

Introduction

A compiler is a program that can read a program in one language - the source language and translate it into an equivalent program in another language - the target language.

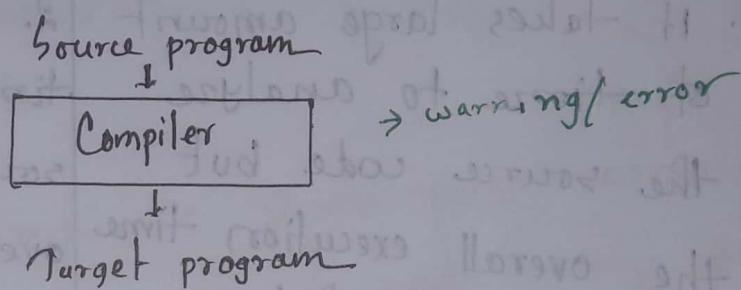


fig: Compiler

An interpreter is another common kind of language processor. Instead of producing a target program as a translation, an interpreter appears to directly execute the operations specified in the source program on inputs supplied by the user.

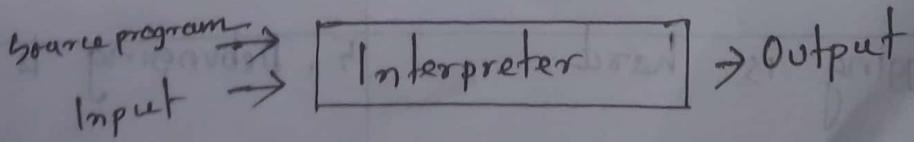


fig: Interpreter

QV

Compiler vs Interpreter

Compiler	Interpreter
1. Scans the whole program and translates it as a whole into machine code.	1. Translates one statement at a time.
2. It takes large amount of time to analyse the source code but the overall execution time is unparatively faster.	2. It takes less amount of time to analyse the source code but the overall execution time is slower.
3. Generates intermediate object code with which further requires linking, hence requires more memory.	3. No intermediate object code is generated, hence memories are efficient.
4. It generates the error message only after scanning the whole program.	4. Continuously translates the program until first error is found.
5. Debugging is hard	5. Debugging is easy
6. C, C++	6. Python, Ruby, Matlab

7. High speed	7. Low speed
8. Compilers are larger in size	8. Smaller in size
9. Error detection is difficult	9. Error detection is easier

#Types of a compiler

- ① Single pass compiler
- ② Multipass compiler

⇒ Single pass compiler : Single pass compiler that processes the source code only once.

⇒ Multipass compiler : It is a type of compiler that processes the source code multiple times to convert the source code to multiple code.

Two major parts of a compiler

Analysis phase :- It breaks up the source code into small parts and creates an intermediate representation of the source phase.

Synthesis phase :- Output of the 1st analysis part acts as input here, and generates the target output.

usage with

code

in general are called

crit

the miss or not both ready

to

the 1st analysis part

①

the 2nd analysis part

②

the 3rd analysis part

③

the 4th analysis part

④

the 5th analysis part

⑤

the 6th analysis part

⑥

the 7th analysis part

⑦

the 8th analysis part

⑧

Language processing system

High level language programs are fed into a series of tools and os components to get the desired code that can be used ~~to~~ by the machine. This is called language processing system.

Preprocessor

Source program inserted into preprocessor. Preprocessor modifies the source program. It removes all the preprocessor directives, and add the respective files.

Example: `#include <stdio.h>` ignore those the header file.

It will do macro expansion, operator conversion.

$$a-- \Rightarrow a = a - 1;$$

Preprocessor

Source program

Preprocessor

Modified source program

Compiler

Target assembly program

Assembler

Relocatable machine code

Linker/ Loader

Target Machine code

Compilers

Compiler during compilation converts pure high level language into assembly language. It also gives errors and warnings.

- Assembly language is a low-level. It is not in a binary form.

Assemblers

for every platform we have a assembler. A assembler for one platform will not work for another platform.

Assembler code $\xrightarrow{\text{converted}}$ Executable machine code.

Code & can be loaded anywhere in the memory.

Loader & Linker

A linker links different object files into a single file/executable file. and

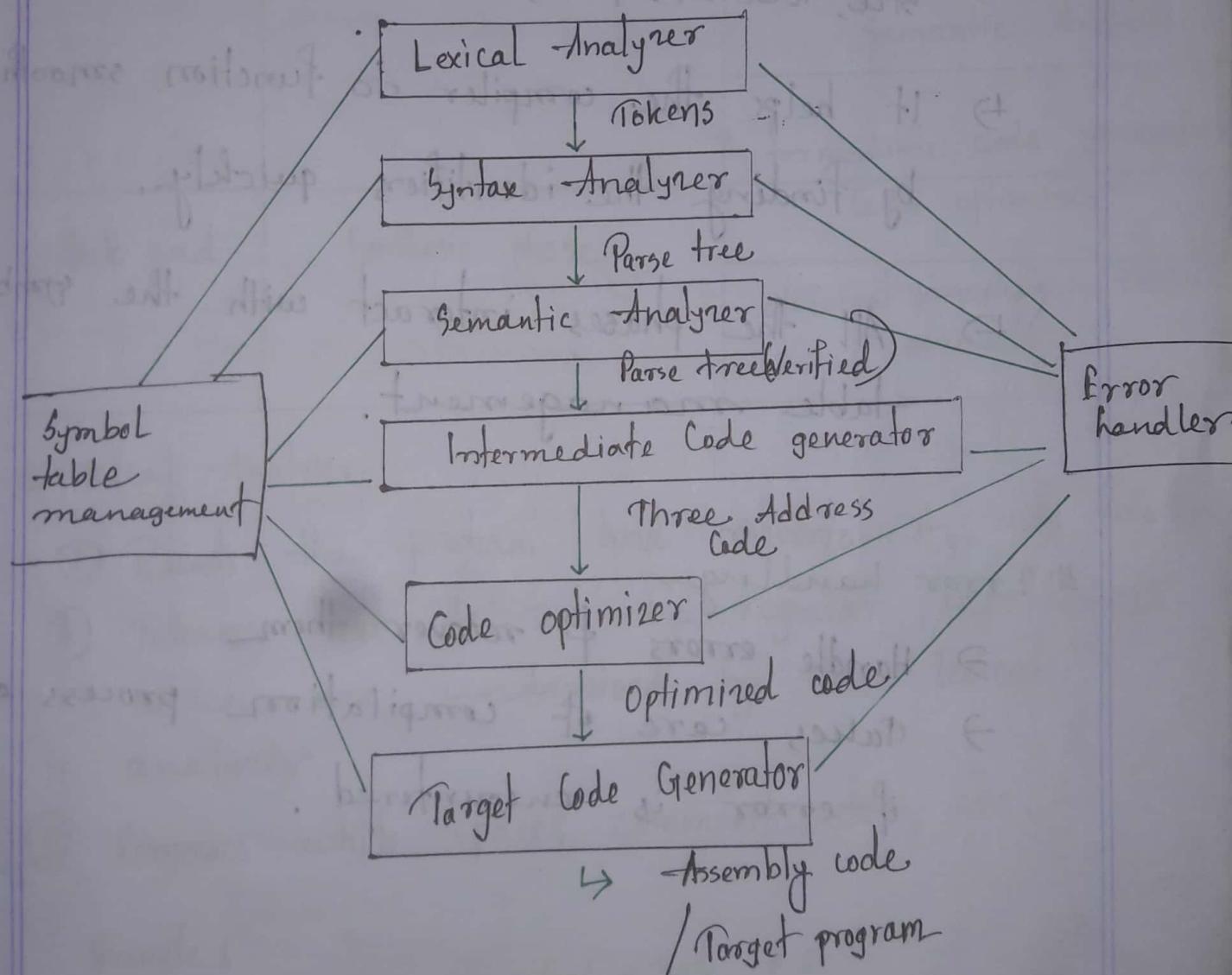
Loader loads that executable file in the memory and executes it.

OK

Preprocessor, assembler and loader/linker are
cousins of a compiler.

#Phases of a compiler

Source program



Symbol table management

- ↳ A data structure which is used by the compiler to store all the information related to identifiers. Example: types, scope, size, location, name, etc.
- ↳ It helps the compiler to function smoothly by finding the identifiers quickly.
- ↳ All the phases interact with the symbol table management.

Error handling

- Handles errors & recover them
- Takes care of compilation process even if error is encountered.

Part	Major phases	Phases
front end	Analysis phase	Lexical Analysis Syntax Analysis Semantic Analysis
		Intermediate code generator
Back end	Synthesis phase	Code optimizer Target machine code generator

Lexical Analyzer

- ① Reads the program and converts it into tokens
- ② Tokens are defined by regular expressions which are understood by the lexical analyzer.
- ③ Removes white spaces, comments, tabs, etc.

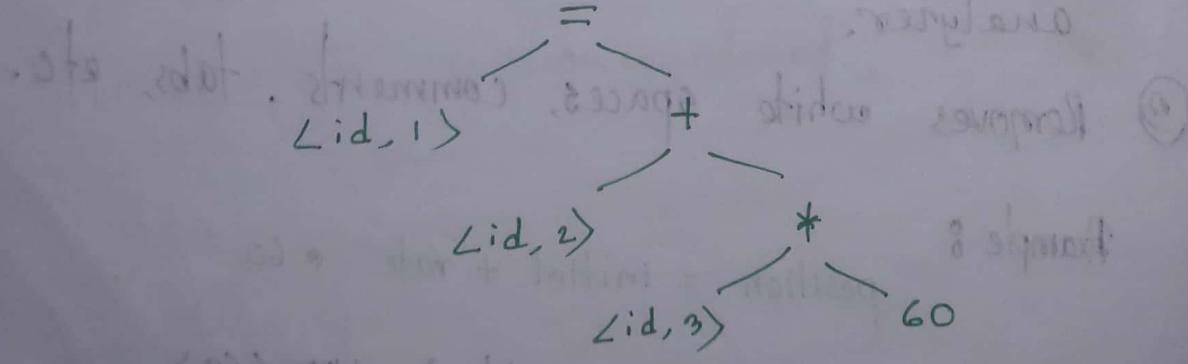
Example :

position = initial + rate * 60

$\Rightarrow \langle id, 1 \rangle \leftrightarrow \langle id, 2 \rangle \leftrightarrow \langle id, 3 \rangle \langle * \rangle \langle 60 \rangle$

#Syntax Analyser

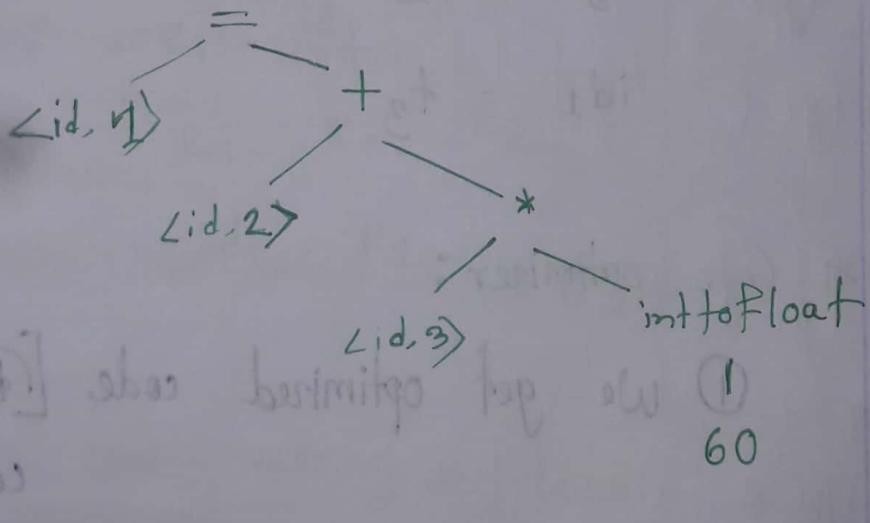
- ① Constructs the parse tree.
- ② Takes the tokens one by one, and uses CFG to construct the parser tree.
- ③ Using the productions of CFG we can represent what the program actually is. The input has to be checked whether it is in the desired format or not.
- ④ Syntax errors can be detected by it if the input is not accurate to the grammar given.
- ⑤ So, it takes input as tokens and checks their syntactic correctness.



Semantic Analyzer

- ① Verifies whether the parse tree is meaningful or not.
- ② It uses the parse tree and info in the symbol table to check the source program for semantic consistency with the language definition.

Type checking, definition of the variables used, label checking, flow control, etc.



✓

Intermediate Code Generator

① Generates the intermediate code.
→ 3 address code

② This code converts into machine language using last two phases of the platform.

$$t_1 = \text{inttofloat}(60)$$

$$t_2 = id_3 * t_1$$

$$t_3 = id_2 + t_2$$

$$id_1 = t_3$$

Code optimizer:

① We get optimized code [Removing unnecessary code]

② Code optimization phase attempts to improve the intermediate code so that it runs faster and works efficiently.

$$t_1 = id_3 * 60.0$$

$$id_1 = id_2 + t_1$$

Target Code Generator

① Final phase of the compiler which generates the target code.

LDF R₂, id₃

MULF R₂, R₂, #60.0

LDF R₁, id₂

ADDF R₁, R₁, R₂

STF id₁, R₁

Grammar and Context Free Grammar

In normal human languages, a grammar is a set of rules that used for proper conversation with each other. So, in the same way for writing computer languages there is a mathematical model of grammar which is used to write the languages in the correct way.

Noam Chomsky gave a mathematical model of grammar for writing computer languages effectively.

Type	Grammar Accepted	Languages Accepted	Automation
Type-0	Unrestricted grammar	Recursively Enumerable language	Turing machine
Type-1	Context sensitive grammar	Context sensitive language	Linear bounded automation
Type-2	Context free grammar	Context free languages	Pushdown automata
Type-3	Regular grammar	Regular language	finite state automation

15

Regular grammar

A grammar 'G' can be formally described using 4 tuples as $G = (V, T, S, P)$ where,

V = set of variables or Non terminal symbols

T = set of terminals

S = start symbol

P = Production rule

A production rule has the form $\alpha \rightarrow \beta$ where α and β are strings on VUT and atleast one symbol of α belongs to V

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

Here

$V = S, A, B$

$T = a, b$

$S = S$

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

$S \rightarrow AB$

$S \rightarrow aB$

$S \rightarrow \underline{ab}$

Types of a regular grammar :

Regular grammars are two types as follows :-

1. Right linear grammar

2. Left linear grammar

Right linear grammar : A grammar is said to be right linear if all the non-terminals are in the right side.

$$\left. \begin{array}{l} S \rightarrow xB \\ S \rightarrow x \end{array} \right\} V \rightarrow T^* V \text{ or } V \rightarrow T^*$$

Left linear grammar : A grammar is said to be left linear if all the non-terminals are in the left side.

$$\left. \begin{array}{l} S \rightarrow Bx \\ S \rightarrow x \end{array} \right\} V \rightarrow V T^* \text{ or } V \rightarrow T^*$$

A grammar is regular, if it is either ~~regular~~
right linear or left linear. This means, all
productions in the grammar have to be
completely left linear or completely right linear
but not mixed left-linear and right linear.

