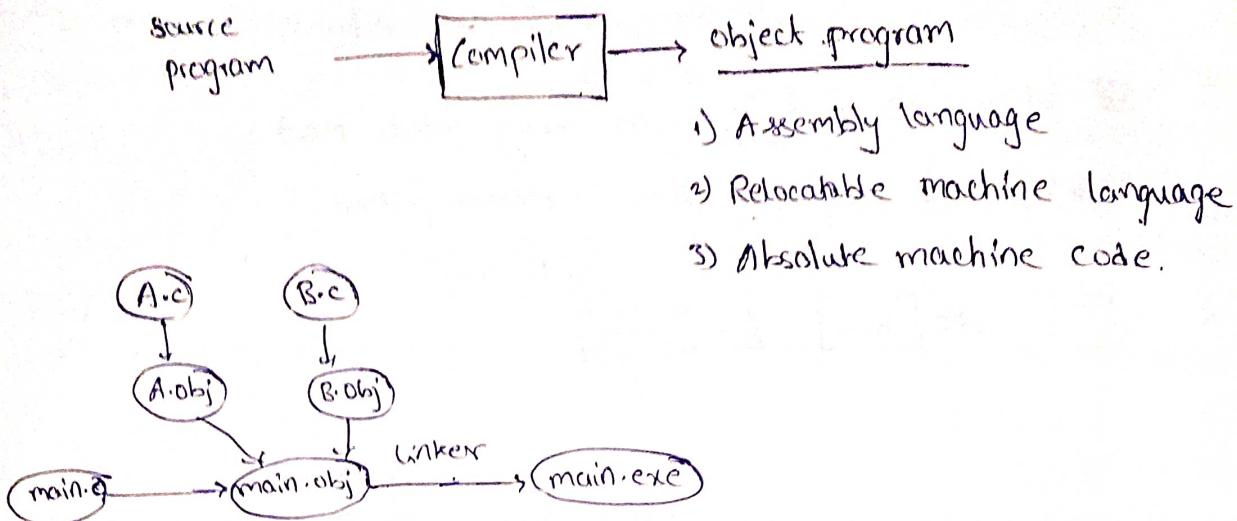


16/06/20

Compiler Design

166

Compiler: It is a language translator which translates from one language to others.



Contents:

- 1) Lexical Analysis
- 2) Parsing
- 3) S.D.T { 70%
- 4) Intermediate code generation

4 marks:

2 - 1 marks

1 - 2 marks

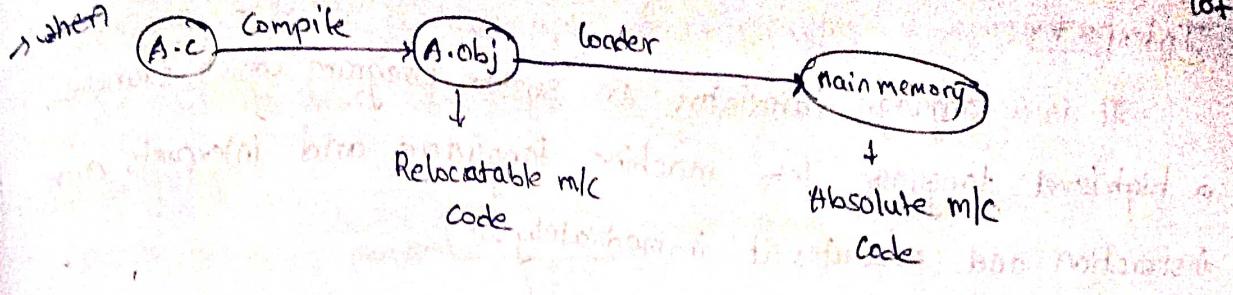
Reference: Ullman

Introduction:

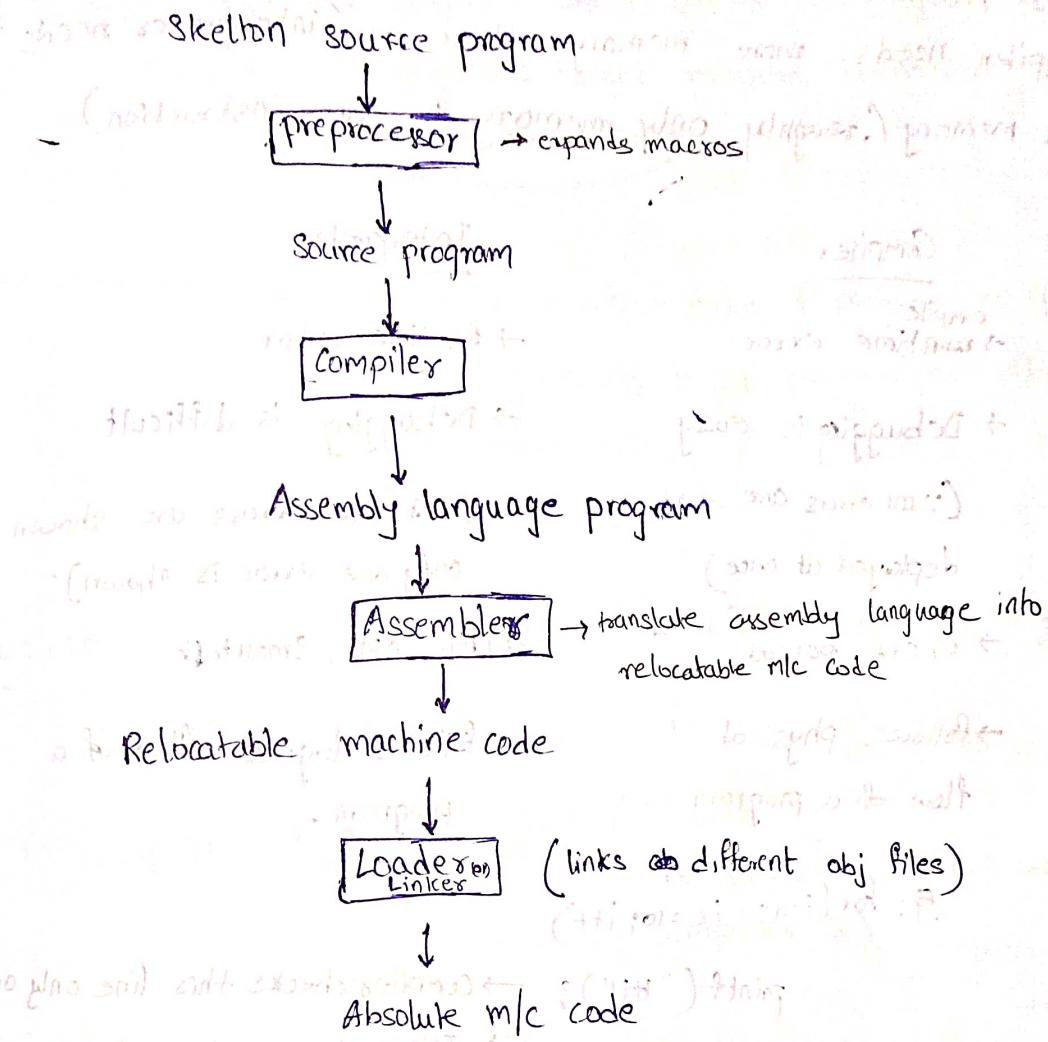
Compiler: It is a language processor which translates source program into object program.

Object program may be any of below :-

- i) Assembly language
- ii) Relocatable machine code
- iii) Absolute machine code



Language Processing system:



Linker: It makes a single program from several files of relocatable machine code.

Loader: It takes relocatable machine code

and altering the relocatable addresses and placing code & data into main memory. This process is called loading.

Interpreter:

It is a language translator or system program which translates a high-level language into machine language and interprets each instruction and executes it immediately.

→ Execution of compiled languages is faster than interpreted languages.

→ Compiler needs more memory whereas interpreter needs less memory (roughly only memory for one instruction)

Compiler

→ ~~Runtime error~~
compile time error

→ Debugging is easy

(as all errors are displayed at once)

→ C, C++, Pascal

→ follows physical flow of a program

Eg: `for(i=1; i<=10; i++)`

`printf("Hi");` → Compiler checks this line only once
 ~~stays~~ → Interpreter considers the line 10 times

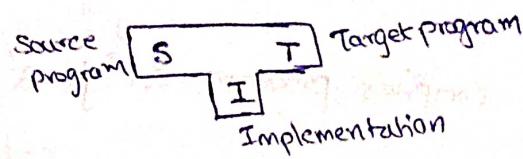
→ Interpreters are mainly used when size of main memory is less

(so that we can execute each instruction line by line)

→ All web based programming languages are interpreted programming languages because they need to be run different devices with different memory.

