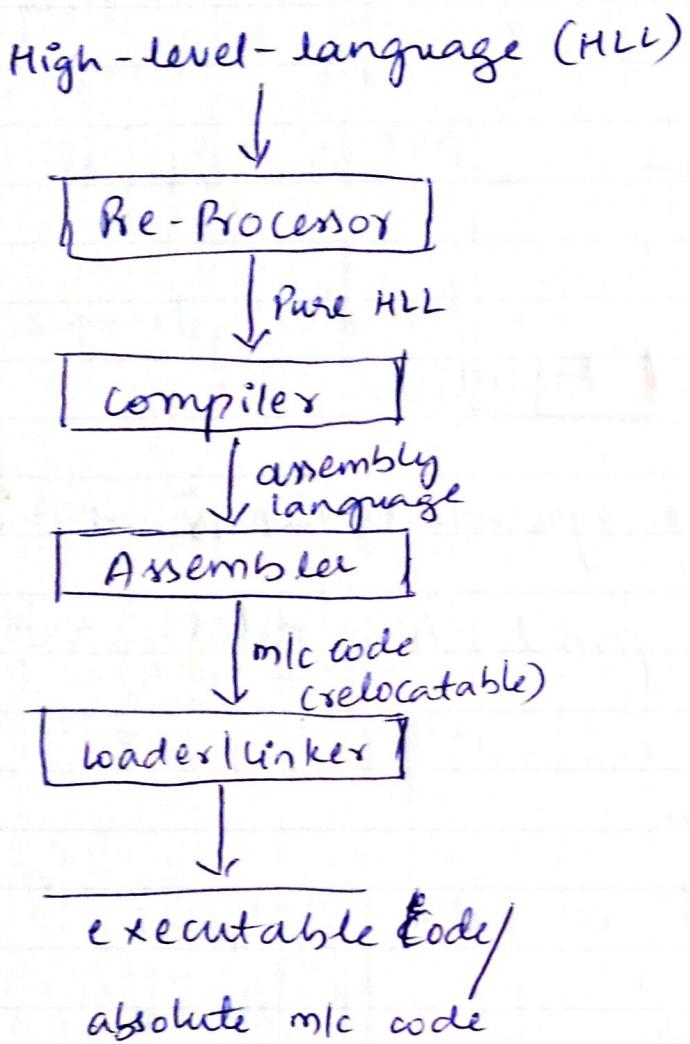
$a \rightarrow$  terminal symbol

$S \rightarrow$  Start symbol

$\epsilon \rightarrow$  empty string

It can appear only if ' $\epsilon$ ' is in  $L(G)$ .

## Compiler Design :-



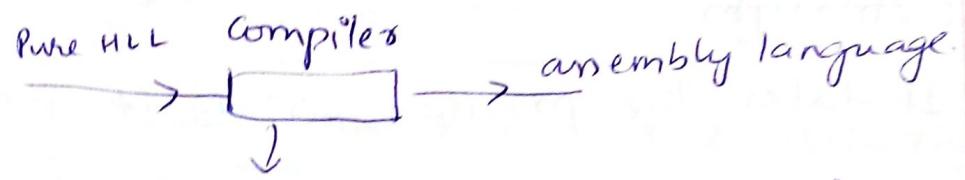
HLL → Pre-Processor: - generates

Pure HLL :-

- Rep example in 'c' pgm, it replaces all  
file inclusion # include < > → <sup>includes</sup> original file, for ex → conio.h
- # define → replaces ~~with~~ the defined constant in  
all locations or references in code.  
macro  
expansion.