Linear grammar : Grammars, in which each rule is in right linear or left linear form, i.e right linear rules and left linear rules can be mixed, is called linear.

$$S \rightarrow aX$$

$$X \rightarrow Sb$$

$$S \rightarrow \epsilon$$

It is a linear grammar but neither left linear nor right linear.

Context Free Grammar

A context free grammar (CFG) consisting of a finite set of grammar rules where it has 4 tuples.

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

V = Non terminals / Variables

T = Terminals

P = Production rules

S = Starting state symbol

Production rule P will have the relation, $A \rightarrow a$, where $a = (V \cup T)^*$ and $A \in V$

That means, only nonterminals are allowed on the left side.

$$G_1 = \{ (S, A), (a, b), S, (S \rightarrow aAb, A \rightarrow aAb | \epsilon) \}$$

$$V = (S, A)$$

$$T = (a, b)$$

$$S = S.$$

$$P: S \rightarrow aAb, A \rightarrow aAb | \epsilon$$

Context free grammar is used for syntax analysis.

Derivation and Parse tree

#Derivation

A derivation is basically a sequence of production rules, in order to get the input string. During parsing, we take two decisions for some sentential form of input :

- ⇒ Deciding the non-terminal which is to replaced.
- ⇒ Deciding the production rule by which, the non-terminal will be replaced.

To decide which non-terminal to be replaced with production rule, we can have two options.

- ⇒ Left most derivation
- ⇒ Right most derivation

⇒ Left most derivation :- If the sentential form of an input is scanned from left to right, it is called left most derivation.

⇒ Right most derivation : If we scan and replace the input with production rules, from right to left, it is known as right most derivation.

Example : Production rules :

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow id$$

String : id + id * id ,

Left most	Right most
$E \rightarrow E * E$	$E \rightarrow E + E$
$E \rightarrow E + E * E$	$E \rightarrow E + E * E$
$E \rightarrow id + id * id$	$E \rightarrow E + E * id$
$E \rightarrow id + E * E$	$E \rightarrow E + id * id$
$E \rightarrow id + id * E$	$E \rightarrow id + id * id$
$E \rightarrow id + id * id$	

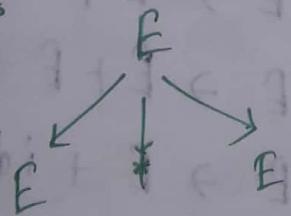
Parse tree

A parse tree is a graphical depiction of a derivation. It is convenient to see how strings are derived from the start symbol. The start symbol of the derivation becomes the root of the parse tree. Let us see this by an example from the last topic.

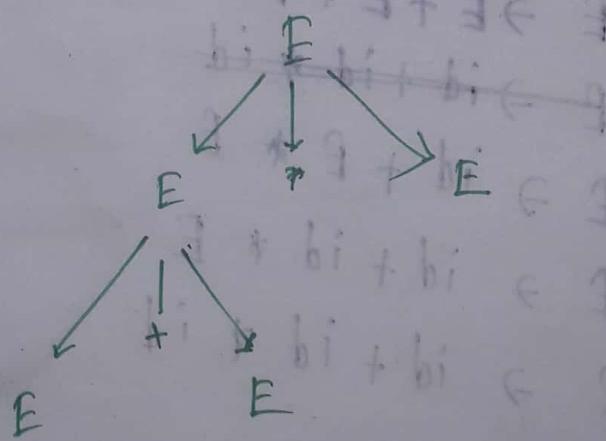
Grammar :

$$\begin{aligned} E &\rightarrow E + E \\ E &\rightarrow E * E \\ E &\rightarrow id \end{aligned}$$

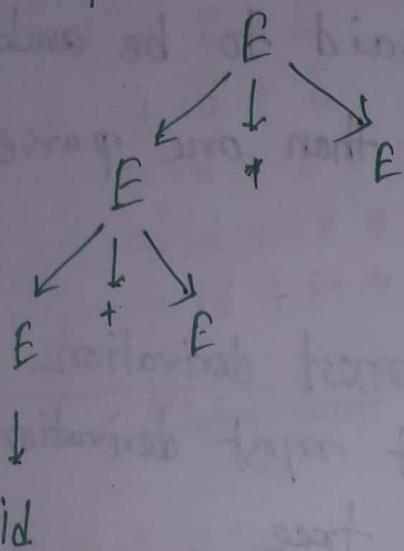
Step 1 :



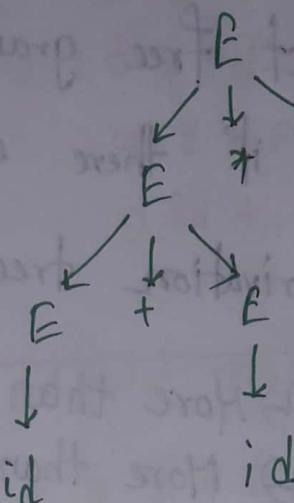
Step 2 :



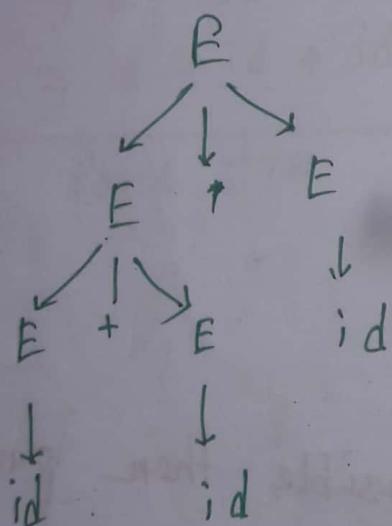
Step 3:



Step 4:



Step 5:



Properties of a parse tree

1. All leaf nodes are terminal.
2. All interior nodes are non-terminals
3. Start from the root node, which is the starting symbol.

24

Ambiguous Grammar

A context free grammar is said to be ambiguous grammar if there exists more than one parse tree / derivation tree.

- ↳ More than one left most derivation
- ↳ More than one right most derivation
- ↳ More than one parse tree

Grammar:

$$E \rightarrow E+E$$

$$E \rightarrow E * E$$

$$E \rightarrow id$$

String: $id + id * id$

If we have two parse trees possible then parser will get confused about which one to generate.

left most derivation | Left most derivation

$$\begin{aligned}
 E &\rightarrow E + E \\
 &\rightarrow id + E \\
 &\rightarrow id + E * E \\
 &\rightarrow id + id * E \\
 &\rightarrow id + id * id
 \end{aligned}$$

$$\begin{aligned}
 E &\rightarrow E * E \\
 &\rightarrow E + E * E \\
 &\rightarrow id + E * E \\
 &\rightarrow id + id * E \\
 &\rightarrow id + id * id
 \end{aligned}$$

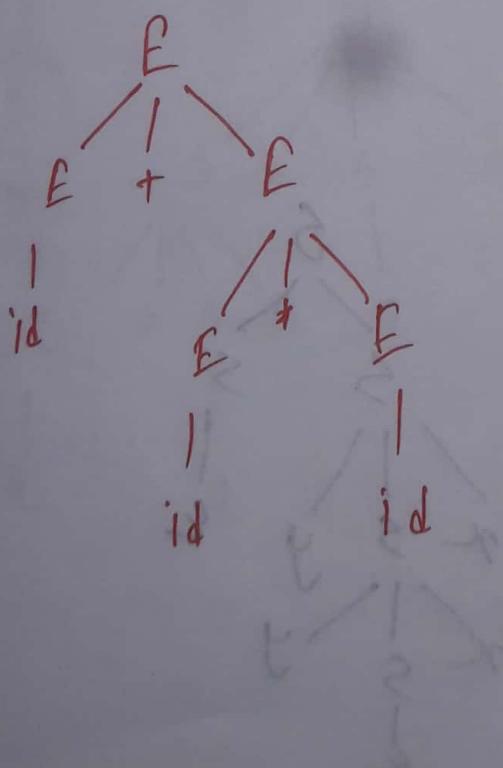
Right most derivation

$$\begin{aligned}
 E &\rightarrow E E + E \\
 &\rightarrow E + E * E \\
 &\rightarrow E + E * id \\
 &\rightarrow E + id * id \\
 &\rightarrow id + id * id
 \end{aligned}$$

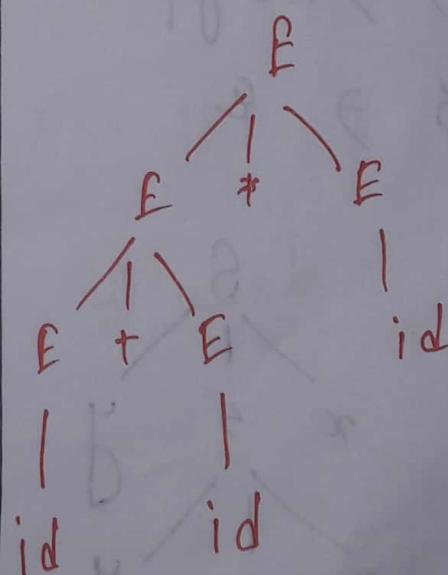
Right most derivation

$$\begin{aligned}
 E &\rightarrow E * E \\
 &\rightarrow E * id \\
 &\rightarrow E + E * id \\
 &\rightarrow E + id * id \\
 &\rightarrow id + id * id
 \end{aligned}$$

Parse tree



Parse tree



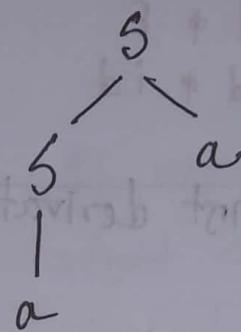
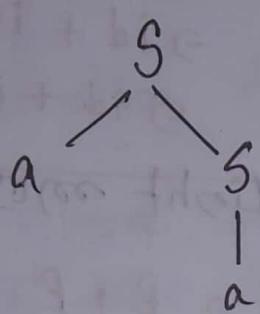
26

68

Check the grammar is ambiguous or not

$$(i) \quad S \rightarrow aS \mid Sa \mid a$$

string $\Rightarrow aa$

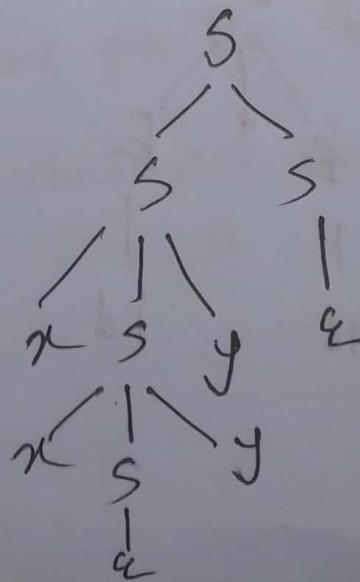
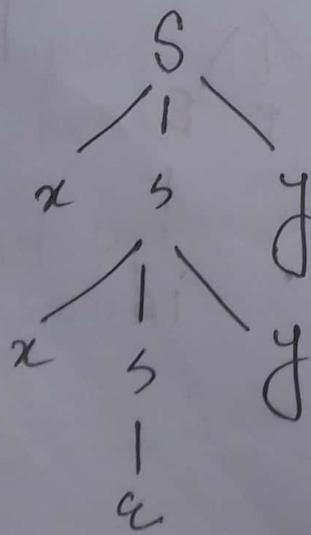


We have got two parse tree. So, this particular grammar is ambiguous.

$$(ii) \quad S \rightarrow xS \mid y \mid ss$$

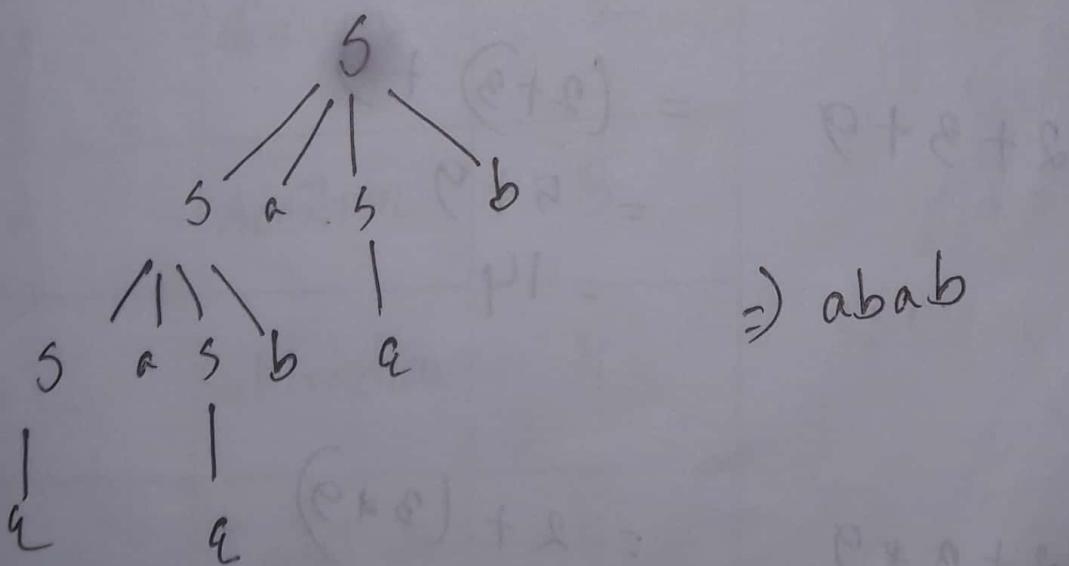
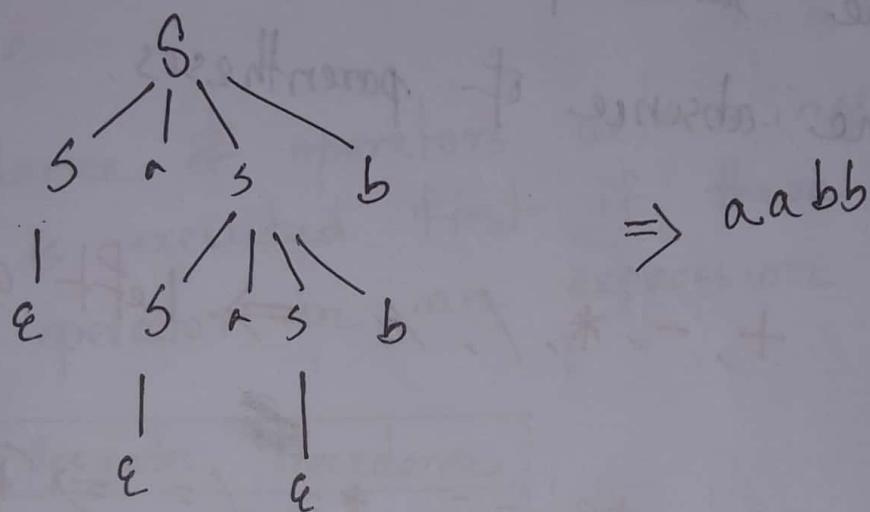
string $\Rightarrow xxyy$

$$S \rightarrow e$$



Ambiguous grammar.