104

Bootstrapping: It is a concept of compiling a compiler program in its own language to generate the compiler for a new language.

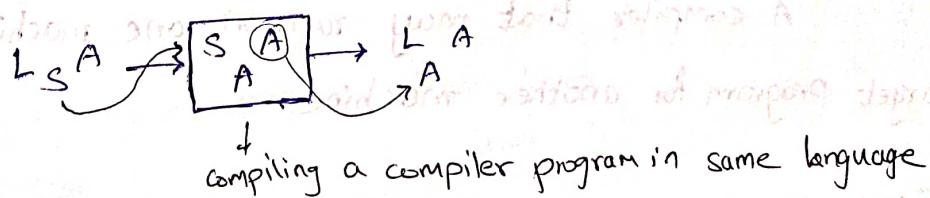


Existing compiler: $S_A \xrightarrow{A}$ compiler which takes 'S' program and generates target program for m/c A, which runs on A.

Required compiler: $L_A \xrightarrow{A}$ compiler which takes L program as input and generates target program for m/c A and which runs on A.

Step 1: Write a compiler program for 'L' language in language 'S' targeted for m/c A.

Step 2: Compile L_A in S_A ($\because L$ is written in S to compile L we use compiler S)



Eg: Existing: C win Compiler for 'C' for generating target program for windows and it runs on windows

Required: Consider we need a compiler for java

i.e., java win
 win

Step 1: Write a compiler program for "java" using 'c' targeted for windows

java win

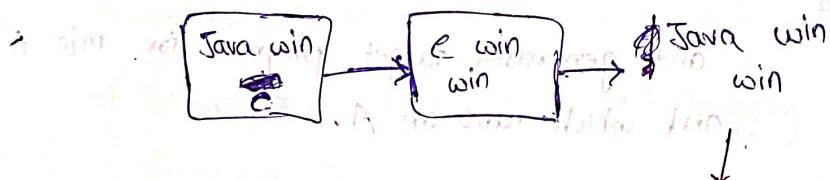
C

Step 2: Compile

Java win using C win

C

win



Java compiler
targeting windows and
running on windows

Self-hosting Compilers:

It is compiler which is capable of compiling its own compiler program.

Cross compilers:

A compiler that may run on one machine and generates target program for another machine.

Existing: L_AA (compiler for language L, targeted for m/c 'A' and runs on machine A)

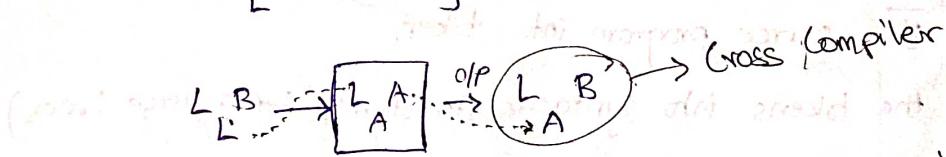
Required: L_BB (compiler for language L, targeted for m/c 'B' and runs on machine B)

Step 1: write compiler program for ~~language 'L'~~ 'L',
targeted for machine B using L.

i.e., $L_L B$

Step 2: Compiler

$L_L B$ using $L_A A$

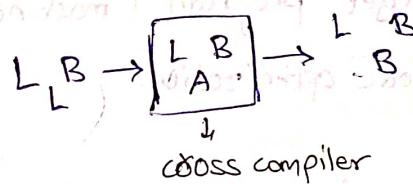


↳ Compiler for language 'L'

targeted for m/c 'B'
and runs on m/c 'A'

Set Step 3:

Compile $L_L B$ into $L_A B$



cross compiler

above, entire process is known as cross compilation.

Existing

Java win
\\ win

Req: Java linux (figured out already)
linux

Step 1: write a compiler program

Java
linux

Step 2:

Java
linux
Java
→ Java
win
win

Java
linux
win

↳ cross compiler

Step 3:

Java
linux
win
Java
linux
Java
win
win

Java
linux
win

Java
linux
linux

Compilation Process:

It is divided into 2 parts

i) Analysis part

ii) Synthesis part

Analysis part

(i) divide the source program into tokens

(ii) group the tokens into syntactic structures (~~parse trees~~)

(iii) verify the semantics

(iv) generate intermediate code

(v) M/C independent code optimization

Synthesis part:

1) generate the target program (machine dependent)

2) M/C dependent code optimization.

→ Each task in analysis and synthesis part is represented by phase.

phase: It is a logically cohesive operation

i.e., logically interlinked operation.

Phases in Compilation:

1) Lexical Analysis

2) Syntax Analysis

3) Semantic Analysis

4) Intermediate Code generation

5) M/C independent optimization (only if required)

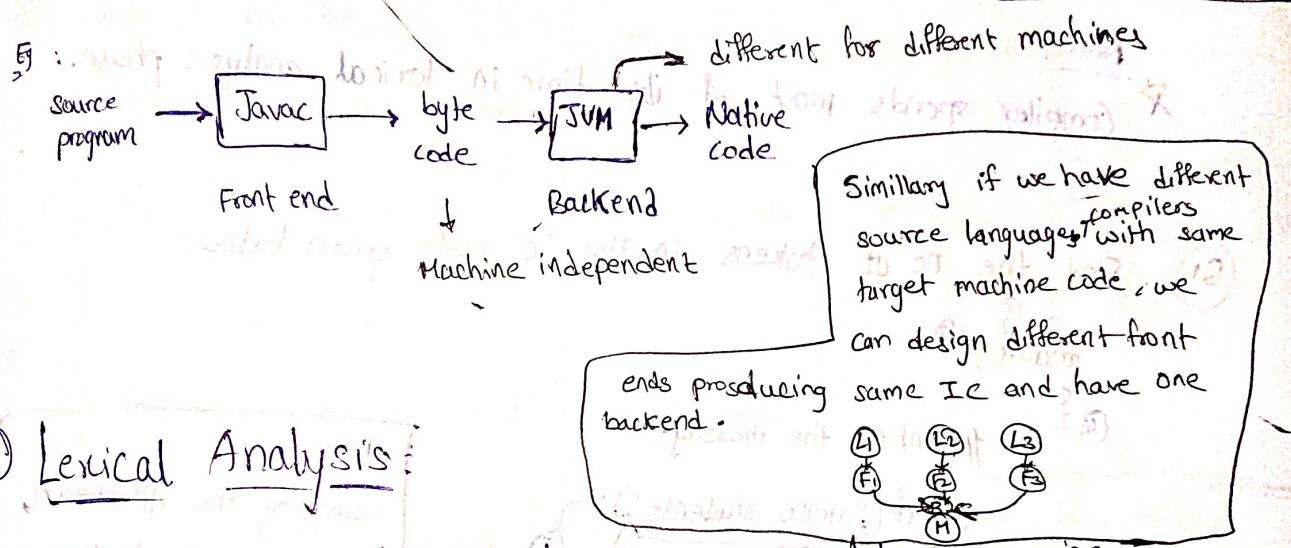
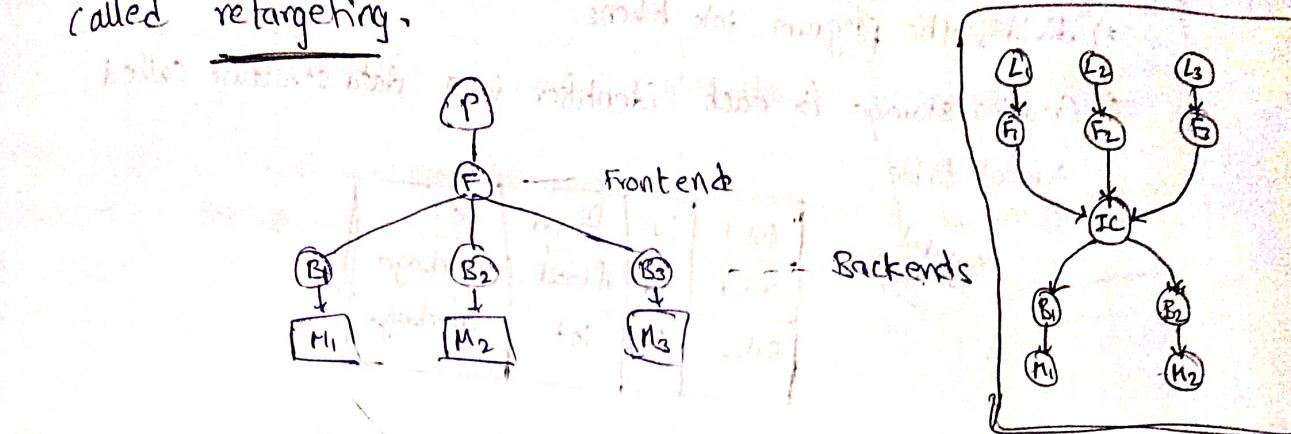
6) Code generation

7) M/C dependent code optimization

} front end

} back end

Consider we need to make a compiler of language 'P' for 3 machines M_1, M_2, M_3 . We can design one front end and design 3 backends for each machine. This process is called retargeting.