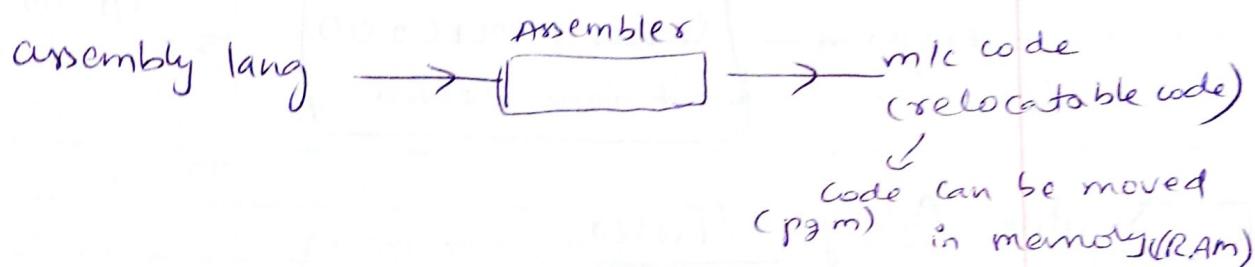
HLL → Pure HLL

②



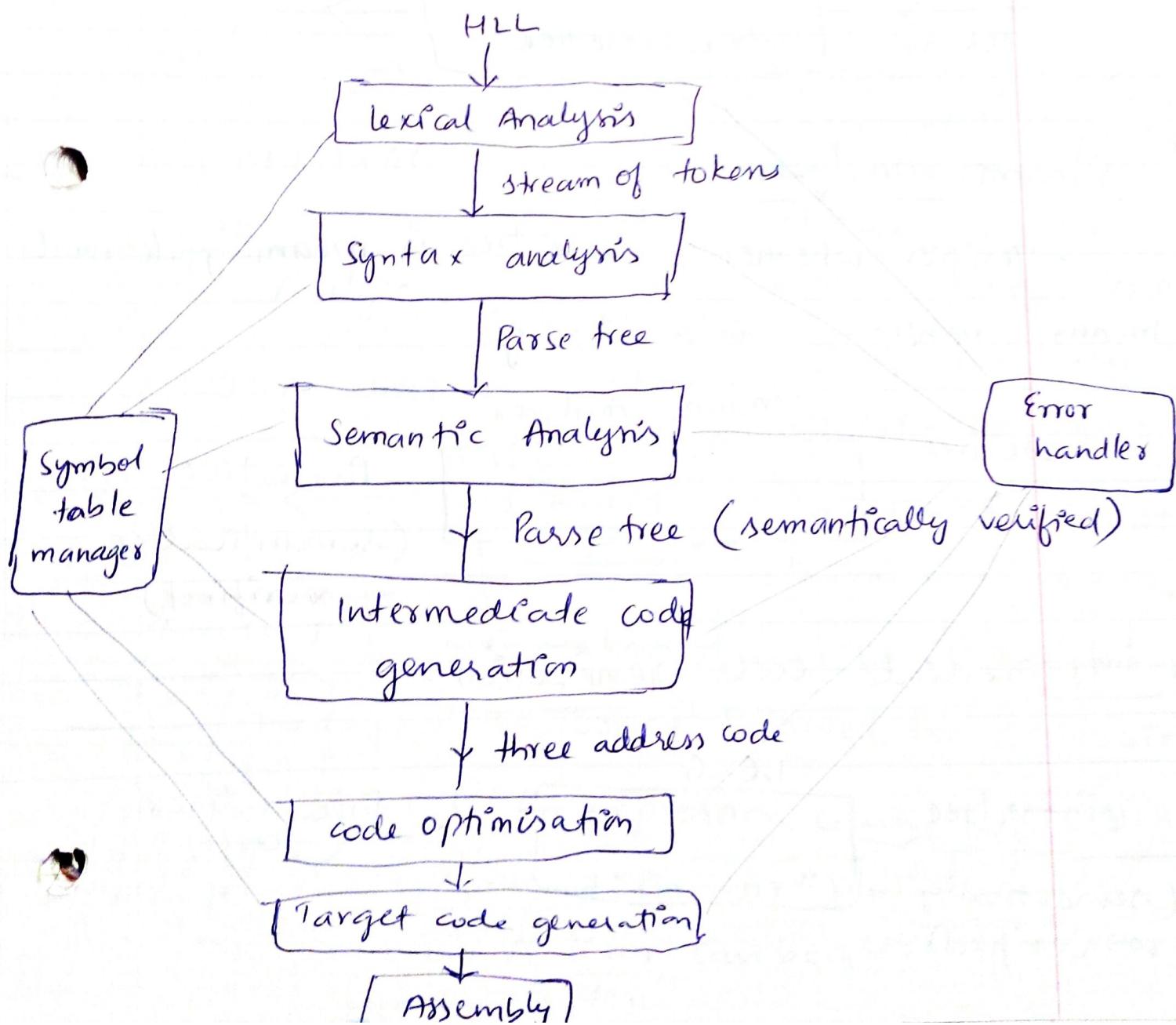
generates platform (Linux, Unix) specific code