(iii) $S \rightarrow S_aSb \mid \epsilon$



Unambiguous grammar

#Associativity

The associativity of an operator is a property that determines how operators of the same precedence are grouped in the absence of parentheses.

$+, -, *, /, \% \Rightarrow$ Left associative

$=, !, +\!, -\!, *\!, /\!, \% \Rightarrow$ Right associative

$$2+3+9 = (2+3)+9 \\ = 5+9 \\ = 14$$

$$2+3*9 = 2+(3*9) \\ = 2+27 \\ = 29$$

Here, + and * both have the left are left associative. So, we will consider the precedence of operator to decide which operator will take progress first.

#Precedence

The precedence of operators determines which operator is executed first if there is more than one operator in an expression.

Symbol	Operator	Precedence
*	Multiplication	1
/	division	2
+	Addition	3
-	Subtraction	4
=	Equal to	5

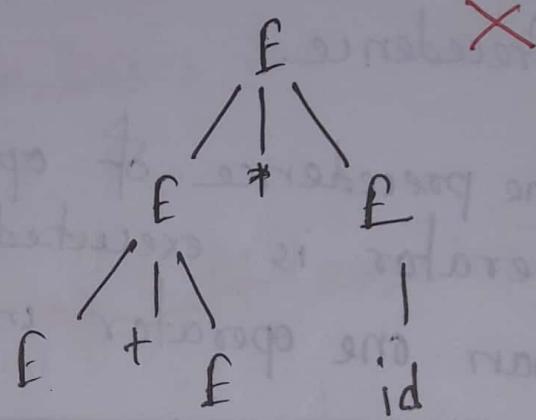
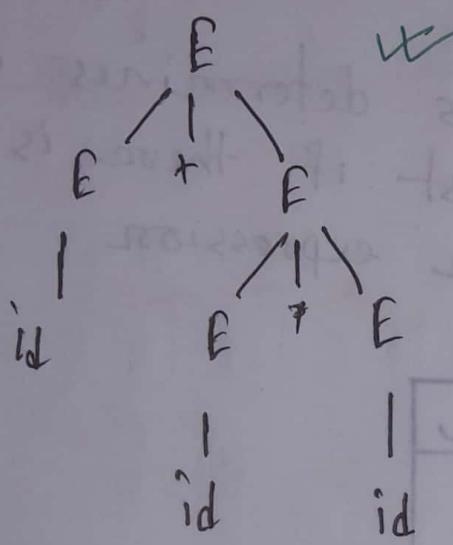
30

Ambiguous grammar to unambiguous grammar.

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

(id + id * id)

This is an ambiguous grammar. because.



$$\begin{aligned} & id + (id * id) \\ & 1 + (3 * 4) \\ & = 1 + 12 = 13 \end{aligned}$$

$$\begin{aligned} & (id + id) * id \\ & (1 + 3) * 4 \\ & = 4 * 4 \\ & = 16 \end{aligned}$$

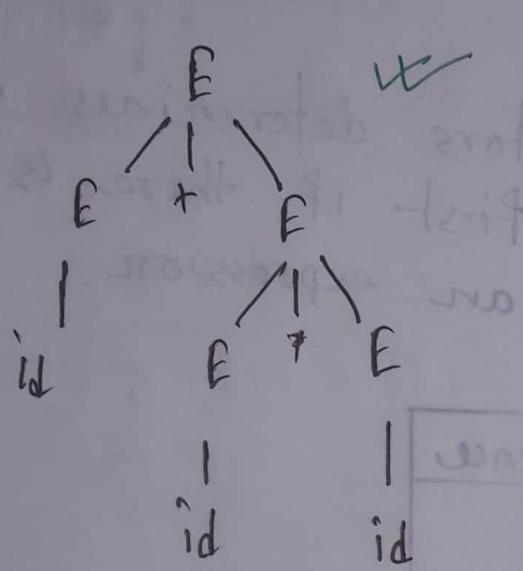
Compiler becomes confused when of internal with ambiguous grammar. So, it is necessary

30

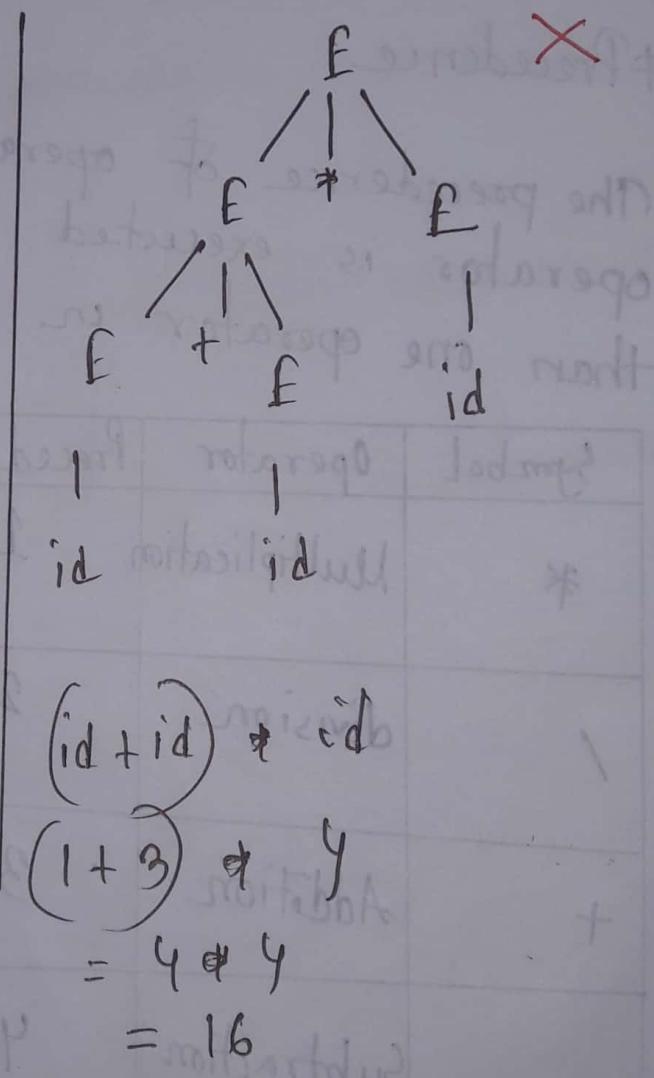
Ambiguous grammar to unambiguous grammar.

$$E \rightarrow E+E \mid E * E \mid id \quad (id + id * id)$$

This is an ambiguous grammar, because,



$$\begin{aligned}
 & id + (id * id) \\
 & 1 + (3 * 4) \\
 & = 1 + 12 = 13
 \end{aligned}$$



Compiler becomes confused when it interacts with ambiguous grammar. So, it is necessary

for us to improve the grammar and
remove ambiguity.

Depending on two criterias we convert an ambiguous grammar into an unambiguous grammar.

→ Associativity

→ Precedence

We will introduce recursion for solving associativity problem. Left recursion for left associative and right recursion for right associative.

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow id$$

Left Recursion

If the variable (RHS) is same as the leftmost variable then we say that particular grammar has left recursion.

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

Two types:

⇒ Direct left recursion $[A \rightarrow A^i]$

⇒ Indirect left recursion $[X \rightarrow Ab]$
 $A \rightarrow Xt$

Top-down parser can not accept left recursion

that's why we have to remove it. But we

have to keep in mind that after removing

the left recursion the grammar should

be the one as it was before.

33

Recursion two types:

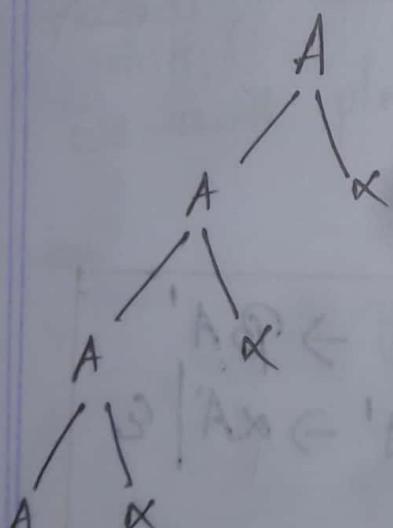
\Leftarrow Left recursion

\Rightarrow Right recursion

Left recursion

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

$$\begin{array}{c} A \\ \downarrow \\ \beta \end{array}$$



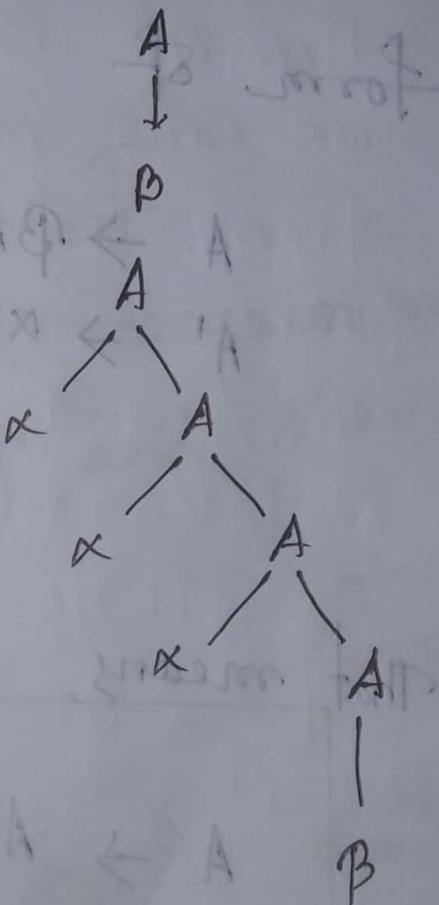
f()

Infinite loop

{ f();
others;

Right recursion

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



f()

} other;
f();

m

Conversion of left recursion to right recursion :

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

This grammar is a left recursive grammar.
To convert it we will write it in the form of

$$\begin{aligned} A &\rightarrow \beta A' \\ A' &\rightarrow \alpha A' | \epsilon \end{aligned}$$

That means,

$$A \rightarrow A\alpha | \beta \equiv A \rightarrow \beta A' \quad A' \rightarrow \alpha A' | \epsilon$$

95

Example : 01

$$E \rightarrow E + T \mid T$$

$$e | (1) \leftarrow s$$

$$e | a \leftarrow s$$

Sol 2 :

This is a left recursive grammar. We will compare it with $A \rightarrow A\alpha \mid \beta$

So we can see,

$$\frac{E}{A} \rightarrow \frac{E + T}{A} \mid \frac{T}{\beta}$$

According to the rule of left recursion removing we will place, $A \rightarrow \beta A'$ and $A' \rightarrow \alpha A' \mid \epsilon$

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

CAns

36

Example 02

$$A \rightarrow A\alpha | \beta \Rightarrow A \rightarrow \beta A' \\ A' \rightarrow \alpha A' | \epsilon$$

$$S \rightarrow (L) | x$$

$$L \rightarrow L, S | S$$

Sol 2:

Here in the first line of the grammar we can see that the production is free from left recursion. So, we don't have to apply anything there. But in the second line left recursion is available and we have to remove it.

$$S \rightarrow (L) | x$$

$$L \rightarrow SL'$$

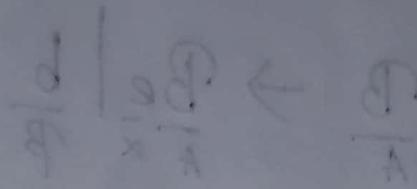
$$L' \rightarrow , SL' | \epsilon$$

(2)

37

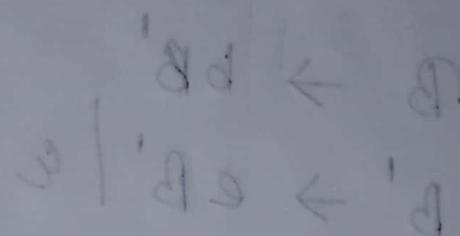
Example 803

$$S \rightarrow S051/01$$

Solⁿ

$$S \rightarrow 01S'$$

$$S' \rightarrow 0S1S' | \epsilon$$



Example 804

$$A \rightarrow ABd | A\alpha | a$$

$$B \rightarrow Be | b$$

$$\overline{A \rightarrow A\alpha | P}$$

$$\overline{\begin{array}{l} A \rightarrow BA' \\ A' \rightarrow \alpha A' | \epsilon \end{array}}$$

Solⁿ

$$\frac{A}{A} \rightarrow \frac{ABd}{A\alpha} | \frac{A\alpha}{A} | \frac{a}{B}$$

So,

$$A \rightarrow aA' | aA' | \epsilon$$

$$A' \rightarrow BdA' | \epsilon$$

36

Again,

$$\frac{B}{A} \rightarrow \frac{Be}{A\bar{\alpha}} \mid \frac{b}{B}$$

$$A \rightarrow A\alpha \mid \beta$$

$$A' \rightarrow \alpha A' \mid e$$

$$B \rightarrow bB'$$

$$B' \rightarrow eB' \mid e$$

So,

$$A \rightarrow aA'$$

$$A' \rightarrow BdA' \mid aA' \mid e$$

$$B \rightarrow bB'$$

$$B' \rightarrow eB' \mid e$$

DA

39

Example 05

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T \times F \mid F$$

$$F \rightarrow id$$

Sol 2^o

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow \times FT' \mid \epsilon$$

$$F \rightarrow id$$



Left Factoring

If RHS of more than one production starts with the same symbol, then such a grammar is called as grammar with common prefixes.

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \alpha\beta_3$$

→ Such grammar creates a problematic situation for top-down parsers

⇒ Top down parsers can not decide which production must be chosen to parse the string in hand.

To remove this confusion, we use left factoring.

41
Left factoring is a process by which the grammar with common prefixes is transformed to make it useful for topdown parsers.

- ⇒ We make one production for each common prefix.
- ⇒ The common prefix may be a terminal or a nonterminal or a combination of both.
- ⇒ Rest of the derivation is added by new productions.