1) Lexical Analysis:

* also known as linear analysis, scanner, token recognizer

* lexical analysis is a phase in which stream of character

making the source program is read and grouped into tokens

that are sequences of characters having a collective meaning

such as identifiers, keywords, constants, special characters etc.

$$2): \text{Total} = \text{sub1} + \text{sub2} * 2.0$$

Butter

total ← identifizier

= < operator

in StudienfA. Zul. Subj → identifizieren

~~function~~ ~~if~~ ~~using~~ ~~the~~ ~~operator~~

`Sub2 ← identifier`

* ↑ open to

2.0 ← number

; \leftarrow special character

function of LA:

1) Removes whitespaces and comments

(If there are more than 1 consecutive whitespaces, then it is reduced to one whitespace)

2) In \hookrightarrow single ~~space~~ spaces are not removed

3) divides the program into tokens.

4) Creating storage for each identifier in a data structure called symbol table.

Symbol table:

total	id	float	0.0
sub1	id	float	garbage
sub2	id	int	garbage

Note:

* Compiler spends most of its time in lexical analysis phase.

Q1

Find the no of tokens in the 'c' code given below:

main() {

 printf("Hello Students");

} // printing the message

string constant

∴ 10 tokens

Q2

Find no of tokens in the below program

while(i > 0)

{

 printf("%d", i);

 i++;

}

∴ 18 tokens

Note:

Scanning the I/O stream,

(source program) the

LA groups ~~the~~ the characters

into meaningful sequences

called lexemes, For

each lexeme, a token

is produced as an o/p.

The form of token is

<token-name, attribute-value>

token-name is an abstract

symbol, that is used in

Syntax analysis. Attribute-value

points to an entry in symbol

table for this token.

Q3 which of the following statements can be recognized as a token without depending on next symbol

- a) + b) < c) >> d) none

Sol:

for +, + is (b) + may be possible

for <, < may be possible,

for >>, >> may be possible

; none.

The process of lexical analysis is done by input buffering scheme

Q4 which of the following can be recognized as a token without depending on next symbol

- a) int b) average c) char d) all e) none

Q5 Find no of tokens in the following C code

printf("i=%d, %f\n", i, f);

- a) 3 b) 26 c) 10 d) 21

Sol:

& is an operator

so it must be counted separately

printf ("i=%d, %f\n", i, f);
① ② ③ ④ ⑤ ⑥ ⑦ ⑧ ⑨ ⑩

10 tokens

Difference b/w lexeme and token:

→ lexeme is token name in the source program

→ token is of form
token-name,
attribute value

which is generated by Lexical Analyzer

5: position = initial rate * 60;

For the lexeme position,

token is

$\langle id, 1 \rangle$

id is for identifier

1 points symbol table entry for position.

Similarly for lexeme initial, the token is $\langle id, 2 \rangle$

For lexeme =,

the token is $\langle \Rightarrow \rangle$. It does not

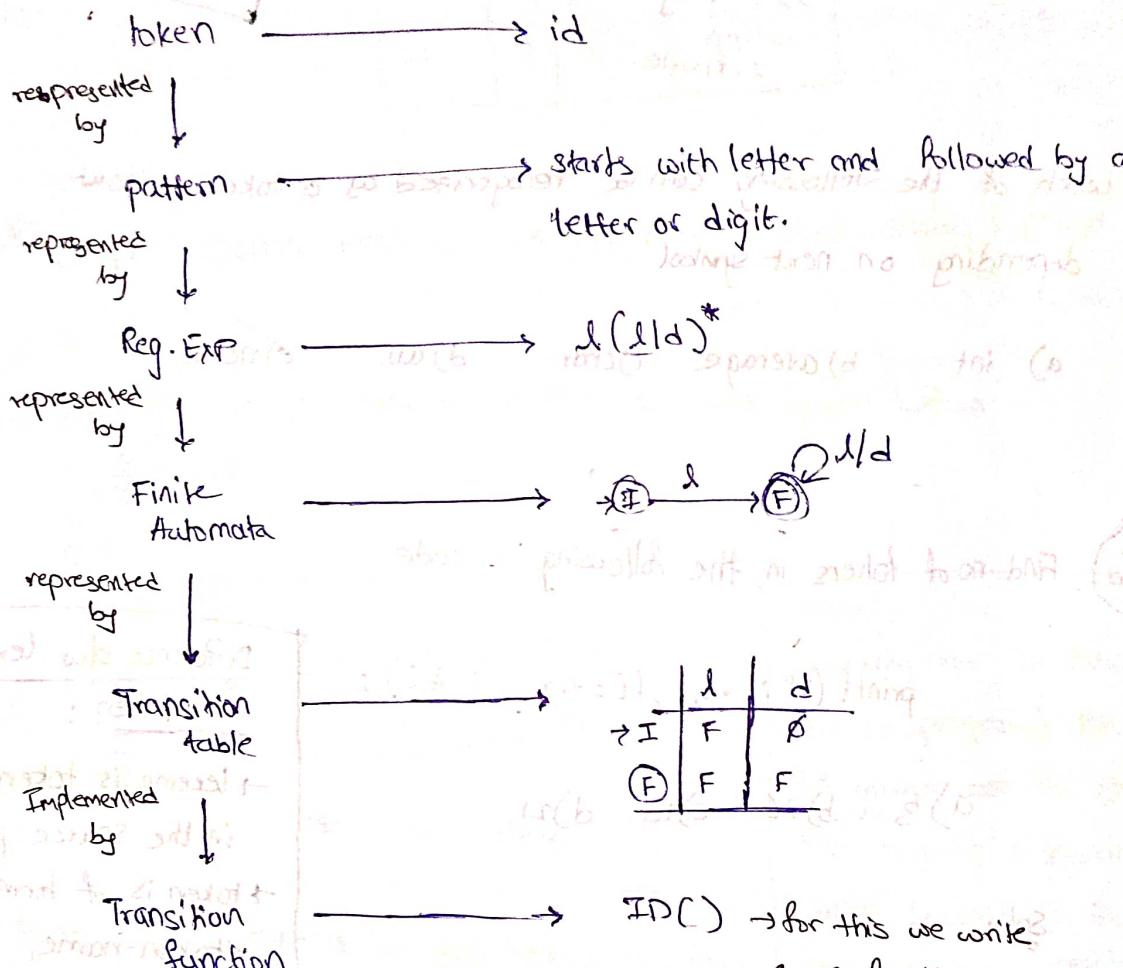
need any attribute value.
token of form $\langle \text{assign} \rangle$ can also be generated

Design of Lexical Analyzer:

→ we can either write a manual program for lexical analyzer or use tool (LEX - tool).

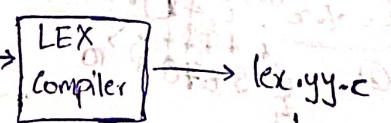
Manual design :

Consider we have a token of identifier type.

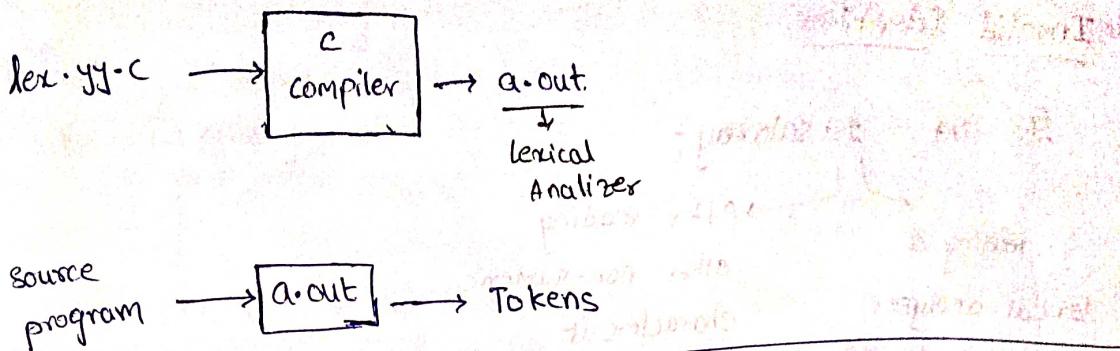


Designing with tool (not needed for gate)

Tool for lexical analyzer
LEX or Compiler.
here we specify token-type and regular exp.