③



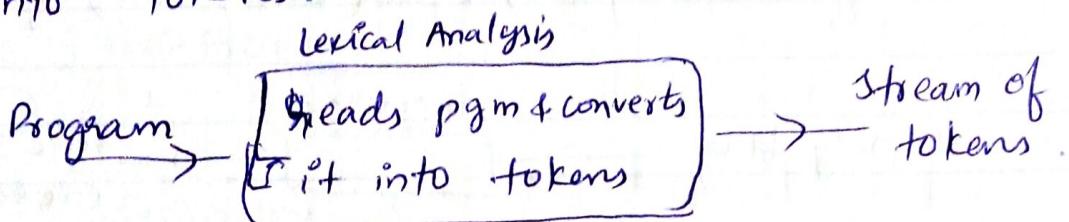
Code can be moved  
(p2m) in memory(RAM)

## Compiler Phases:

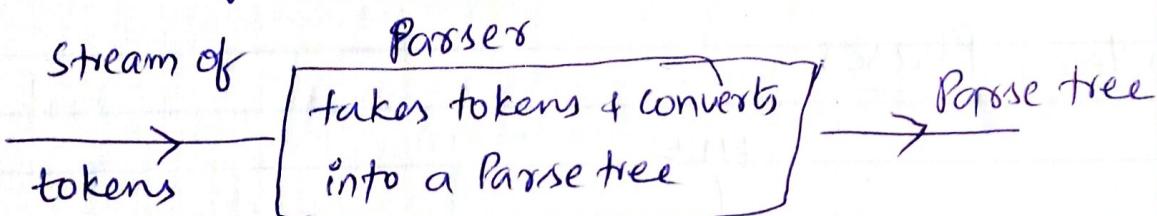


## ① Lexical Analysis:

- It takes the program, reads it and converts it into tokens.