The grammar we obtained after the process is called as left factored grammar.

v²

Examples:

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \alpha\beta_3$$

Soln

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2 \mid \beta_3$$

or

$$S \rightarrow iEtS \mid iEtSeS \mid a$$

$$E \rightarrow b$$

Soln:

$$S \rightarrow iEtSS' \mid a$$

$$S' \rightarrow e \mid es$$

$$E \rightarrow b$$

43

$$\underline{02} \quad A \rightarrow aAB \mid aBc \mid aAc$$

Sol 2

$$A \rightarrow aA'$$

$$A' \rightarrow AB \mid Bc \mid Ac$$

Again, there is a common prefixes,

$$A \rightarrow aA'$$

$$A' \rightarrow Bc \mid AA''$$

$$A'' \rightarrow B \mid c$$

03

$$S \rightarrow bSSaaS \mid bSSaSb \mid bSb \mid a$$

Sol 2

$$S \rightarrow bSS' \mid a$$

$$S' \rightarrow SaS \mid S \cdot Sb \mid b$$

$$S \rightarrow bSS' \mid a$$

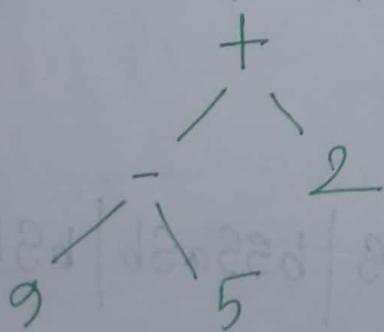
$$S' \rightarrow \cancel{S} \cancel{S} \mid S a S'' \mid b$$

$$S'' \rightarrow aS \mid Sb$$

44
#Abstract Syntax: A useful starting point for designing a translator is a data structure called an abstract syntax tree.

In an abstract syntax tree for an expression, each interior node represents an operator; the children of the nodes represent the operand of the operator.

What are the significant parts of the expression?



A sum expression has its two operand expressions as its significant parts.

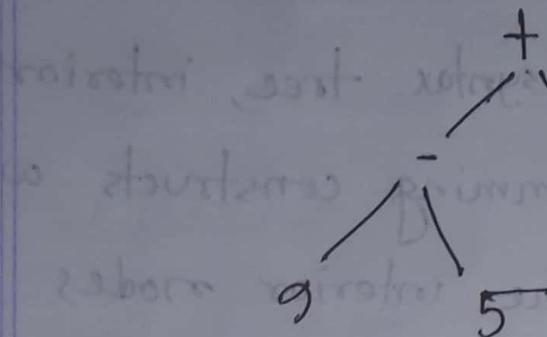
45
Concrete Syntax : Abstract syntax trees, or simply syntax trees, resemble parse trees to an extent. However, in the syntax tree, interior nodes represent programming constructs while in the parse tree, the interior nodes represent nonterminals.

Many nonterminals of a grammar represent programming constructs, but others are helpers of one sort or another such as those representing terms, factors, or other variations of expressions.

In the syntax tree, these helpers typically are not needed and are hence dropped. To emphasize the contrast, a parse tree is sometimes called a concrete syntax tree, and the underlying grammar is called a concrete syntax for the language.

46

what does the expression look like?



The same sum expression can look in different ways:

infix	$2 + 3$
prefix	$(+ 2 3)$
postfix	$(2 3 +)$
JVM	bipush 2 bipush 3 iadd
English	the sum of of and 3

Removal of white space and comments //

The expression translator sees every character in the input, so extraneous characters, such as blanks, will cause it to fail.

Most languages allow arbitrary amounts of white space to appear between tokens. Comments are likewise ignored during parsing, so they may also be treated as white space.

If white space is eliminated by the lexical analyzer, the parser will never have to consider it.

```
i = 0;
while (string[i] != '\0')
```

```
{ if (string[i] == whitespace or tab or newline)
```

```
    i++;
```

```
}
```

#Reading Ahead: A lexical analyzer may need to read some characters before it can decide on the token to be turned to the parser.

For example, a lexical analyzer for C or Java must read ahead after it sees the character >. If the next character is =, then > is part of the character sequence >=, the lexeme for the greater than or equal to operator. Otherwise, > itself forms the greater than operator, and the lexical analyzer has read one character too many.

#Common Programming Errors

1. Lexical errors include misspelling of identifiers, keywords, etc.
2. Syntactic errors include misspaced semicolon, extra or missing braces.
3. Semantic errors include type mismatch.
4. Logical errors can be anything by a programmer.

#Tokens

A token is a sequence of characters that can be treated as a unit / single logical entity.

Typical tokens are : keywords, identifiers, separators, etc.

`int a = 15;`

5 tokens here.

Keywords : for, if, int, while, etc

Identifiers : variable name, function name, etc.

Operators : +, -, *, =, etc

Separators : , ; etc

Constants : Any numbers

Lexeme

A sequence of input characters that comprises a single token is called a lexeme.

uint "a", "=", "18"

Pattern

Pattern is a rule describing all these lexemes that can represent a particular token in a source program.

These are some predefined rules for every lexeme to be identified as a valid token. These rules are defined by grammar rules, by means of pattern.

62

Examples:

```
int main()
```

```
{ int b;  
}
```

⇒ Lexical analyzer read int and finds it to be valid and accepts as token.

⇒ read main and finds it to be a valid function name and accepts as token.

⇒ Similarly (,), {, }, int, b, ;, - }.

Total tokens: 9

⇒ int main ()

// 8 variables declared below.

int a, b;

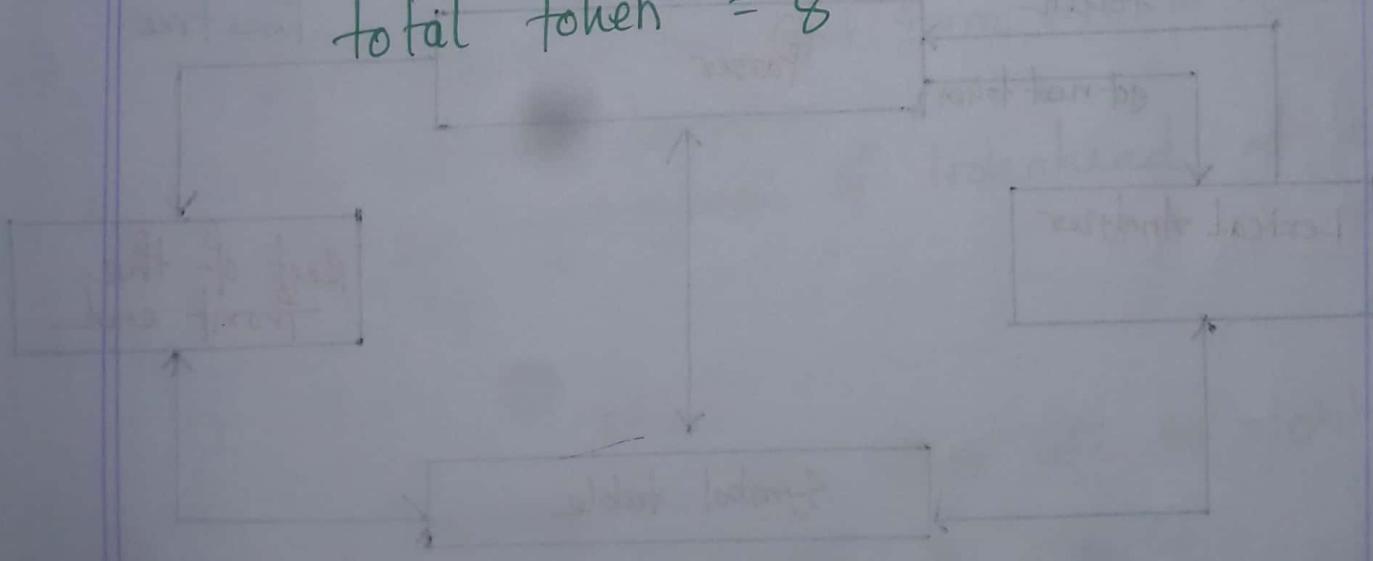
a = 10;

return 0; total tokens = 18

}

⇒ printf("%d", &j);

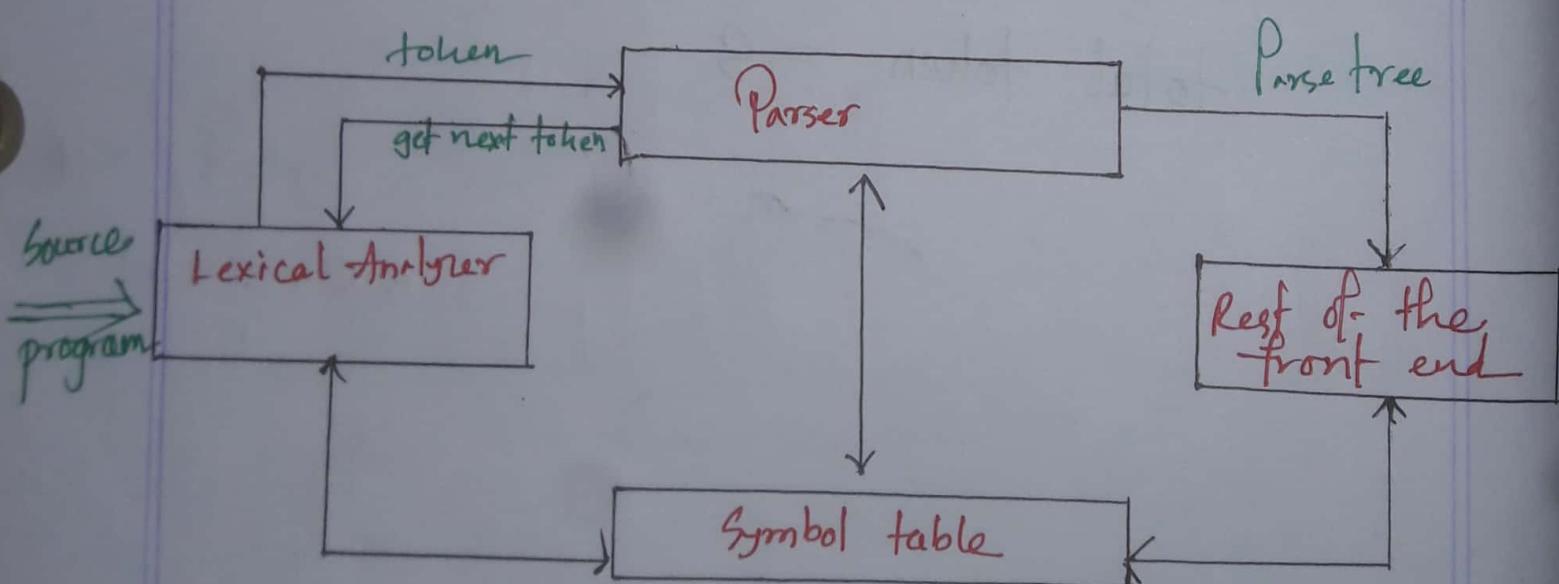
total token = 8



Syntax Analysis

() nisan tri

- ⇒ The parser obtains a string of tokens from the lexical analyzer.
- ⇒ Verifies that the string can be the grammar for the source language.
- ⇒ Detects and reports any syntax errors and produces a parse tree from which intermediate code can be generated.



55
some important topics required to understand parsing properly.

1. Context free grammar
2. Ambiguity
3. Left Recursion
4. Left factoring

Types of Parsers

LL(1)

L \Rightarrow Scanning the input from left to right

L \Rightarrow Leftmost derivation

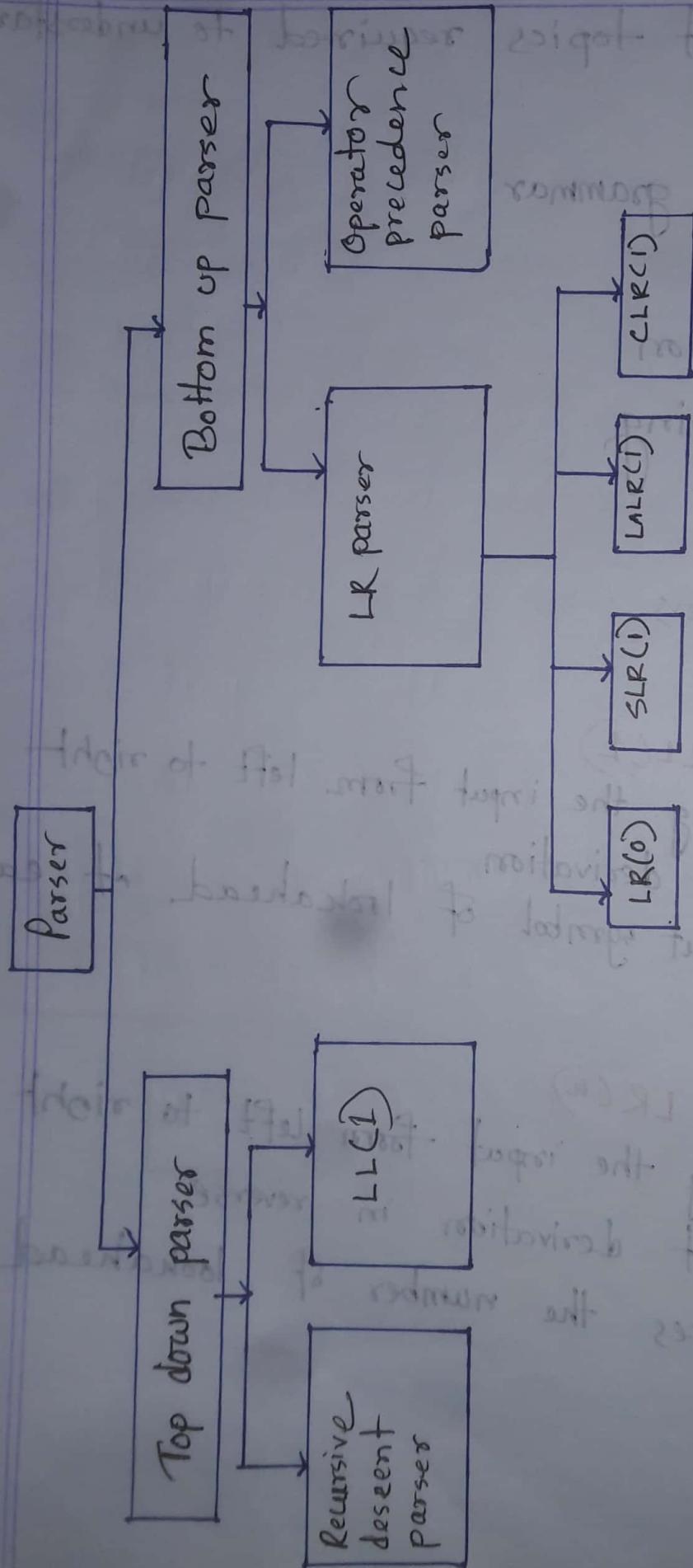
1 \Rightarrow One input symbol of lookahead at each step

LR(k)

L \Rightarrow Scanning the input from left to right

R \Rightarrow Rightmost derivation in reverse

k \Rightarrow k denotes the number of lookahead



$SLR(1) \Rightarrow$ Simple LR parser

$LALR(1) \Rightarrow$ Look-Ahead LR parser

$CLR(1) \Rightarrow$ Canonical LR parser

Top down parser

Top down parsing can be viewed as the problem of constructing a parse tree for the input string, starting from the root and creating the nodes of the parse tree in preorder. Equivalently, top down parsing can be viewed as finding a leftmost derivation for an input string. To design a top-down parser, a grammar must not be in a left recursive form and also free from common prefixes problem.