It contains FA, its transition table and transition function.



Types of Lexical errors:

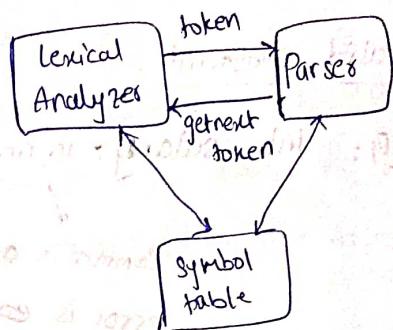
(i) Unterminated Comments error:

Eg : /* Line1
Line2
Line3

Here comment is not closed

Hence unterminated comment error occurs.

Process of lexical analysis

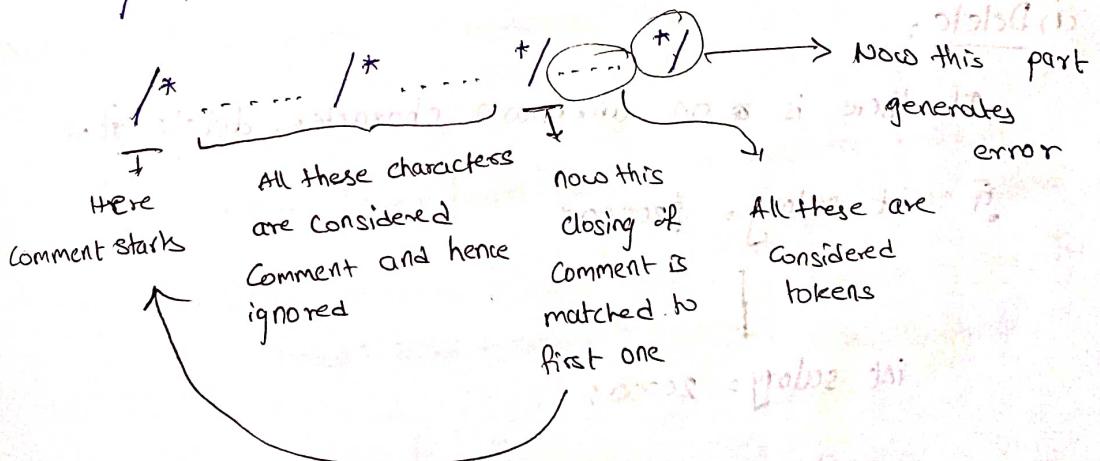


Thus parser and lexical analyzer run in parallel if needed.
If needed all front end can

be implemented in one module
There no definik
boundary b/w
phase of front
end

(ii) Nested Comments:

Eg : /* ~~.....~~ */



(iii) Length of longest identifier:

Eg:

int interest_on_principle_amount_and_compound_interest;

generally in C programming length of identifier shouldn't exceed 31

characters. So we ~~say~~ consider above line ~~as~~ as error.

(iv) Invalid Identifier

Eg: int 1st salary;
reading 1
lexical analyser
assumes it to be
a constant
After reading
other non-numeric
characters it
generates error.

(v) Invalid constant:

Eg: int salary = 10,000;
↓
Comma is operator or hence
error is generated

(vi) Illegal symbols:

Eg: int salary = \$20,000;

Error Recovery in LA:

i) Delete:

If there is an unknown character, delete it.

Eg: int salary = \$20000;

Unknown
character
↓
int salary = 20000;

ii) Insert:

Eg: /* (unterminated comment)

Correction
↓

/*



→ where to insert can't be
found though.

③ Transpose:

Eg: int \leftarrow swapped
int salary;

↓
int swapped salary;

④ Replace:

Eg: int salary = \$20000;
↓

It is replaced by a character allowed in the language.

If no appropriate character is found then it just replace '\$' with empty space (i.e., like delete).

Syntax Analysis (parser):

* To define the syntax of a language, CFG is used

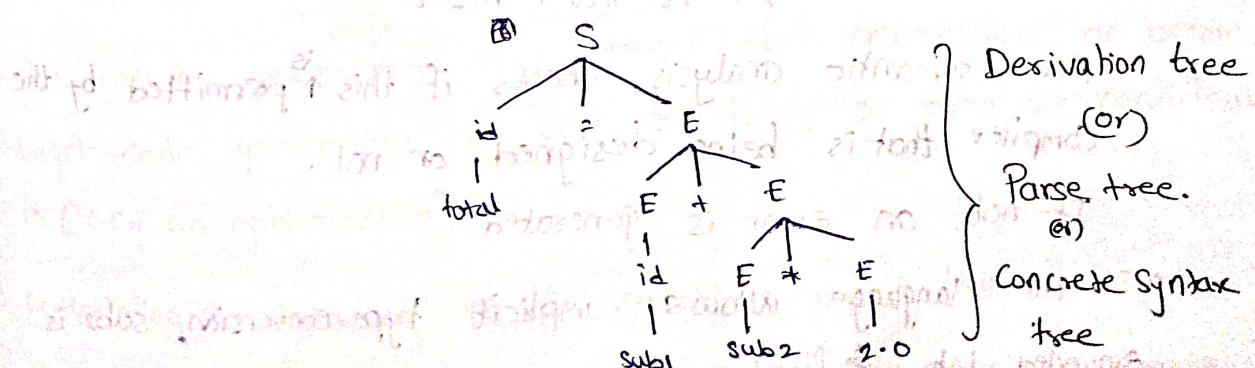
* Consider below ~~general~~ grammar to recognize certain

mathematical assignment statements.

Eg: $S \rightarrow id = E$

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

now consider deriving below stmt. Concrete syntax for the language



* The obj of syntax analysis is either parse tree or syntax tree.

* Parse tree and syntax tree are different from each other. (P)

Syntax tree:

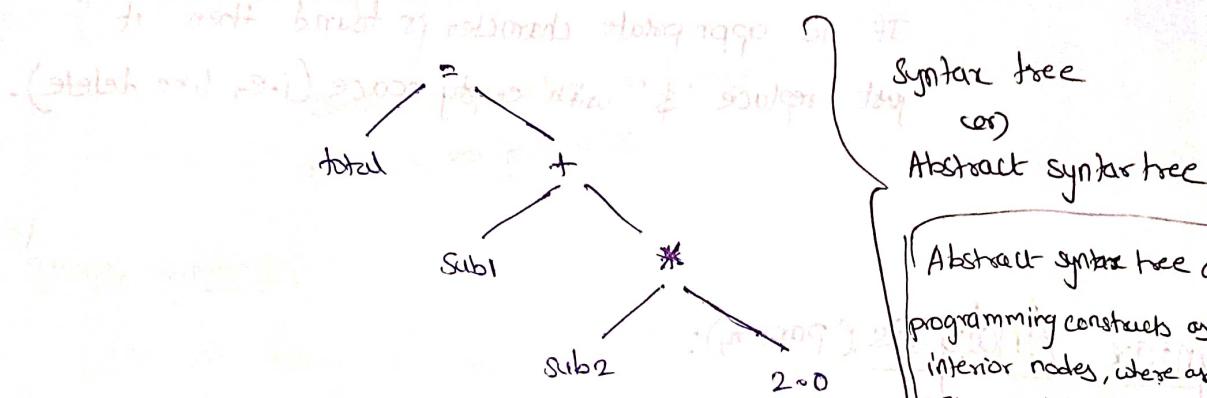
It is a compressed form of parse tree.

~~for pr~~

→ In syntax tree all the operators are present at internal nodes and operand are present at leaf nodes.

The syntax tree of previous example's parse tree is

as shown below



This syntax tree is given as input to ~~syntax~~ semantic analysis.

Semantic Analysis:

* Important component of semantic analysis is typechecking

Eg: At the bottom most level of the previous syntax tree

we are multiplying sub_2 with 2.0 .