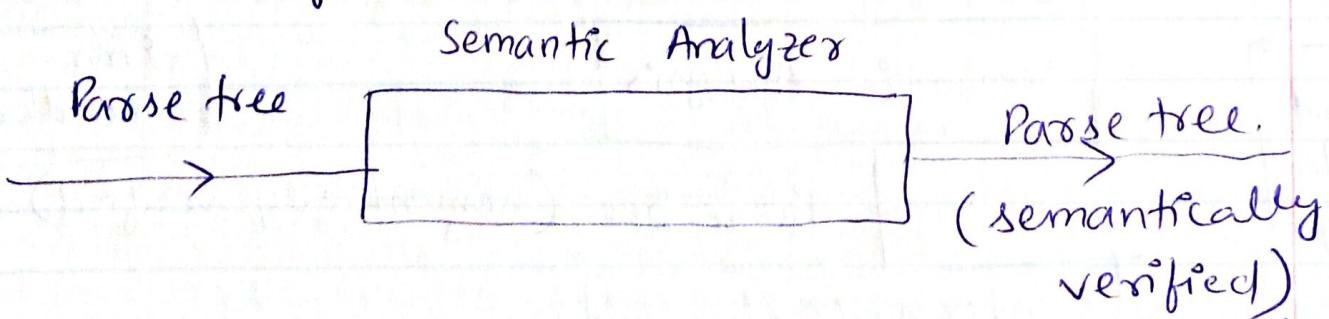


## ② Syntax Analyzer/Parser:

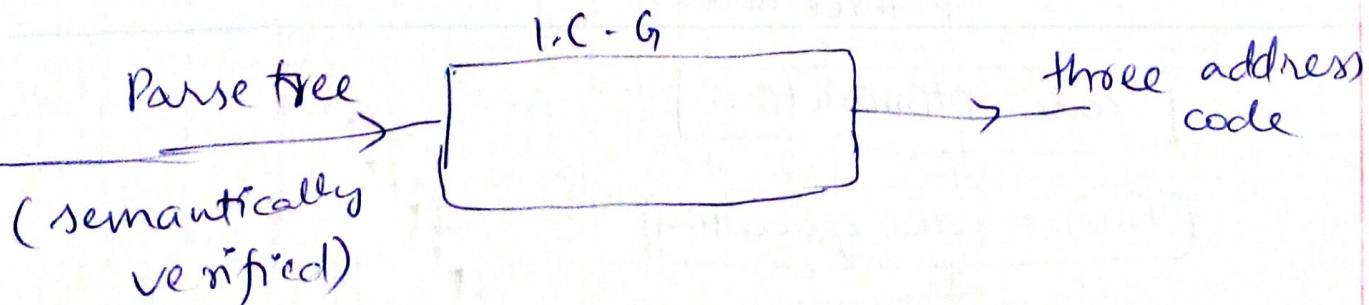


## ③ Semantic Analyzer

- Verifies whether , Parse tree is meaningful or not means, verifies, semantically .



## ④ Intermediate Code - Generation:-



(5)

## Code Optimisation:

- Code optimization like reducing # of address codes (or) lines is done.

Example:-

$x = a + b * c \quad (* \text{ comment } *)$

lexeme

↓ converts stream of lexemes → stream of tokens

lexical analyser

↓

$id = id + id * id$

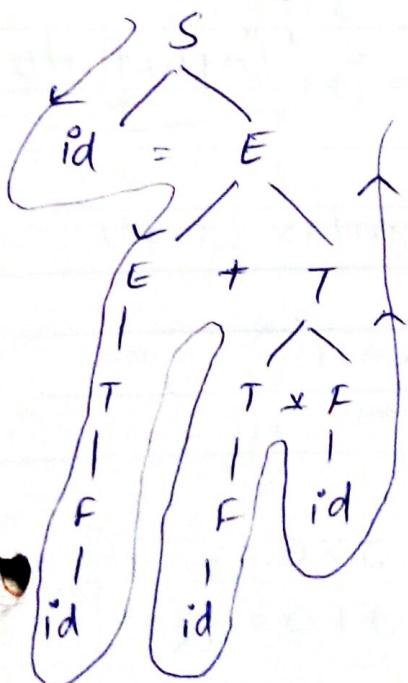
→ 1. A programming language  
rules can be represented using ~~a number of rules~~.

of production rules

2. Syntax Analyser checks.

Syntax Analyser

↓  
(Parse tree)



①  $S \rightarrow id = E;$   
 ②  $E \rightarrow E + T / T$   
 ③  $T \rightarrow T * F / F$   
 ④  $F \rightarrow id$

2. These set of rules are grammar, mostly context-free grammar.

3. For example, the first statement, represents,

① A statement can be, Identifier = Expression followed by semi-colon.

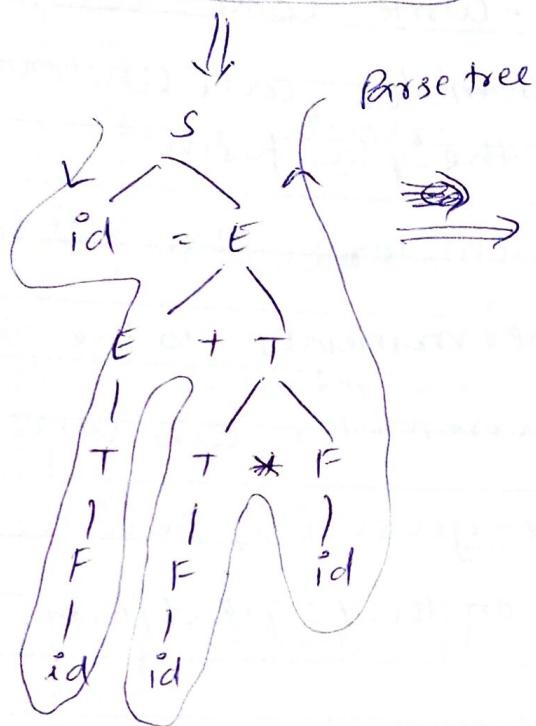
② An expression can be expression plus term (or) term

③ A term can be term into factor (or) factor

④ A factor can be simply a identifier.