Design a top down parser for the following grammar where the input is $\text{id} + \text{id} * \text{id}$

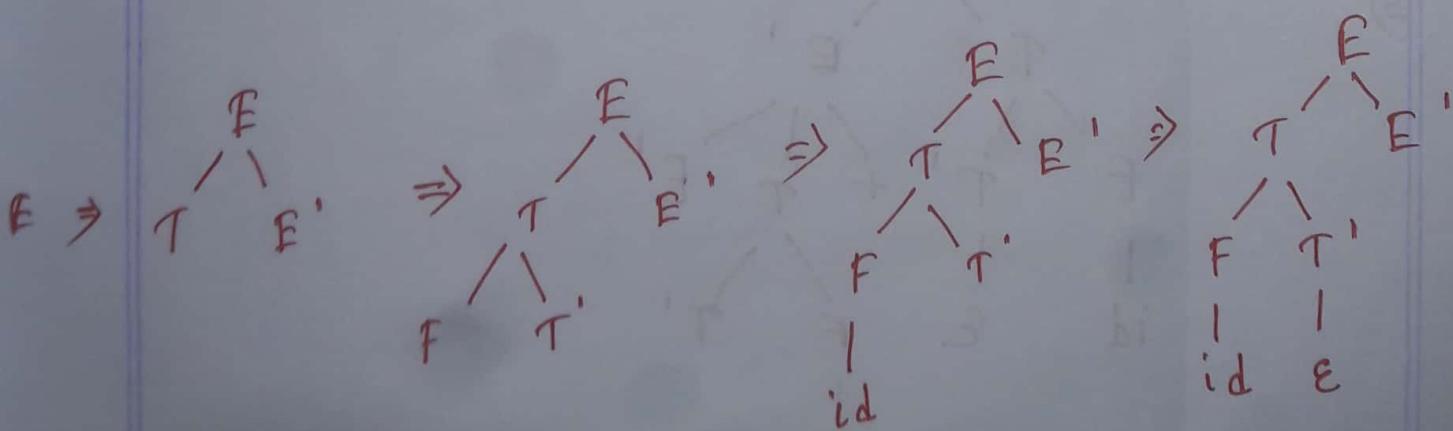
$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \epsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \epsilon \\ F &\rightarrow (\text{E}) \mid \text{id} \end{aligned}$$

Soln :- The grammar is already left factored.

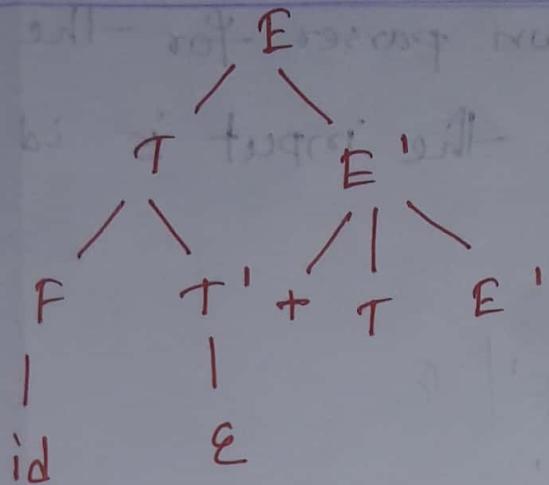
And no left recursion available here.

So, we can apply top-down parsing

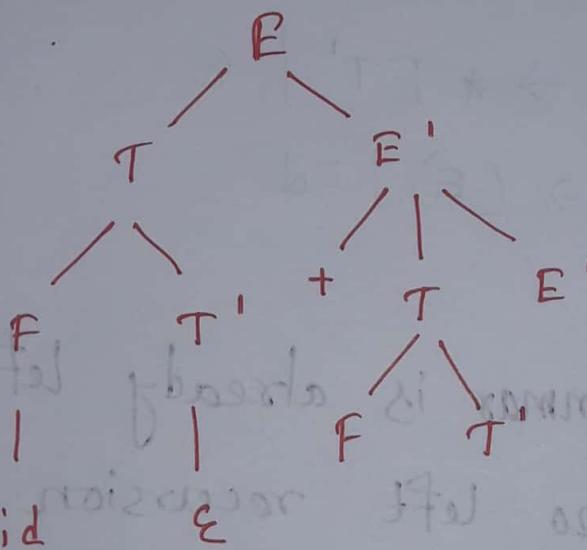
here.



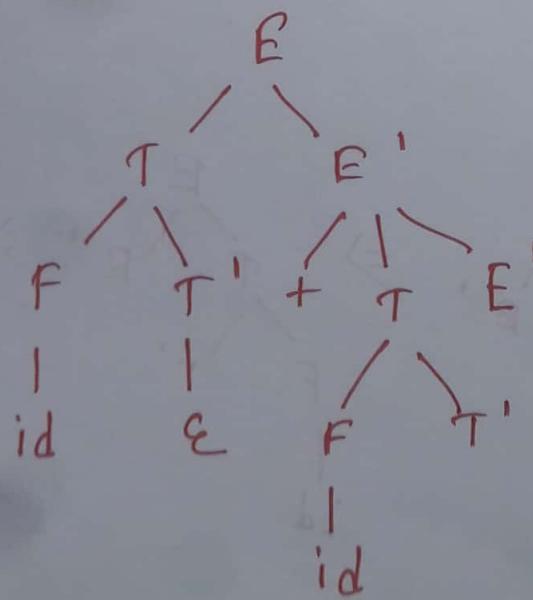
60



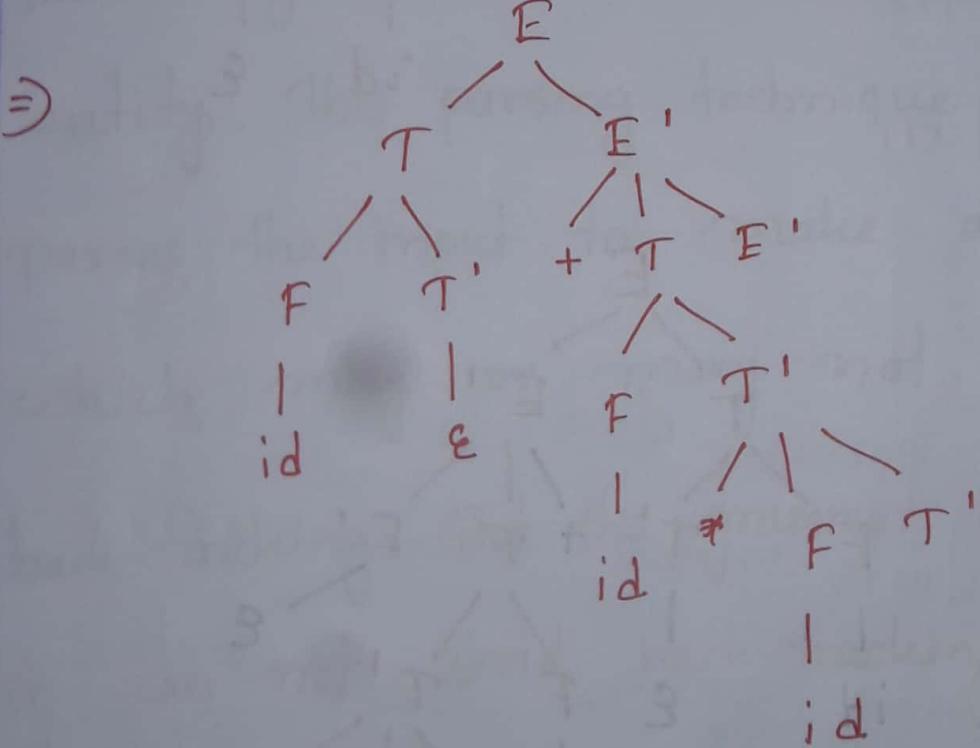
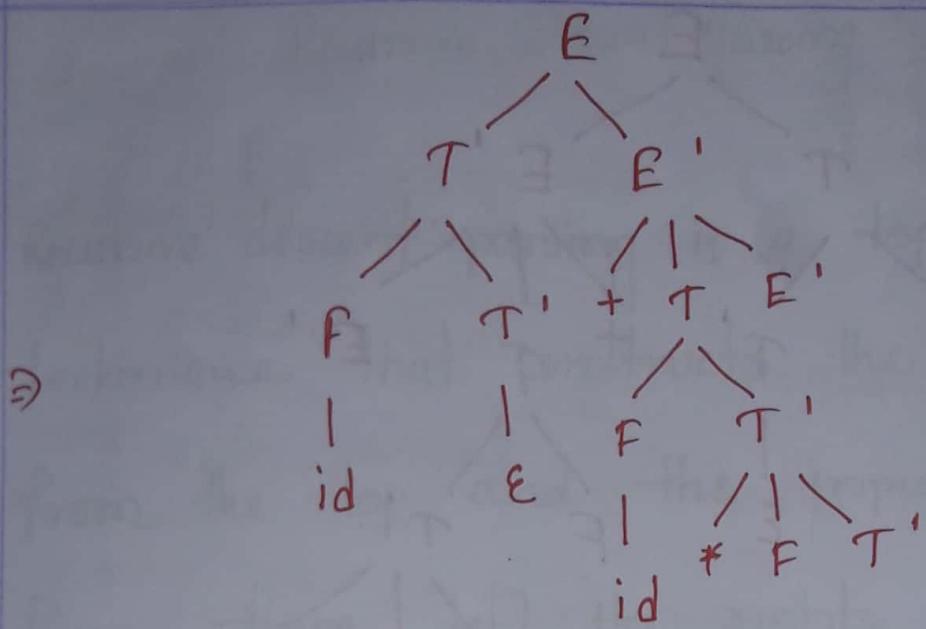
②



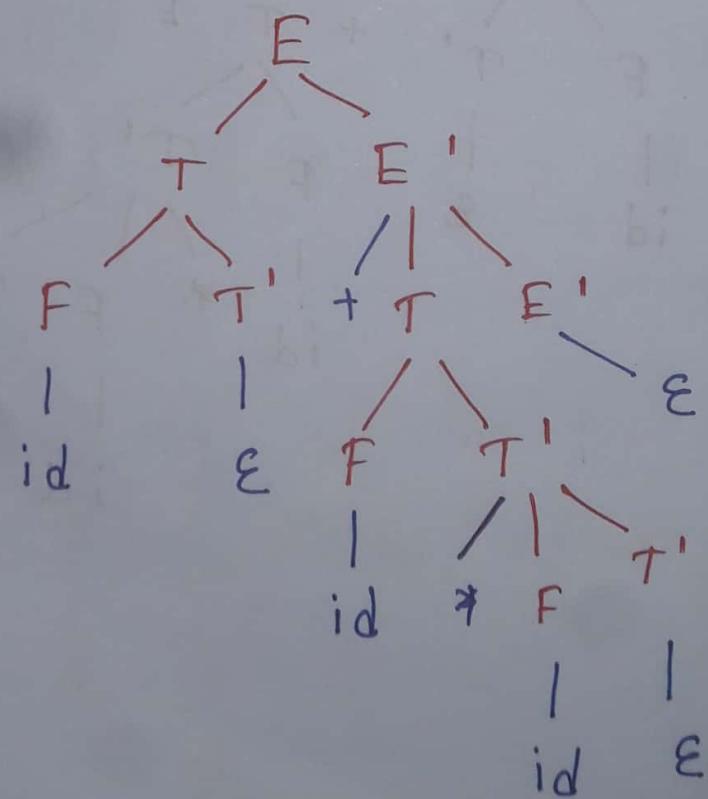
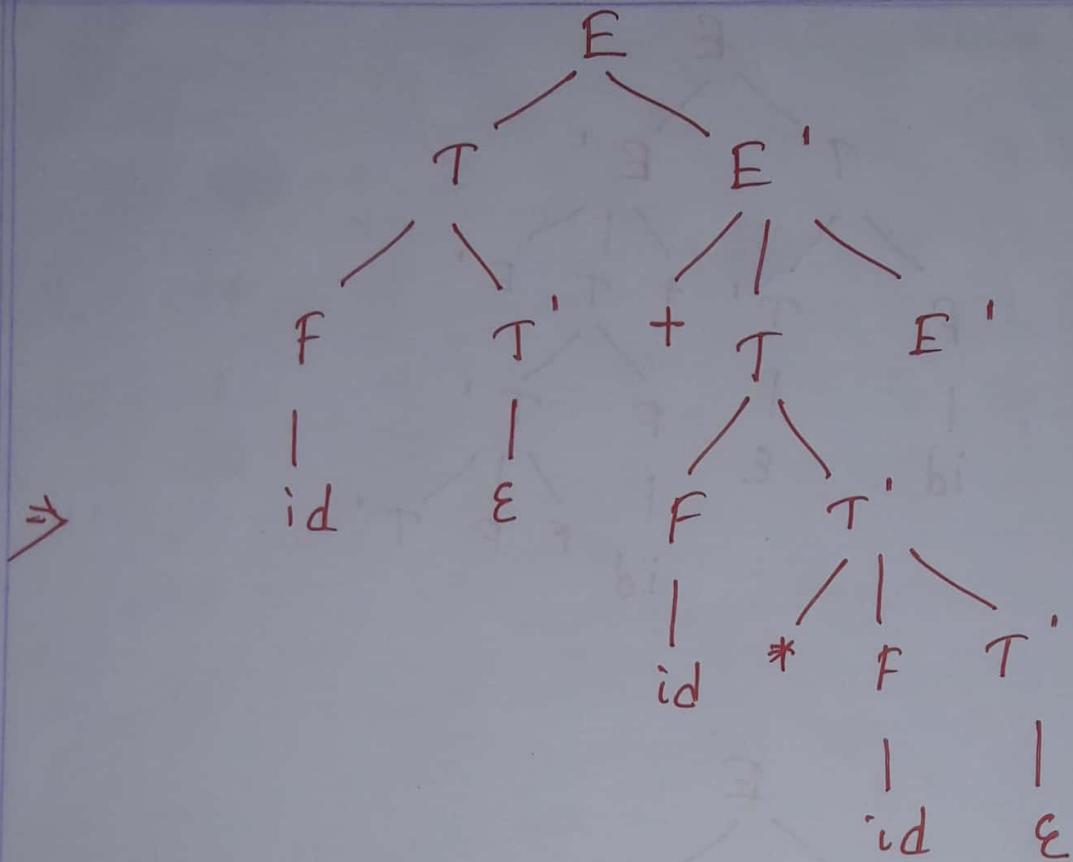
③



61



62



Recursive Descent Parsing

Recursive descent parsing is a top-down parsing technique that constructs the parse tree from the top and the input is read from left to right. It uses procedures for every terminal and non-terminal entity. This parsing technique recursively parses the input to make a parse tree, which may or may not require back-tracking. But the grammar associated with it can not avoid back-tracking (if not left factored).

Recursive Descent Parsing

Recursive descent parsing is a top-down parsing technique that constructs the parse tree from the top and the input is read from left to right. It uses procedures for every terminal and non-terminal entity. This parsing technique recursively parses the input to make a parse tree, which may or may not require back-tracking. But the grammar associated with it can not avoid back-tracking (if not left factored).