Assume data type of sub_2 is integer;

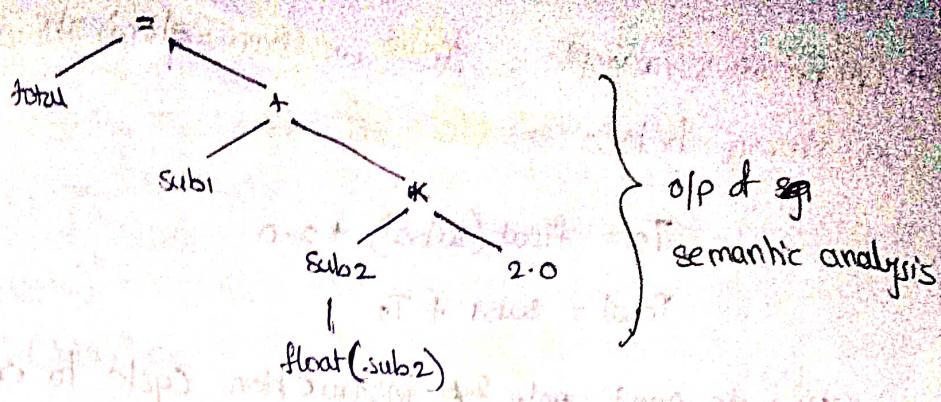
2.0 is real number.

→ Now semantic analysis checks if this is permitted by the compiler that is being designed, or not.

If not, an error is generated.

→ If the language allows implicit type conversion, sub_2 is converted into floating type.

→ Errors related to scope, undeclared or multiple declaration are detected at semantic analysis.



→ O/P of semantic analysis is called optimized syntax tree, or
semantically verified syntax tree.

→ This O/P is given as I/P to intermediate code generation phase.

Intermediate Code Generation:

- * There are different ways of representing IC.
- * Three Address Code (TAC) is mostly used representation of IC.
- * TAC's statements contain almost 3 memory addresses.

Eg: The IC generation for previous example is as shown

$$T_1 = \text{float}(\text{sub2});$$

$$T_2 = T_1 * 2.0;$$

$$T_3 = \text{sub1} + T_2$$

$$\text{Total} = T_3;$$

Code Optimization:

~~Optimization is not here in gate~~

→ This can be performed after target-code generation or before target code generation or both times or even we can leave it (i.e., no optimization).

→ Optimizing is reducing no of variables, no of instructions and some other operations without changing the logical meaning.

→ In previous example's IC we can eliminate T_3 and T_1 as shown below

$$T_2 = \text{float}(\text{sub2}) * 2.0$$

$$\text{Total} = \text{sub1} + T_2$$

Now we need only 2 instruction cycle to execute above code.

Also we have eliminated two temporary variable (T_1, T_2) which frees ~~some~~ some memory.

Code Generation:

Code generation
is not needed for
gate

optimization should balance

below 3 things:

- i) faster execution
- ii) shorter code
- iii) less power consumption

Here we generate object code (in this say assembly language program).

The below is target code for optimized IC

```
assembly
MOVF sub2, R0
MULF 2.0, R0
MOVF sub1, R1
ADDF R1, R0
STORE R0, Total
```

Later this is given as:

I/P to assembler.

Symbol table is constructed by lexical, syntax & semantic analysis phases. Creating entries in symbol table is mostly done by syntax analyzer

Two Additional phases:

(i) Symbol table:

It is a data structure containing a record for each identifier along with some other attributes for storing information. Procedure names are also stored in symbol table.

(ii) Error Handler:

* Errors are possible at every phase.

* After detecting an error, the phase must deal with error.

and compilation should not terminate and should continue the rest of the process of compilation.

Q6 6-18 A lexical analyzer uses the following patterns to recognize three tokens T_1, T_2, T_3 over the alphabet $\{a, b, c\}$.

$$T_1 : a? (b|c)^* a$$

$$T_2 : b? (a|c)^* b$$

$$T_3 : c? (b|a)^* c$$

Note that ' $x?$ ' means 0 (or) 1 occurrence of symbol x . Note also that the analyzer outputs the token that matches the longest possible prefix. If the string $bbaacabc$ is processed by the analyzer, which one of the following is the sequence of tokens it outputs?

- a) $T_1 T_2 T_3$ b) $T_1 T_1 T_2$ c) $T_2 T_1 T_3$ d) $T_3 T_3$

So:

$$a) T_1 T_2 T_3$$

$$b) T_1 T_1 T_3$$

$$\begin{matrix} bb | a & c & a & c & b & c \\ \text{---} & T_1 & T_2 & T_3 \end{matrix}$$

$$\begin{matrix} bb & a & c & a & c & b & c \\ \text{---} & T_1 & T_2 & T_3 \end{matrix}$$

c) $T_2 T_1 T_3$

$$\begin{matrix} b & b & a & a & c & a & c & b & c \\ \text{---} & T_2 & T_1 & T_3 \end{matrix}$$

$$bb | a a | c a c b c$$

Since we need to consider longest possible prefix

∴ opt (d)

Q7 6-20 Consider the following statements.

I. Symbol table is accessed only during lexical analysis and syntax analysis

II. Compilers for programming languages that support recursion necessarily need heap storage for memory allocation in the runtime environment

III. Errors violating the condition 'any variable must be declared before its use' are detected during syntax analysis.

which of the above statements is/are true?

- a) I only
- b) I and III only
- c) II only
- d) None of I, II and III

Sol: Statement II is true as stack is required to detect undeclared variables.

I: symbol table is used in all the phases

II: stack is enough to support recursion

III: This error is detected during semantic analysis.

∴ None

Syntax Analysis:

* This phase is also known as parsing.

* The syntactic rules are defined using a CFG.

CFG:

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

form of production is

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

Eg: $\frac{P}{S \rightarrow AB}$

$$A \rightarrow bA \mid E$$

$$B \rightarrow C$$

Now it is represented as

$$G = \left\{ \{S, A, B\}, \{b, c, E\}, P, S \right\}$$

Let bbc be i/p to above grammar

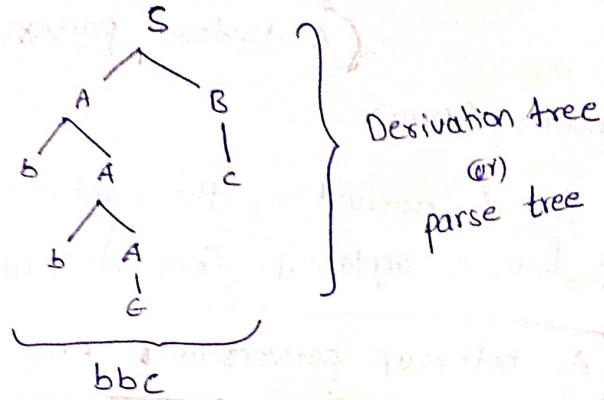
$$\frac{S \rightarrow AB}{\Rightarrow}$$

$S \rightarrow AB$
 $\rightarrow bAB$
 $\rightarrow bAc$
 $\rightarrow bbAc$
 $\rightarrow bbc$

Derivation

Here derived string bbc
is called sentence.

We can represent the derivation pictorially by derivation tree (parse tree)



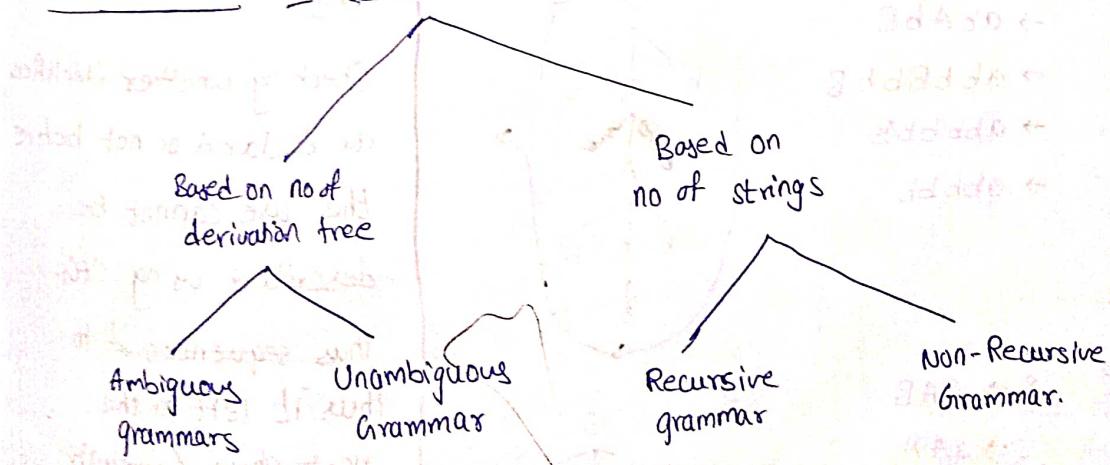
Derivation:

Derivation provides a means of generating the sentences of a language.