- A Syntax Analyzer, checks, whether we have given input, according to the format we needed (i.e., CFG, (set of prod. rules of a programming lang))
- while constructing parse tree, it<sup>①</sup> starts with 'S', start symbol, ② statement is Identifier & Expression
- ③ An expression = expression + term
- ④ A term is term \* factor
- ⑤ A factor can be identifier.
- We can say the input is in required format only if, "yield of parse tree" = "input to Syntax Analyzer". else ~ Syntax Error"

## Syntax Analysis



## Semantic Analyzer

↓  
(Parse tree semantically verified)

- Type checkings
- Variable declarations etc

## Intermediate C.G.

→ generates address codes

- Popular one = 3-address code
- 3-add- codes per line atm,

$$t_1 = b * c$$

$$t_2 = a + t_1$$

$$x = t_2$$

Till here :- Compi Platform  
- independent.

## Code-Optimizer

↓ reduce - # of lines

$$t_1 = b * c ;$$

~~$$t_2 = a + t_1 ;$$~~

## ↓ Target - Code - Generator

mul R<sub>1</sub>, R<sub>2</sub>  
add R<sub>0</sub>, R<sub>0</sub>  
mov R<sub>2</sub>, \*

a → R<sub>0</sub>  
b → R<sub>1</sub>  
c → R<sub>2</sub>

- write code which assembler can understand on the platform
- Assemblers, are ~~code~~ like user manuals to the platform, if assembler can understand code-generated, then we can run on that platform.

- 
- Syntax - Analyzers, heart - of - the - compiler

- mostly, front - end phases → (starting part of the compiler) are implemented by → Syntax Analyzer  
Parser

Lexical Analyzer

Syntax Analyzer

Semantic Analyzer

Intermediate Code G

} Syntax Analyzer / Parser

"Front End"

- Later stages  $\Rightarrow$  Backend
  - Code optimiser
  - Target-code generator

## Lexical Analyser:

- It reads program by "character by character" and "line by line" and generates tokens.
- It removes "comments", "whitespaces" & if any errors, it shows up.

Ex:

int| max(|x|,|y|)

int| x|,| y|;

/\* find max of n and y \*/

{|

return| (|x|>|y| ? |x| : |y|);|

|3|

not counted as tokens.

→ Totally, there are '25' tokens.

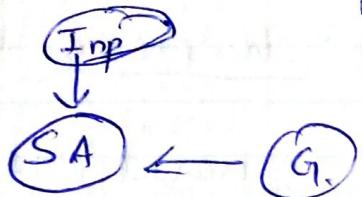
Ex2:

printf| ((|%d| Hai|)|,|{||+2||});| →

→ entire string literal is one-token  
→ so, '8' tokens in total.

## Syntax Analyzer:

- For Syntax Analyzers, we need Grammar & Input to verify whether input is correct or not



$$G = \{ V, T, P, S \}$$

$$P = \{ E \rightarrow E + E \mid E * E \mid id \}, V = \{ E \}, T = \{ +, *, id \}$$