Predictive Parsing

A form of recursive-descent parsing that does not require any back-tracking is known as predictive parsing. This parsing technique is regarded as it uses context free grammar which is recursive in nature.

Backtracking

Top down parsers start from the root node (start symbol) and match the input string against the production rules to replace them (if matched).

65

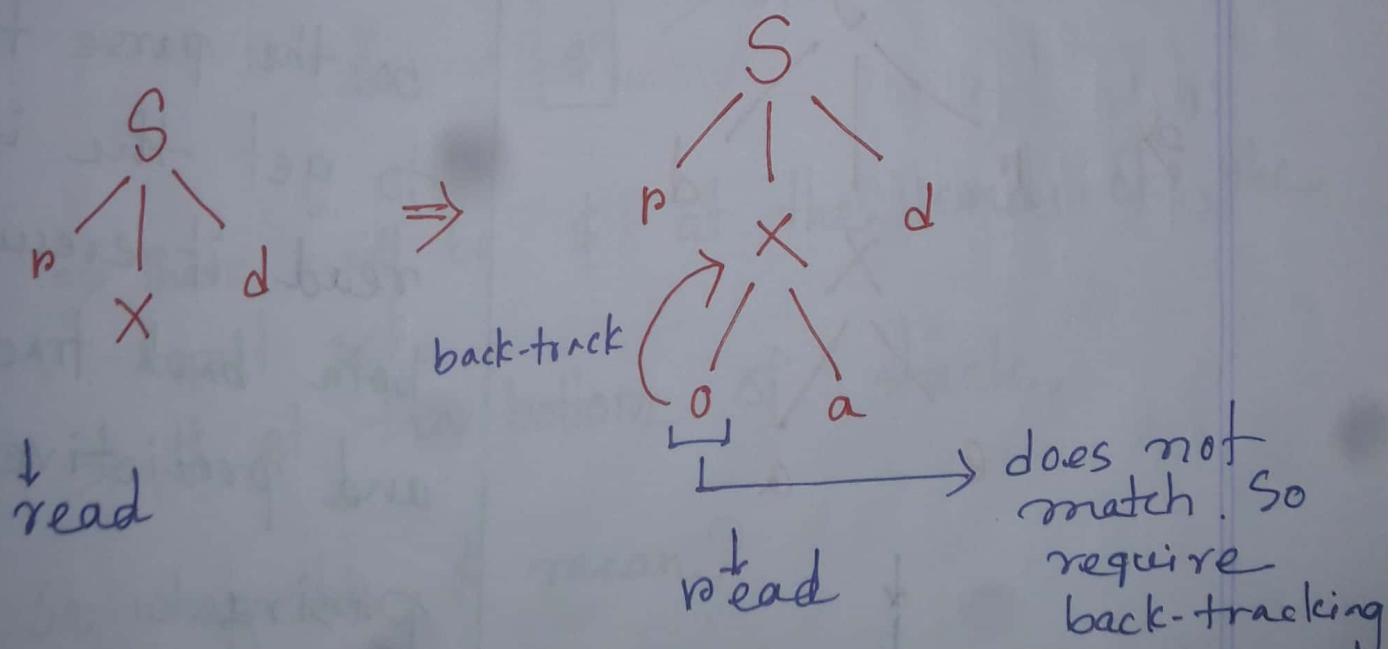
To understand this, take the following example of CFG.

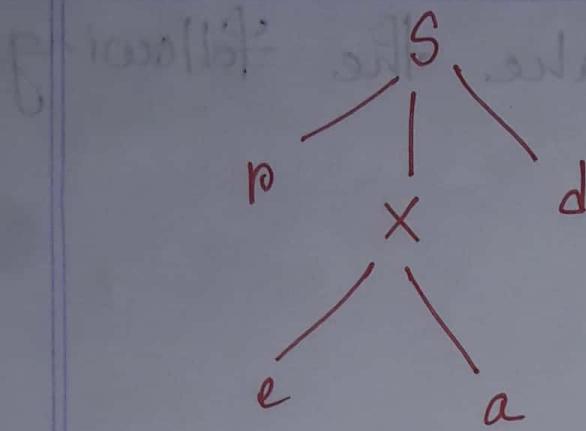
$$S \rightarrow rXd \mid r^2d$$

$$X \rightarrow oa \mid ea$$

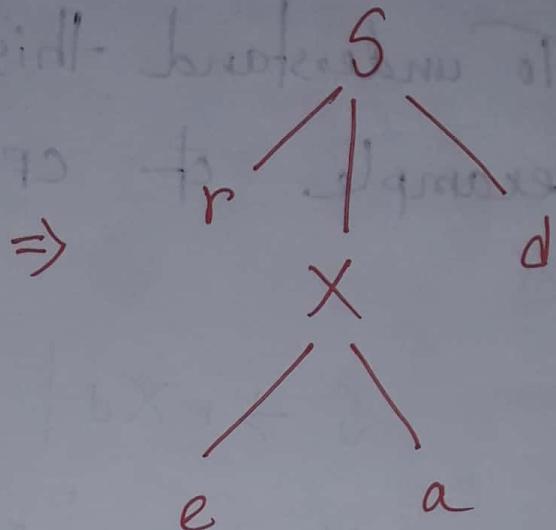
$$Z \rightarrow a^i$$

Consider input = read

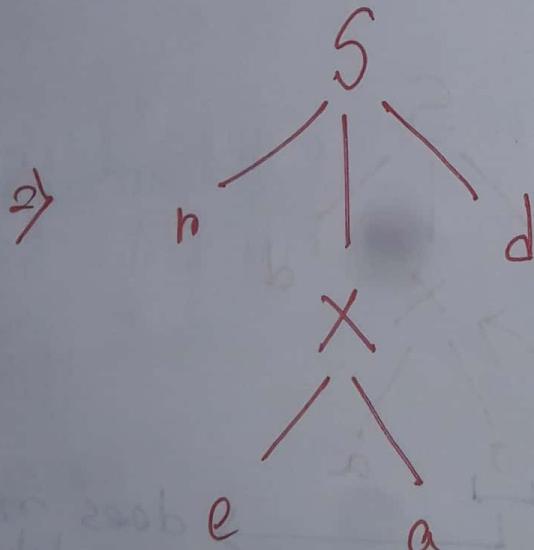




read



retad



read

So, the parse tree to get the input read, it execute both back-tracking and predictive parsing.

67

LL(1) Parser

L \Rightarrow Left to right

L \Rightarrow Left most derivation

(1) \Rightarrow No. of lookahead

how many symbol to see
to make a decision

input set : 

abcd

LL(1) parser :  \Rightarrow stack

\Rightarrow There will be \$ at the end of the input set.

\Rightarrow \$ at the bottom of stack.

So, what this \$ mean?

It means we are done with the input set also with the stack. Then we can generate a parsing table.

To generate the parsing table we will learn a function called FIRST & FOLLOW.

FIRST (A) \Rightarrow is a set of terminal symbols that begin after derivation from A .

FOLLOW (A) \Rightarrow is the terminal which could follow a variable in the process of derivation.

Rules to derive ~~FOLLOW~~:
FIRST

Rule 1 :- for a production rule $X \rightarrow \epsilon$,

$$\text{FIRST}(X) = \{\epsilon\}$$

Rule 2 :- for any terminal symbol $X \rightarrow a$

$$\text{FIRST}(x) = \{a\}$$

Rule 3 :- for any production rule $X \rightarrow Y_1 Y_2 Y_3$

$$\text{FIRST}(X)$$

\Rightarrow If $a \notin \text{FIRST}(Y_1)$, then $\text{FIRST}(X) = \text{FIRST}(Y_1)$

\Rightarrow If $a \in \text{FIRST}(Y_1)$, then $\text{FIRST}(X) =$

$$\{\text{FIRST}(Y_1) - \epsilon\} \cup \text{FIRST}(Y_2 Y_3)$$

$$\text{FIRST}(Y_2 Y_3)$$

\Rightarrow If $a \notin \text{FIRST}(Y_2)$, then $\text{FIRST}(Y_2 Y_3) = \text{FIRST}(Y_2)$

\Rightarrow If $a \in \text{FIRST}(Y_2)$, then $\text{FIRST}(Y_2 Y_3) =$

$$\{\text{FIRST}(Y_2) - \epsilon\} \cup \text{FIRST}(Y_3)$$

20

Similarly, we can make expansion for any production rule.

$$X \rightarrow Y_1 Y_2 Y_3 \dots Y_n$$

Rules to derive FOLLOW

Rule 01 :- for the start symbol S ,

place $\$$ in $\text{Follow}(S)$.

Rule 02 :- for any production rule, $A \rightarrow \alpha B$,

$$\text{Follow}(B) = \text{Follow}(A)$$

Rule 03 :- for any production rule, $A \rightarrow \alpha B \beta$.

\Rightarrow If $\epsilon \notin \text{FIRST}(\beta)$, then $\text{Follow}(B) = \text{FIRST}(\beta)$

\Rightarrow If $\epsilon \in \text{FIRST}(\beta)$, then $\text{Follow}(B) =$

$$\left\{ \text{FIRST}(\beta) - \epsilon \right\} \cup \text{Follow}(A) \quad \begin{array}{l} \text{[if } \beta \text{ is the} \\ \text{last element]} \end{array}$$

$\text{Follow}(A) = \text{FIRST}(\beta) - \epsilon \cup \text{Follow}(A)$ replace with ϵ [if β is not the last element]

Important notes:

1. ϵ may appear in the first function of a non-terminal.
2. ϵ will never appear in the follow function of a non-terminal.
3. Before calculating the first and follow functions, eliminate left recursion from the grammar, if present.
4. We calculate the follow function of a non-terminal by looking where it is present on the RHS of a production rule.

{ ϵ , $\$$ } = (S) reg

{ ϵ , $\$$ } = (T) reg

72

Example 01 :

Section 10.2

$S \rightarrow aBDh$

$B \rightarrow cC$

$C \rightarrow bC | \epsilon$

$D \rightarrow EE$

$E \rightarrow g | \epsilon$

$F \rightarrow f | \epsilon$

Soln

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

$$\text{FIRST}(B) = \{c\}$$

$$\text{FIRST}(C) = \{b, \epsilon\}$$

$$\text{FIRST}(D) = \{\text{FIRST}(E) - \epsilon\} \cup \text{FIRST}(F)$$

$$= \{g, f, \epsilon\}$$

$$\text{FIRST}(E) = \{g, \epsilon\}$$

$$\text{FIRST}(F) = \{f, \epsilon\}$$

X^B

$$\text{FOLLOW}(S) = \{ \$ \}$$

$$\text{FOLLOW}(B) = \text{FIRST}(A)$$

$$= \{ g, f, \epsilon \}$$

Here we got an empty. So, replace.

$$= \{ g, f, h \}$$

$$\text{FOLLOW}(D) = \{ h \}$$

$$\text{FOLLOW}(C) = \text{FOLLOW}(B)$$

$$= \{ g, f, h \}$$

$$\text{FOLLOW}(E) = \{ \text{FIRST}(F) - \epsilon \} \cup \text{FOLLOW}(D)$$

$$= \{ f, h \}$$

$$\text{FOLLOW}(F) = \text{FOLLOW}(D)$$

$$= \{ h \}$$

74

Problem 02

$$S \rightarrow A$$

$$A \rightarrow aB \quad | \quad Ad$$

$$B \rightarrow b \quad \{d, f, cf\}$$

$$C \rightarrow g \quad \{df\}$$

Sol2 :- The grammar is left recursive. So we have to eliminate it.

$$S \rightarrow A$$

$$A \rightarrow aBA'$$

$$A' \rightarrow dA' \quad | \quad c$$

$$B \rightarrow b \quad \{df\}$$

$$C \rightarrow g$$

15

$$\text{FIRST}(S) = \text{FIRST}(A)$$
$$= \{a\}$$

$$\text{FIRST}(A) = \{\wedge\}$$

$$\text{FIRST}(B) = \{b\}$$

$$\text{FIRST}(A') = \{d, e\}$$

$$\text{FIRST}(C) = \{g\}$$

FOLLOW

$$\text{FOLLOW}(S) = \{\$\}$$

$$\text{FOLLOW}(A) = \text{FOLLOW}(S)$$

$$\text{FOLLOW}(A') = \text{FOLLOW}(A)$$
$$= \{\$\}$$

$$\text{FOLLOW}(B) = \{\text{FIRST}(A') - \{\$\}\} \cup \{\text{FOLLOW}(A)\}$$
$$= \{d, \$\}$$

$$\text{FOLLOW}(C) = \text{NAR}_{\{d, \$\}}$$

16

Problem 03

$$S \rightarrow (L) \mid ^*$$

$$L \rightarrow SL'$$

$$L' \rightarrow SL' \mid \epsilon$$

Soln

$$\text{FIRST}(S) = \{c, a\}$$

$$\text{FIRST}(L) = \{ \text{FIRST}(S) \} = \{c, a\}$$

$$\text{FIRST}(L') = \{\cdot, \epsilon\}$$

FOLLOW

$$\text{FOLLOW}(S) = \{\$\} \cup \text{FIRST}(L')$$

Here, $\text{FIRST}(L')$ contains

ϵ so,

$$= \{\$\} \cup \{\text{FIRST}(L') - \epsilon\} \cup \{\text{FOLLOW}(L')\}$$

$$= \{ \$, , ,) \}$$

$$\text{FOLLOW}(L) = \{ > \}$$

$$\begin{aligned} \text{FOLLOW}(L') &= \text{FOLLOW}(L) \\ &= \{ > \} \end{aligned}$$

Problem 03 :

$$S \rightarrow AaAb \mid BbBa$$

$$A \rightarrow \epsilon$$

$$B \rightarrow \epsilon$$

Soln :-

$$\text{FIRST}(S) = \text{FIRST}(A) \cup \text{FIRST}(B)$$

Both both $\text{FIRST}(A)$ and $\text{FIRST}(B)$

contain ϵ . So,

$$= \{\text{FIRST}(A - \epsilon)\} \cup \{a\} \cup \{\text{FIRST}(B - \epsilon)\} \cup$$

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

X9

$$\text{FIRST}(A) = \epsilon \quad \{ \dots \} =$$

$$\text{FIRST}(B) = \epsilon$$

$$\{ \dots \} = C \text{ word}$$

FOLLOW

$$\text{FOLLOW}(S) = \{ \$ \} = C \text{ word}$$

$$\text{FOLLOW}(A) = \text{FIRST}\{a, b\}$$

$$\text{FOLLOW}(B) = \{a, b\}$$

: so make it

$$ABP/BAA \leftarrow A$$

$$e \in A$$

$$e \in A$$

A

(A)

$$(A) T2917 \cup (A) T2917 = (A) T2917$$

(A) T2917 for (A) T2917 add for

(A) T2917 for

$$v \{ (A - e) T2917 \} \cup \{ v \} \{ (A - e) T2917 \} =$$

$$\{ d, e, f \} =$$

X9

LL(1) Parsing Table

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \epsilon$$

$$F \rightarrow id \mid (E)$$

Soln: First of all, we will find the FIRST & FOLLOW of this grammar.