Derivation tree (or) Parse tree:

A ^{tree} parse is an equivalent form of showing a derivation which represents derivation graphically or pictorially

Classification of CFG:



Ambiguous Grammars: An ambiguous grammar is one that produces more than one left most derivation or more than one right most derivation for some sentence of a language.

Types of derivations:

There are two types of derivations:

(i) Left most derivation (LMD)

In each step of derivation, the left most non-terminal in the sentential form is replaced. This derivation is called LMD.

All topdown parsers use LMD

(ii) Right most derivation (RMD):

In each step of derivation, the right most non-terminal in the sentential form is replaced. This derivation is called RMD.

All bottom up parsers use "RMD in reverse"

$$\text{Ex: } S \rightarrow aAB$$

$$A \rightarrow bBb$$

$$B \rightarrow A \mid e$$

$$\text{I/P: } abbbb$$

Construct LMD & RMD for above string.

LMD:

$$S \rightarrow aAB$$

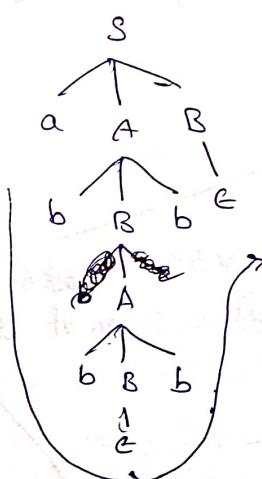
$$\rightarrow abBbB$$

$$\rightarrow abAbB$$

$$\rightarrow abbBbbB$$

$$\rightarrow abbbbB$$

$$\rightarrow abbbb$$



RMD:

$$S \rightarrow aAB$$

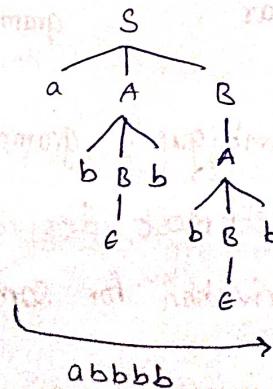
$$\rightarrow aAA$$

$$\rightarrow aAbBb$$

$$\rightarrow aAbb$$

$$\rightarrow abBbbb$$

$$\rightarrow abbbb$$



Note:

Checking whether identifiers are declared or not before the use cannot be described using CFG.

Thus sequences of tokens accepted by a parser forms a superset of the programming language.

Thus it left to the next phase (semantic analysis). Thus sequence of tokens accepted by a parser forms a superset of the programming language.

Find LMD & RMD for the following grammar

$$S \rightarrow aB | bA$$

$$A \rightarrow a | as | bAA$$

$$B \rightarrow b | bs | aBB$$

IP: aaabbabbba

B LMD:

$$S \rightarrow aB$$

$$\rightarrow aaBB$$

$$\rightarrow aaABBB$$

$$\rightarrow aaabBB$$

$$\rightarrow aaabbB$$

$$\rightarrow \cancel{aaabbbs}$$

$$\rightarrow \cancel{aaabbbs}$$

$$\rightarrow aaabbabb$$

$$\rightarrow aaabbabB$$

$$\rightarrow aaabbabB$$

$$\rightarrow aaabbabs$$

$$\rightarrow aaabbabbA$$

$$\rightarrow aaabbabbA$$

RMD:

$$S \rightarrow aB$$

$$\rightarrow aaBB$$

$$\rightarrow aaBAaBB$$

$$\rightarrow aaBaBbs$$

$$\rightarrow aaBaBbbA$$

$$\rightarrow aaBaBbbA$$

$$\rightarrow aaBabbba$$

$$\rightarrow aaABBabba$$

$$\rightarrow aaABbabba$$

$$\rightarrow aaabbabba$$

(Q9) Consider the following grammar

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

Check whether the grammar is ambiguous or not

Sol:

Consider

string id + id * id

LMD:

$$E \rightarrow E + E$$

$$\rightarrow id + E$$

$$\rightarrow id + E^* E$$

$$\rightarrow id + id^* E$$

$$\rightarrow id + id^* E$$

$$E \rightarrow E^* E$$

$$\rightarrow E + E^* E$$

$$\rightarrow id + E^* E$$

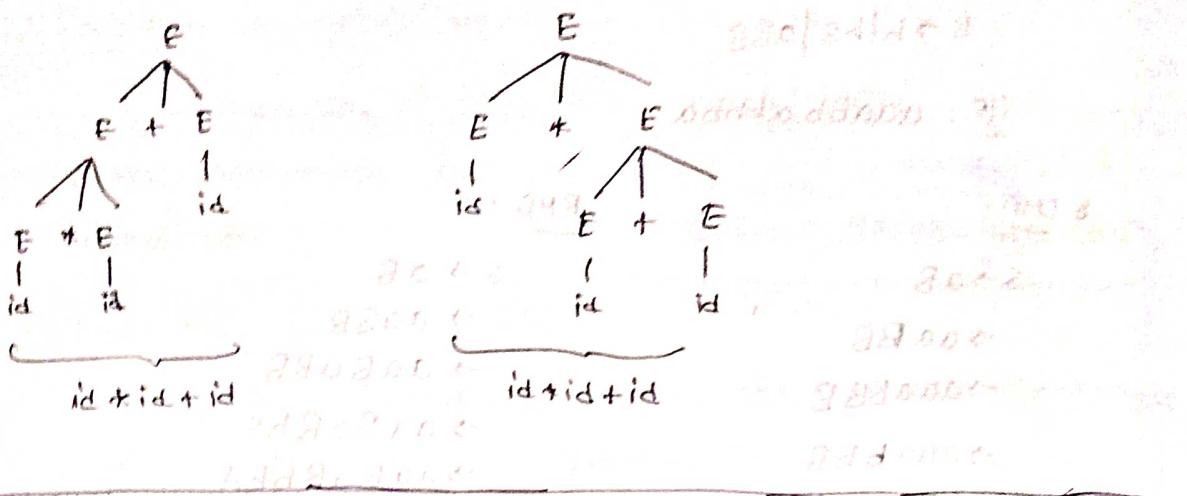
$$\rightarrow id + id^* E$$

$$\rightarrow id + id^* id$$

we have (2 left) most derivations for above string

∴ Ambiguous.

Generally for this type of questions try to construct two different parse tree for same input.



Observation: addition

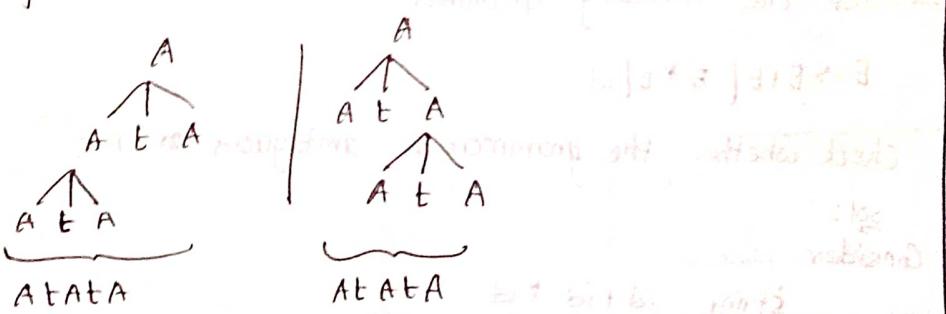
If we have $A \rightarrow A t A$ production of form

$A \rightarrow A t A$, t is terminal

then the grammar is ambiguous.

These productions are called ambiguous productions.

This ambiguous because we can expand left A once and right A once.



(Q10) 1M which of the following terminal strings has more than one parse tree? when parsed according to the grammar below.