Ex:-  $id + id * id$

$$E \Rightarrow E + E$$

$$\Rightarrow id + E$$

$$\Rightarrow id + E * E$$

$$\Rightarrow id + id * E$$

$$\Rightarrow id + id * id$$

→ left-most derivation

→ -

$$E \Rightarrow E + E$$

$$\Rightarrow E + E * E$$

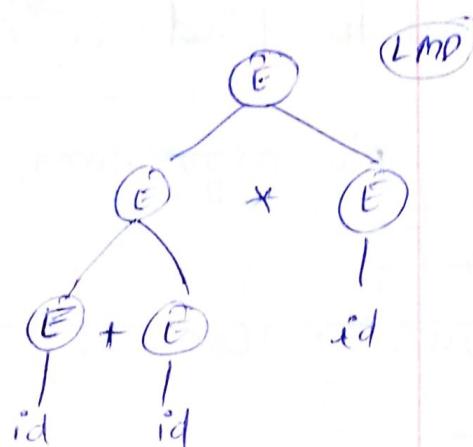
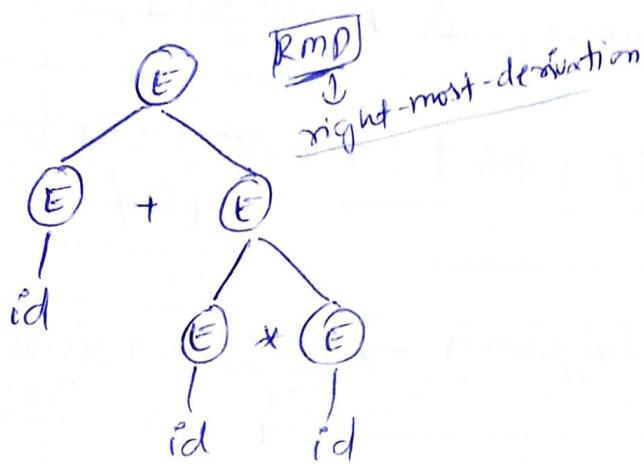
$$\Rightarrow E + E * id$$

$$\Rightarrow E + id * id$$

$$\Rightarrow id + id * id$$

→ (Right-most derivation)

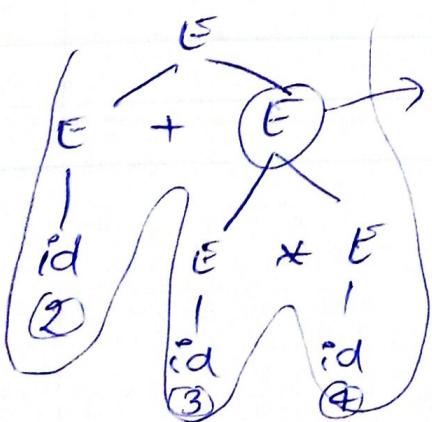
## Parse-tree



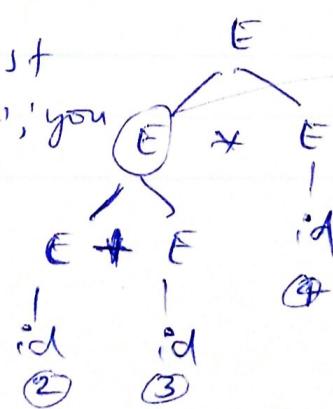
- \* For a given grammar, for a given string, if we have more than one Parse tree, then (or) more than one LMD or RMD

Our grammar is "Ambiguous"

- If we have '2' parse trees, then parser will get confused, as like which one to "generate"
- Even though, o/p looks similar, but there will be difference. Example =  $2 + 3 * 4$



It'll evaluate first  
as to evaluate  $+$ , you  
need it.  
 $\Rightarrow 2 + 3 * 4$   
 $= 14$



but, in this case,  
this will get  
evaluated  
first  
 $\Rightarrow 2 * 3 * 4$   
 $= 20$ ,  
= wrong  
answer

$$2 + 3 * 4$$

## Questions

① How to find a  $\xrightarrow{\text{whether}}$  Grammar is Ambiguous?