Productions		FIRST	FOLLOW
$E \rightarrow TE'$	E	{id, (}	{\$,)}
$E' \rightarrow +TE' \mid \epsilon$	E'	{+, ε}	{\$,)}
$T \rightarrow FT'$	T	{id, (}	{\$,), +}
$T' \rightarrow *FT' \mid \epsilon$	T'	{*, ε}	{\$,), +}
$F \rightarrow id \mid (E)$	F	{id, (}	{\$,), +, *}

	id	+	*	()	\$
E	$E \rightarrow TE'$				$E \rightarrow TE'$	
E'		$E' \rightarrow +TE'$				$E' \rightarrow \epsilon$
T	$T \neq FT'$				$T \rightarrow FT'$	
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$				$F \rightarrow (E)$	

If the grammar has left recursion eliminate it.

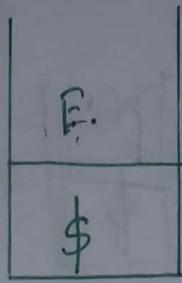
If the grammar has non determinism, apply left factoring to make it deterministic. Then we can make $LL(1)$ parsing table. But not guaranteed that we can generate parse tree.

* If same cell has two productions then $LL(1)$ parse tree impossible

Q

Input : id + id \$

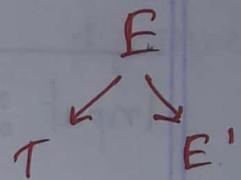
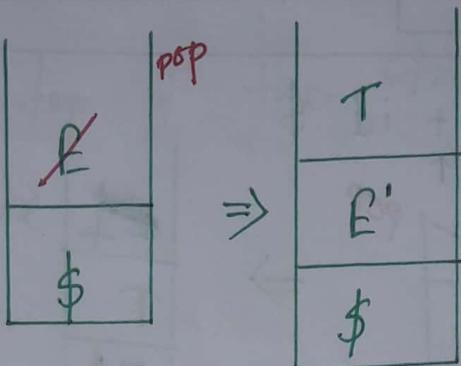
Stack :



E.

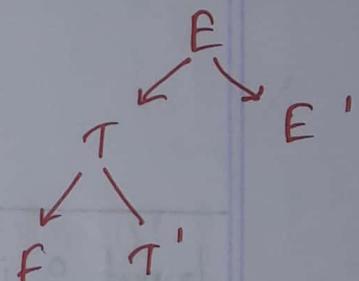
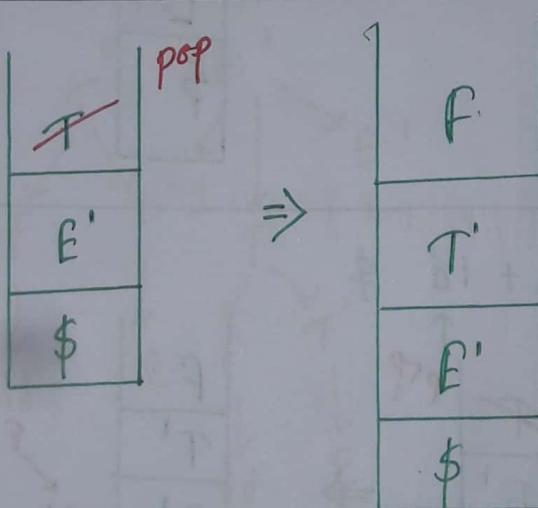
Input : id + id \$

Stack :



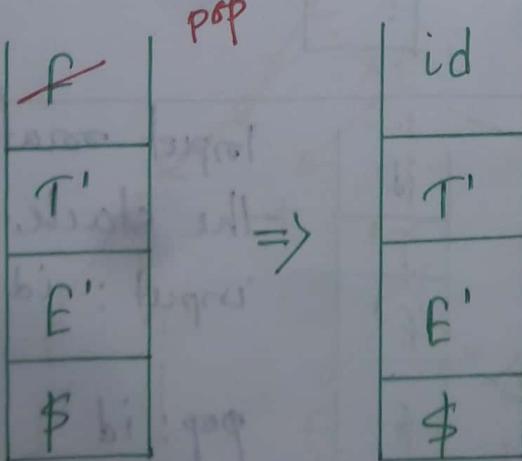
Input : id + id \$

Stack :



Input : id + id \$

Stack :

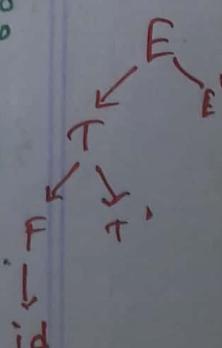


top of the stack & input matches. So, increment the pointer & the input:

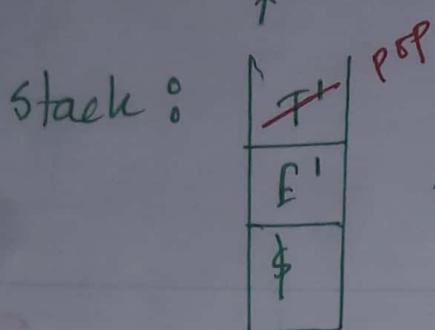
id + id \$

↑

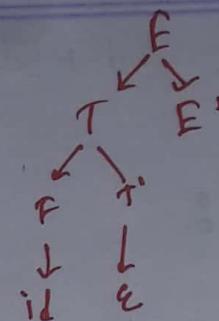
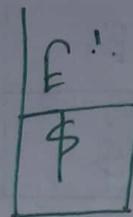
And pop id.



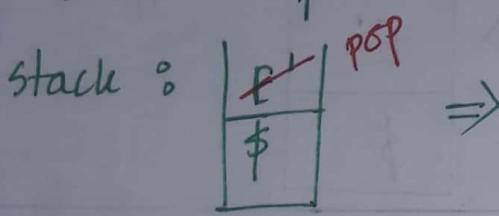
Input : id + id \$



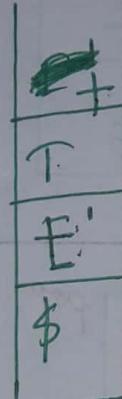
\Rightarrow



Input : id + id \$

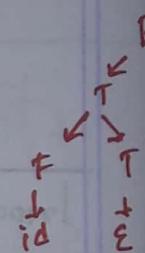


\Rightarrow



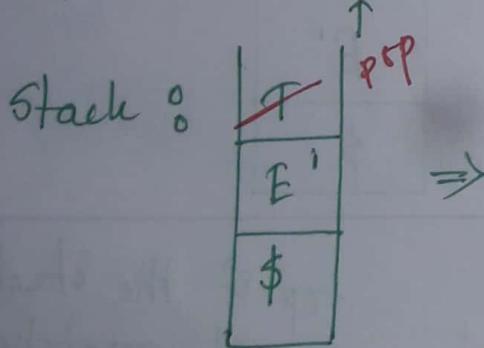
Input matches
top of the stack.

So, input:
id + id \$

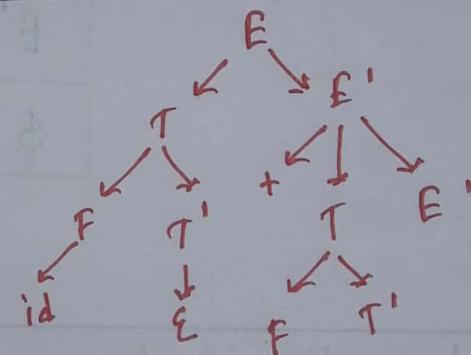
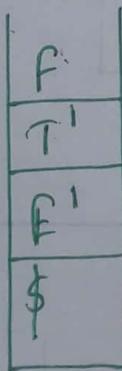


POP +

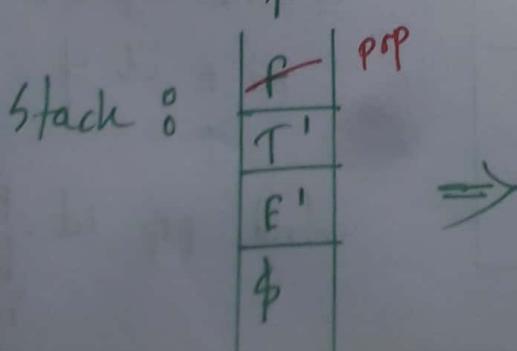
Input : id + id \$



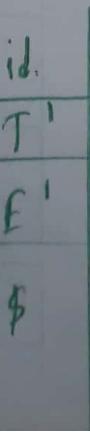
\Rightarrow



Input : id + id \$



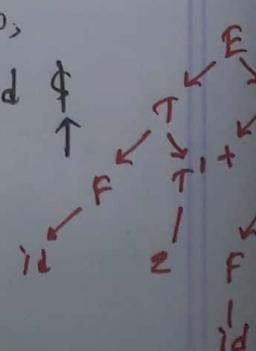
\Rightarrow



Input matches top of
the stack. So,

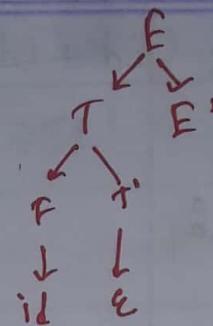
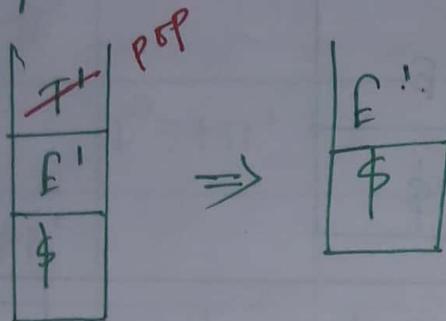
input : id + id \$

pop : id



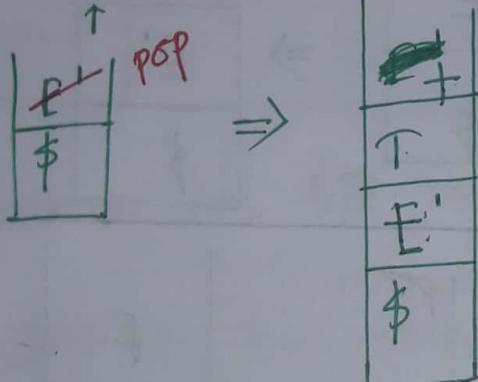
Input : id + id \$

Stack :



Input : id + id \$

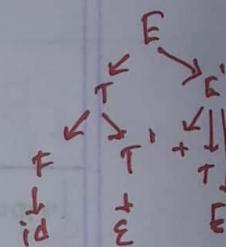
Stack :



Input matches
top of the stack.

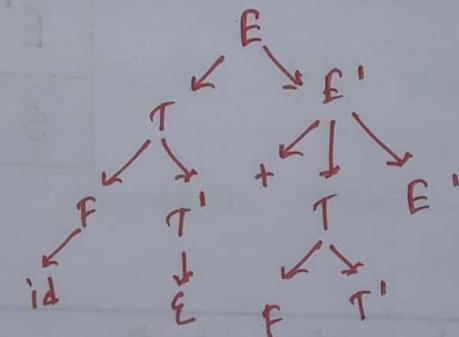
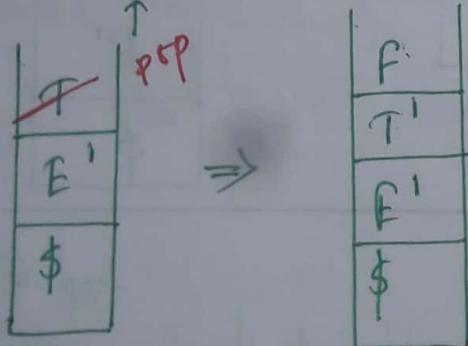
So, input:
id + id \$

pop +



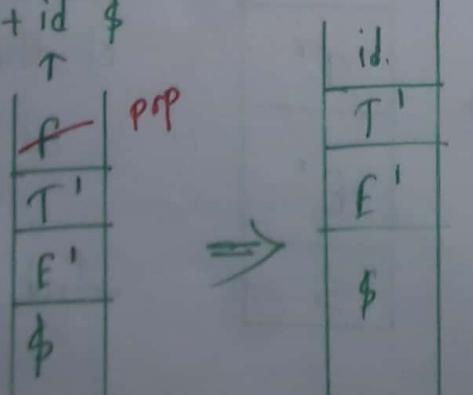
Input : id + id \$

Stack :



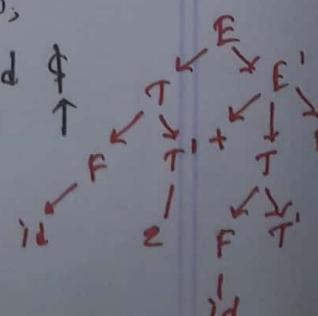
Input : id + id \$

Stack :



Input matches top of
the stack. So,

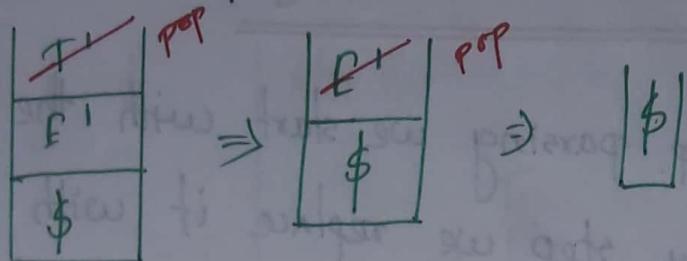
input : id + id \$
pop: id



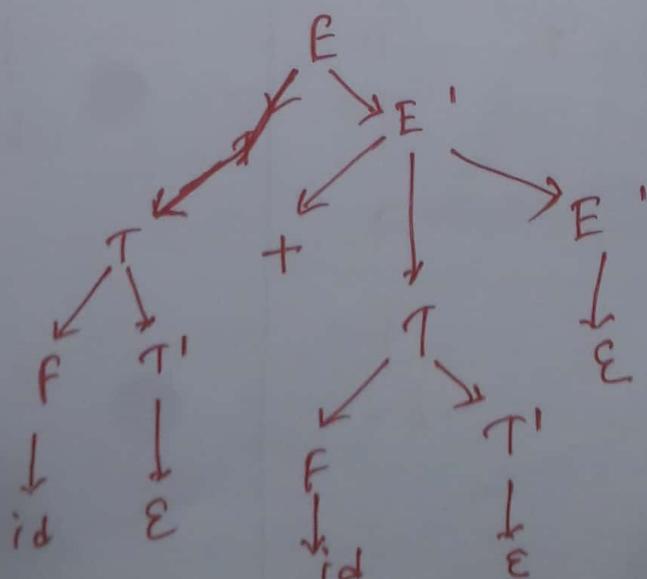
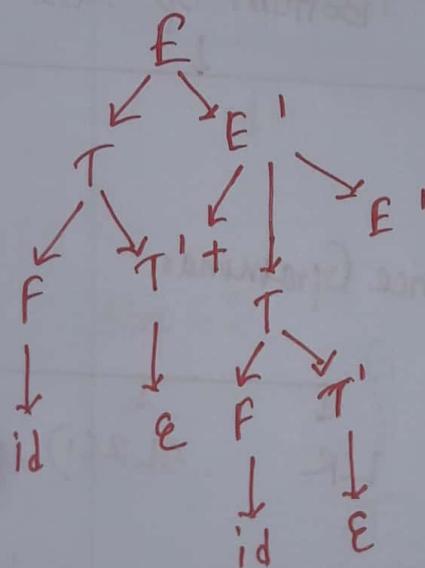
Q3