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

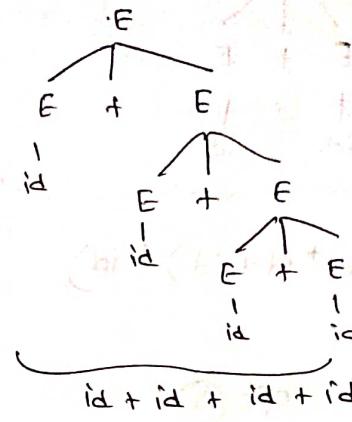
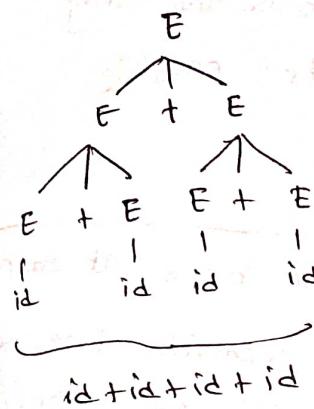
- a) id + id + id + id
- b) id + (id * (id * id))
- c) (id * (id * id)) + id
- d) ((id * id + id) * id)



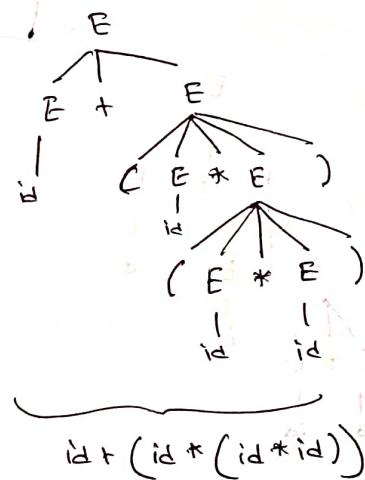
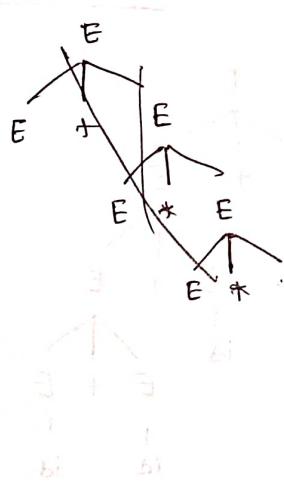
For the correct string with more than one parse tree obtained as solution to the above question how many parse trees are possible?

Sol:

i) a) id + id + id + id

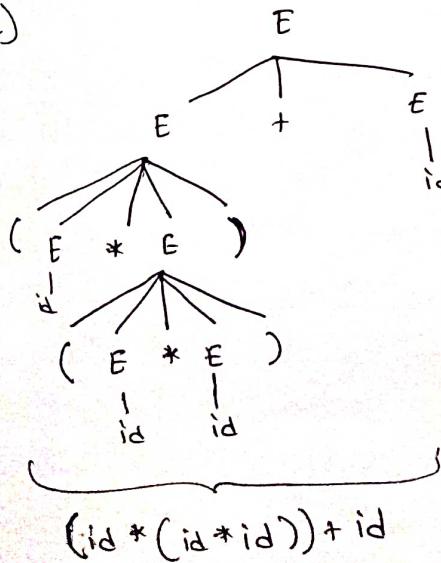


b)

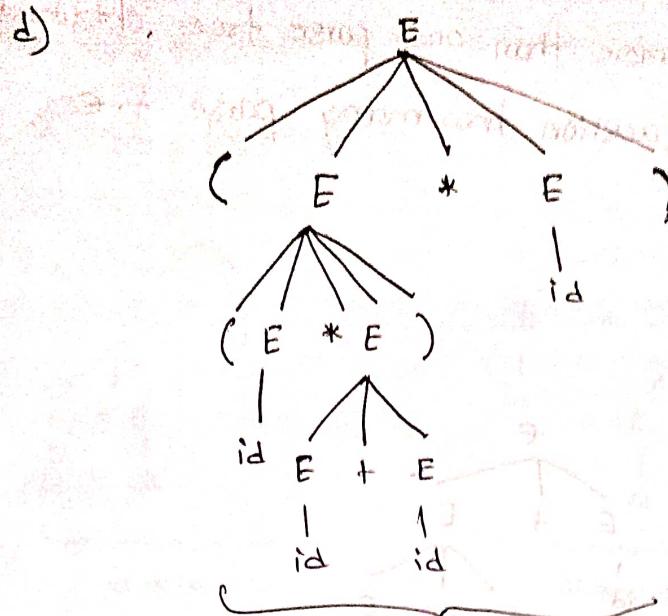


This is the only derivation tree possible

c)



This is the only tree possible



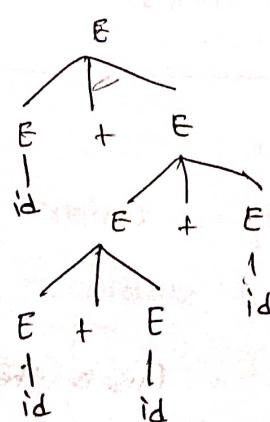
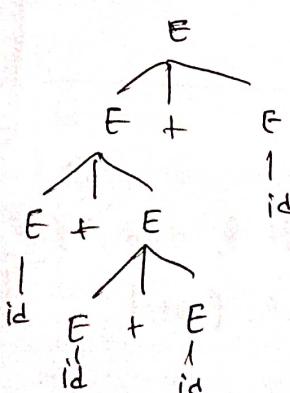
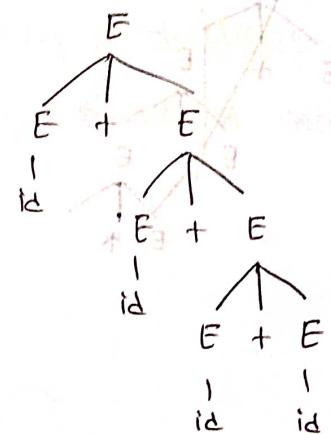
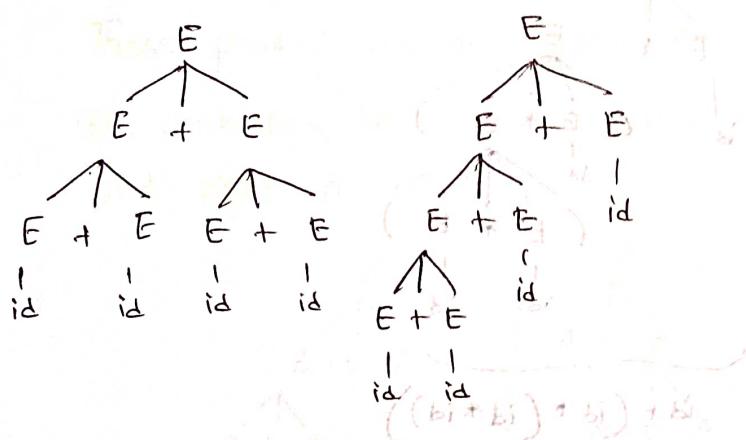
$((id * id + id) * id)$

opt ⑥

This is the only
tree possible

for option b/c/d at each
level of expanding tree, we
can find only one option (way)
 \therefore they have only one posse tree

II)



Consider the following grammar

$$S \rightarrow aS / A$$

$$A \rightarrow aAb / bAa / e$$

i) which of the following string is generated by the above grammar?

- a) aabbaba b) acbaaba c) abababb d) aabbbaa

ii) For the correct answer to the above question, How many steps are required to derive the string and how many parse trees are possible

- a) 6 & 1 b) 6 & 2 c) 7 & 2 d) 4 & 2

Sol:

- a) ~~aabbaba~~

Observing grammar we can say it accepts

any no of a's followed by

$$\begin{array}{c} \cancel{aAb} \\ \cancel{a} \\ \cancel{a} \\ \cancel{a} \\ \cancel{a} \\ \cancel{a} \end{array} \quad \begin{array}{c} aAb \\ bAa \\ e \end{array}$$

only opt ④ is of this form

\therefore opt ④

iii)

~~$S \rightarrow aS \mid Ad \mid aabbaba$~~

$$S \rightarrow aS$$

$$\rightarrow aas$$

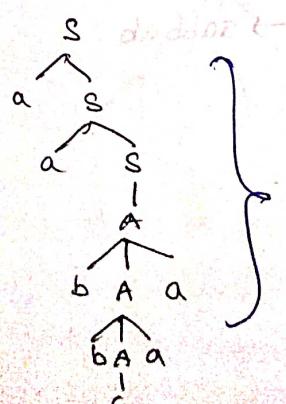
$$\rightarrow aaA$$

$$\rightarrow aabAa$$

$$\rightarrow aab bAaa$$

$$\rightarrow aabbbaa$$

6 steps



At each level of tree's expansion
we can find only one possibility

\therefore only 1 tree

$\therefore 6 \& 1$

19/06/20

132

Q12
G-07

Consider the CFG with $\{S, A, B\}$ as the non-terminal alphabet, $\{a, b\}$ as the terminal alphabet, S as the start symbol and the following set of production rules

$$S \rightarrow aB \quad S \rightarrow bA$$

$$B \rightarrow b \quad A \rightarrow a$$

$$B \rightarrow bS \quad A \rightarrow AS$$

$$B \rightarrow aBB \quad S \rightarrow bAA$$

i) Which of the following strings is generated by the grammar?