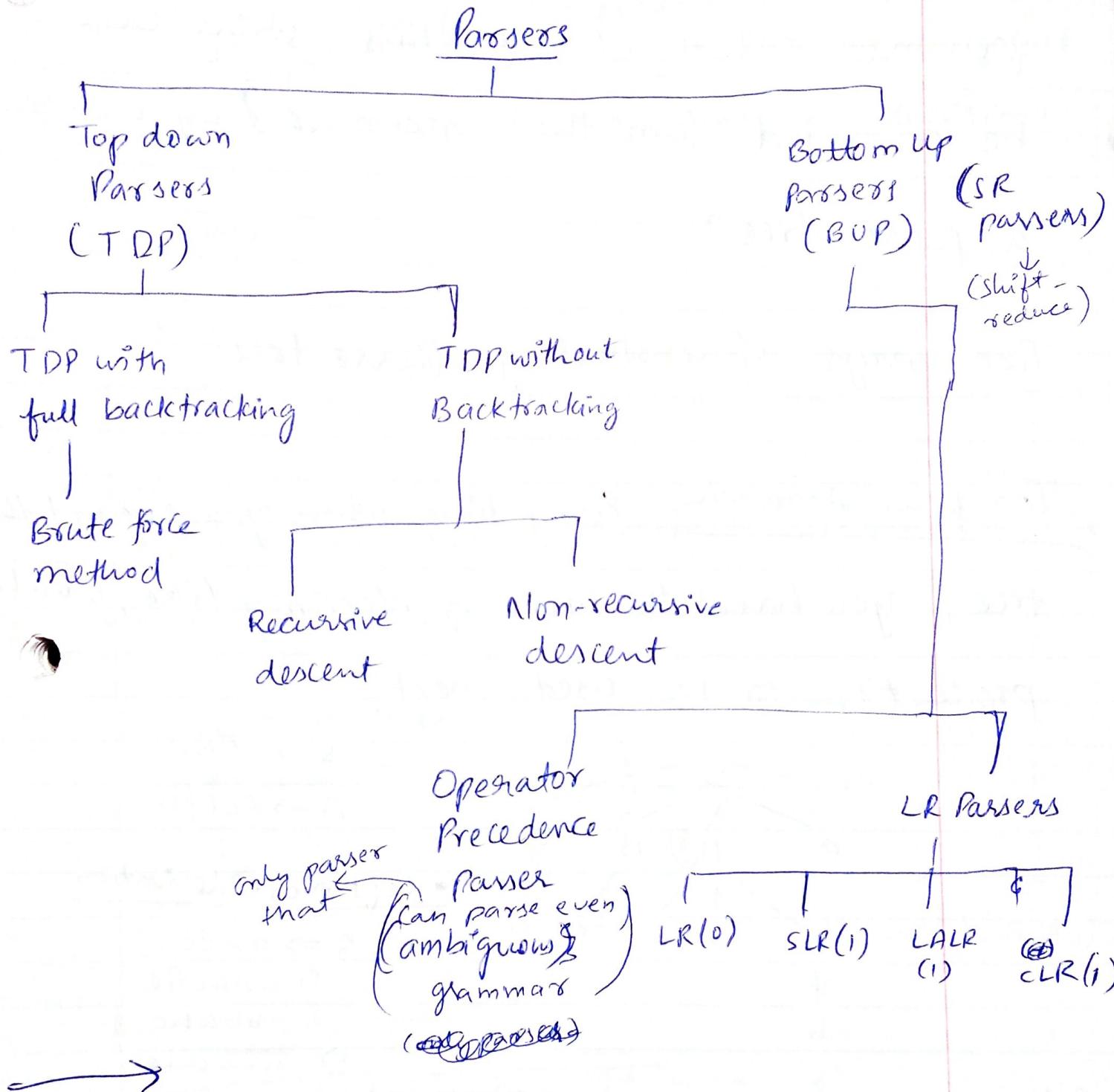
No algorithms, hit & trial manner, we have to find

② How to convert Ambiguous  $\rightarrow$  An Unambiguous grammar?

No algorithms, hit & try

## Passes:

## Parser :-

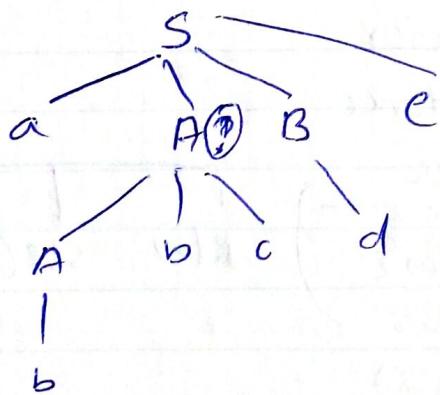


Role/Duty of a Parser: Given a string, given a grammar, check whether this string can be generated from the Grammar by using a parse tree?