Input : id + id $\not\models$

Stack :

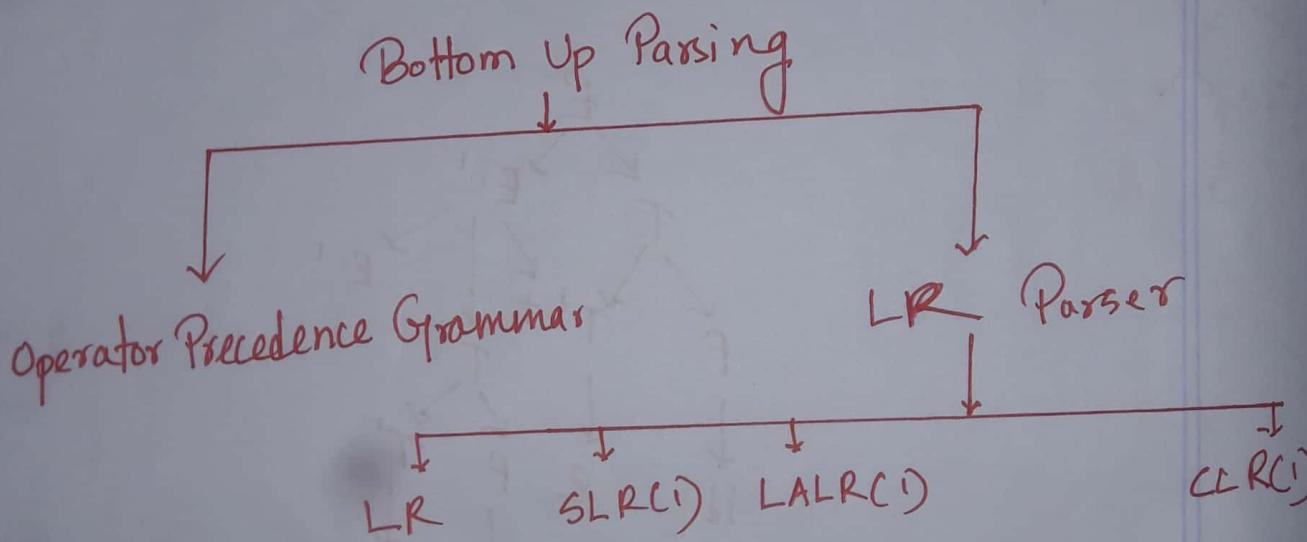


Input matches top of the string. $\not\models$ meets $\not\models$. So, terminate.



Bottom Up Parsing

In bottom up parsing we start with the string symbol and step by step we replace it with its production. Ultimately we end up with the starting state. It is also known as Shift Reducing Parser.

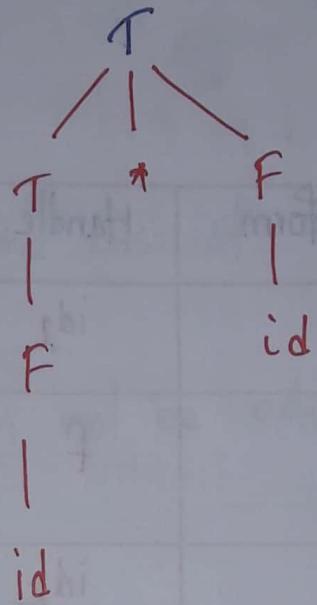


Input string: id * id

$$\begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T * F \mid F \\ F \rightarrow (E) \mid id \end{array}$$

Step 1: id * id

Step 5:

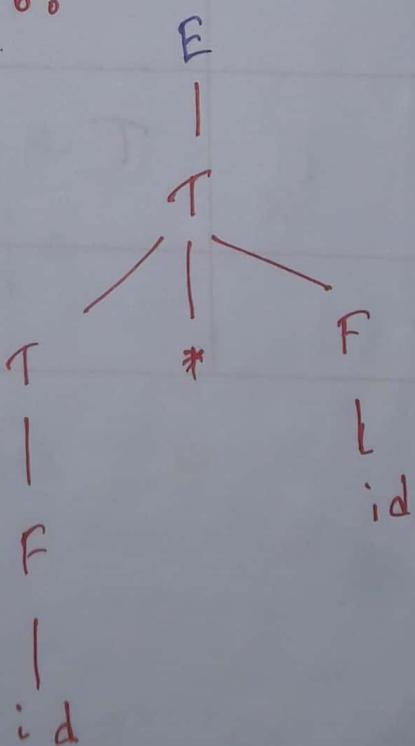


Step 2: F * id
|
id

Step 3: T * id
|
F
|
id

Step 4: T * F
|
F
|
id

Step 6:



Q6P

Handle Pruning

$id_1 * id_2$

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

Right sentential form	Handle	Reducing Production
$id_1 * id_2$	id_1	$F \rightarrow id$
$F * id_2$	f	$T \rightarrow F$
$T * id_2$	id_2	$F \rightarrow id$
$T * F$	$T * F$	$T \rightarrow T * F$
T	T	$E \rightarrow T$
E		

Operator Grammar

Operator grammar is used to define mathematical operators.

Two properties :

- 1 ϵ can not be in the RHS of any production.
- 2 Two nonterminals can not be adjacent to each other.

Examples :

$$E \rightarrow E+E \mid E \cdot E \mid id$$

$$\begin{aligned} x &\rightarrow ABX \mid a \\ B &\rightarrow bXb \mid a \end{aligned}$$

But

$$x \rightarrow AbXbXa \mid Aax$$

ofb

Operator Precedence Relation

3 types of relation.

1 $a > b$ [a has higher precedence than the terminal b]

2 $a < b$ [b has higher precedence than the terminal a]

3 $a \doteq b$ [a and b both have same precedence]

Operator Precedence Table

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

	id	+	*	\$
id	-	>	>	>
+	<	>	<	>
*	<	>	>	>
\$	<	<	<	-

string: id + id * id

- Here we will use a stack.
- We will add \$ at the bottom of the stack as well as at the end of the input string.
- We will compare the value of stack (top) with the input (cursor)
- Stack values are the values from column (left)
- Input values are the values from row (top)

Comparison:

- If the values of stack are smaller than input string then,
 - * push the value into the stack
 - * increment the cursor of input string
- If the values of stack are greater than input string then
 - * pop from the stack
 - * reduce the tree

Step 1:

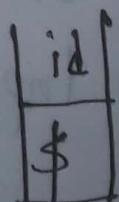
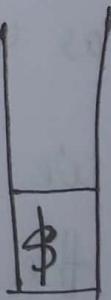
$\underset{\uparrow}{\text{id}} + \text{id} * \text{id} \$$

Compare $(\$, \text{id}) = \$$ smaller so

push the id into the stack &

increment the lookahead/cursor

$\underset{\uparrow}{\text{id}} + \text{id} * \text{id} \$$



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

Step 2: $id + id * id \$$

id
\$

Compare ($id, +$) = id greater so,

pop id from the stack and reduce
the grammar

\$

E
|

\$ bi + bi + bi \$

$id + id * id \$$

↑

Step 3:

Compare ($\$, +$) = $\$$ smaller so push $\$ +$
and increment the cursor

\$
+

E
|

$id + id * id \$$

↑

$+ = (\$, +)$

✓

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

Step 4: Compare $(+, id)$ = + smaller
so push id into the stack & increment the lookahead

id
+
*

E

|
id + id * id \$

E

|
id + id * id \$

+
*

Step 5: Compare $(id, *)$ = id greater
so pop id and reduce the grammar

E

|
id + id * id \$

\$ bi * bi + bi

\$ bi * bi + bi

Step 6: Compare $(+, *)$ = + smaller so push * and increment the cursor.

99
 02 ~~return~~ $E \vdash (\$ + \text{id}) \text{ original} : e$
~~cursor here~~ $\boxed{*}$
 $\boxed{+}$
 $\boxed{\$}$

id + id * id \$
 ↑

Step 7: Compare (id , $\$$) = ~~id~~ * smaller so
 push id & increment cursor

02 ~~return~~ $E \vdash (\$ + \text{id}) \text{ original} : e$
~~cursor here~~ \boxed{id}
 $\boxed{*}$
 $\boxed{+}$
 $\boxed{\$}$

id + id * id \$
 ↑

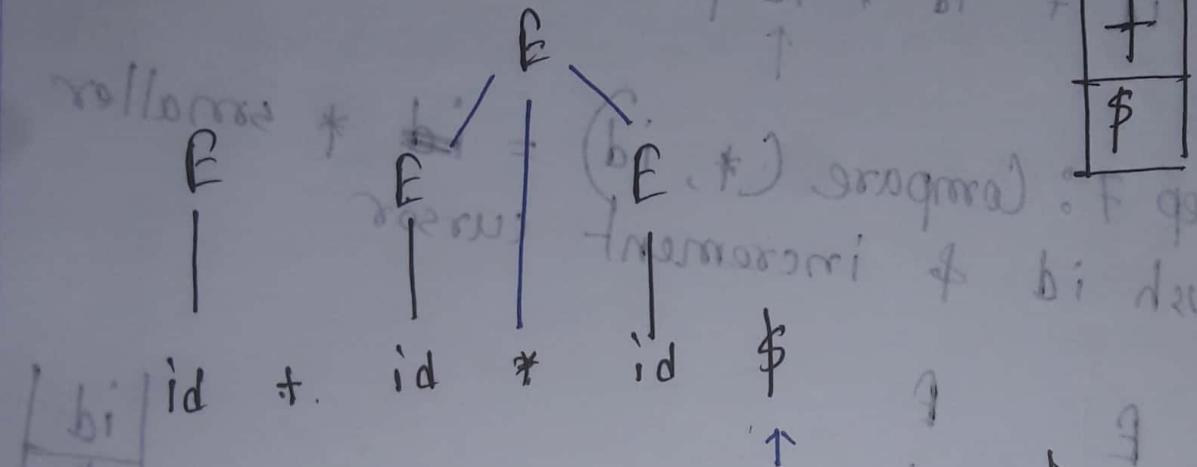
Step 8: Compare (id , $\$$) = id greater so
 pop and reduce

id + id * id \$
 ↑

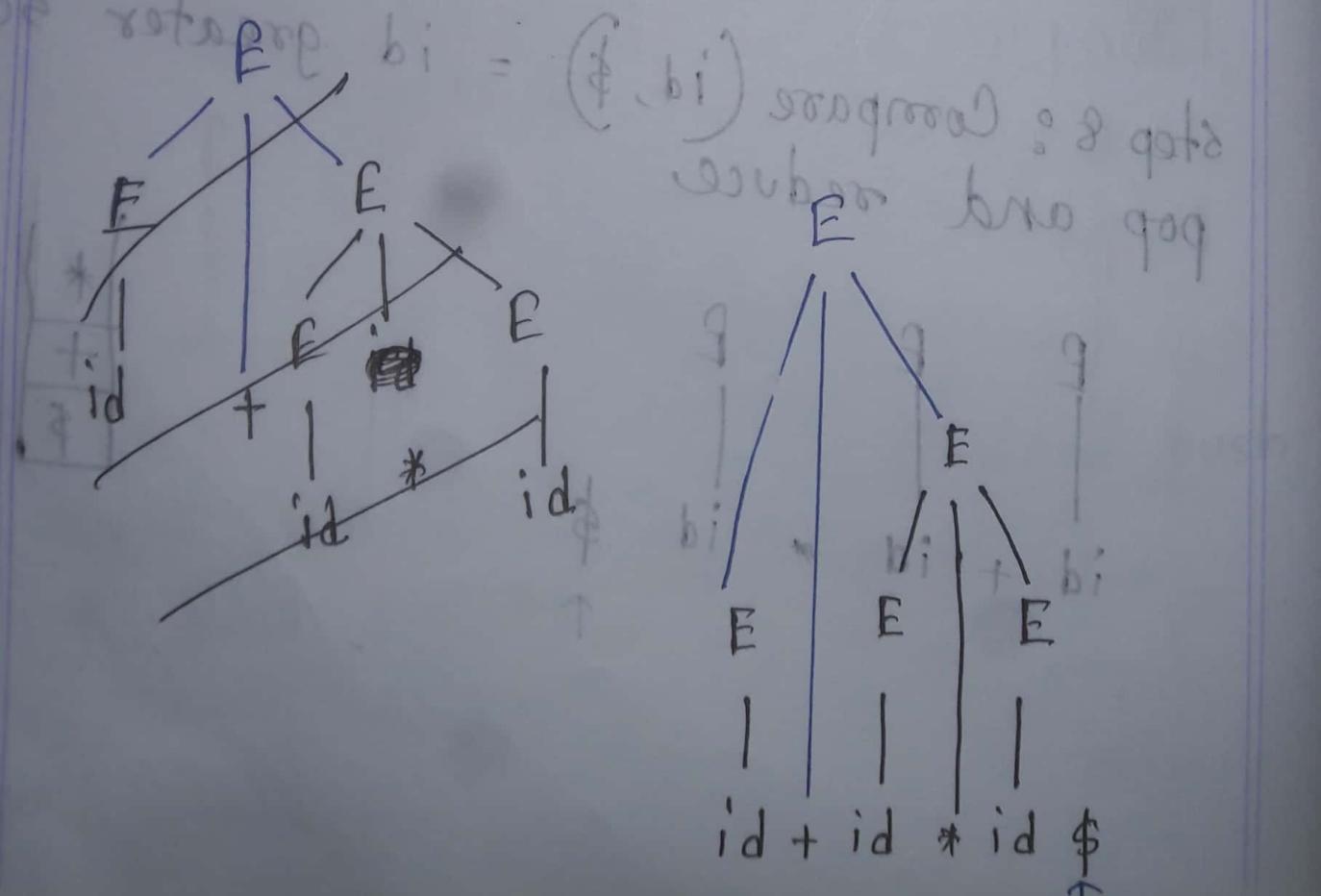
\$ $\boxed{*}$
 $\boxed{+}$
 $\boxed{\$}$

9

step 9 : Compare $(*, \$)$ = * greater so
pop and reduce



step 10 : Compare $(+, \$)$ = + greater so
pop and reduce



So, ϕ meets ψ and we got our desired output.