Two ways of constructing Parse tree

Top-Down Approach: Every time when you expand the tree, you have to make a decision like, which production to be used next?

Ex:



$$S \rightarrow aABe \\ A \rightarrow Abc/b$$

• left most derivation

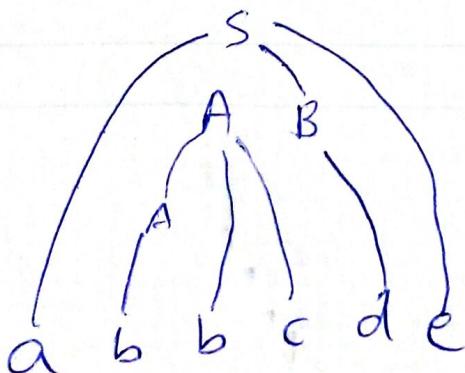
$$S \Rightarrow aABe \\ \Rightarrow aAbcBe \\ \Rightarrow abbcbBe \\ \Rightarrow abbcde$$

Bottom Up: Take the terminals & go up

tasks:

• when to reduce?

• Reverse of Right most derivation



$$S \Rightarrow aABe$$

$$\Rightarrow aAde \\ \Rightarrow aAbcde \\ \Rightarrow abbcde$$

## LL(1) Parser :-

left to right

Lmp.

no. of look ahead's  
(how many symbols you are going to  
see, when you make a decision),

## Operations:

- First()
- Follow()

→ First() :-

• It is the first terminal of the production.  
when you derive all the strings

Ex:-

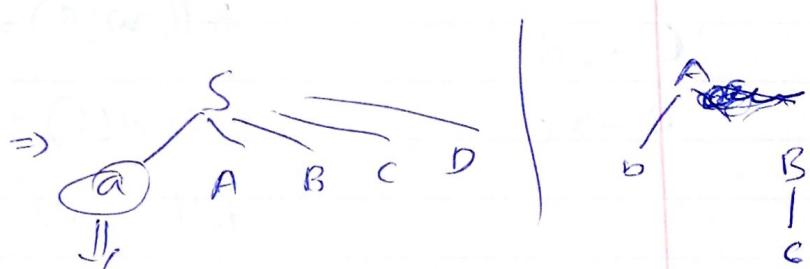
$$S \rightarrow aABC\emptyset$$

$$A \rightarrow b$$

$$B \rightarrow c$$

~~$$C \rightarrow \emptyset$$~~

$$\emptyset \rightarrow e$$



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

$$\text{First}(A) = b; \text{First}(B) = c$$

① If  $S \rightarrow ABCD$   
 $A \rightarrow b/c$

$$B \rightarrow c$$

⇒ First(S) can be 'b' or 'c'.

→ follow():

• what ~~symbols~~ is the terminal which could

follow a variable in the process of derivation

- $s \$$   $\rightarrow$  's' followed by '\$', as ~~the~~ always

input ends with ~~'\$'~~ '\$'

- $s \rightarrow ABCD$

$$\Rightarrow ABCD \$ \Rightarrow \text{follow}(D) = \$$$

$$\epsilon \rightarrow d$$

$$\Rightarrow ABdD \$ \Rightarrow \text{follow}(B) = d$$



- $s \rightarrow ABCD$

$$\text{follow}(s) = d\{\$\}$$

$$A \rightarrow b|\epsilon$$

$$\text{follow}(A) = \text{first}(BCD) = \{c\}$$

$$B \rightarrow C$$

$$\text{follow}(B) = \text{first}(CD) = \{d\}$$

$$C \rightarrow d$$

$$\text{follow}(C) = \text{first}(D) = \{e\}$$

$$D \rightarrow e$$

$$\text{follow}(D) = \{ \$ \}$$