- a) aaaabb b) aabbbb c) aabbab d) abbbba

ii) For the correct answer to above question, how many derivation

trees are there for the string generated

- so a) 1 b) 2 c) 3 d) 4

so:

i) $S \rightarrow aB | bA | bAA$

$A \rightarrow a | aS$

$B \rightarrow b | bS | aBB$

a) aaaabb

b) aabbbb

c) aabbab

$S \rightarrow aB$

$\rightarrow aaBB$

$\rightarrow aaaaBBB$

$\rightarrow aaaaBBBB$

2

$S \rightarrow aB$

$\rightarrow aaBB$

$\rightarrow aabbB$

$\rightarrow aabBS$

$\rightarrow aabbBA$

$\rightarrow aabbABA$

x

$S \rightarrow aB$

$\rightarrow aAB$

$\rightarrow aabbB$

$\rightarrow aabbS$

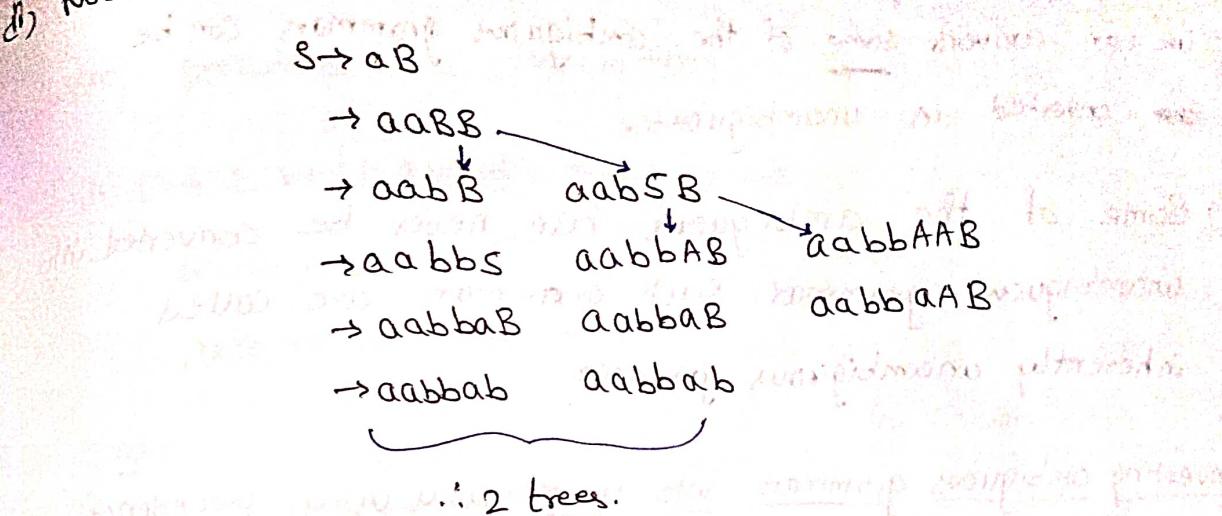
$\rightarrow aabbAB$

$\rightarrow aabbab$

(3)

\therefore opt C

Now we check for different possibilities of derivation



Note :-

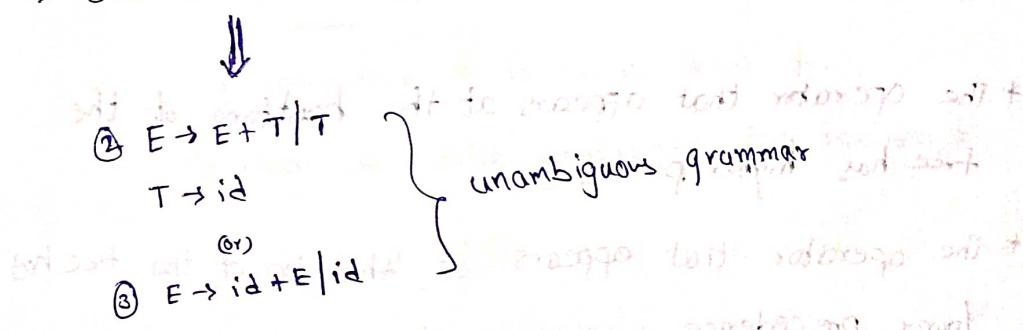
Note : algorithm exists to verify the ambiguity of CFGs.

→ No algorithm exists to convert ambiguous grammar into

* No algorithm exists to convert ambiguous grammar.

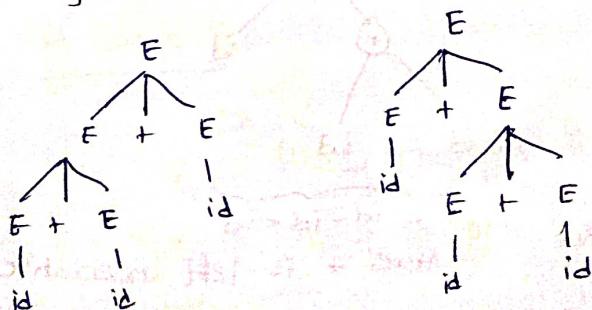
→ By using trial and error method some of the ambiguous grammars can be converted to unambiguous grammars.

$\stackrel{?}{\Rightarrow} : \textcircled{1} E \rightarrow E + E \mid id$ is ambiguous grammar.

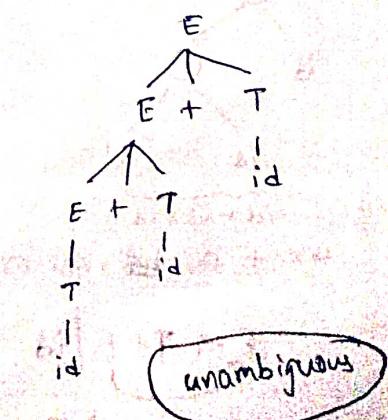


Consider $id + id + id$

deriving with grammar ①



Deriving with (2)

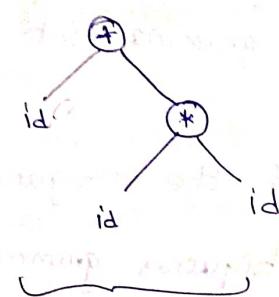


→ By redefining the grammar using precedence and associativity we can convert some of the ambiguous grammars can be converted in unambiguous.

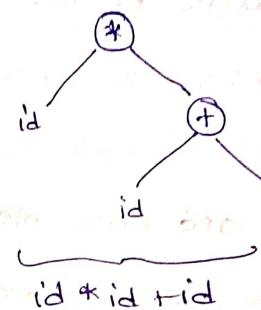
→ Some of the ambiguous can never be converted into unambiguous grammar. Such grammars are called inherently unambiguous grammar.

Converting ambiguous grammars into unambiguous using precedence and associativity:

→



$id + id * id$



$id * id + id$

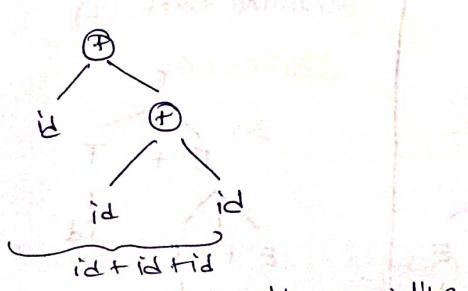
$P(*) > P(+)$

$P(+) > P(*)$

* The operator that appears at the bottom of the tree has higher precedence.

* The operator that appears at the top of the tree has lower precedence.

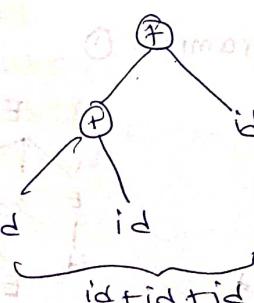
→



Here + is ~~left~~ right associative

$E \rightarrow T + E / id$

right ~~left~~ most symbol of RHS is equal to L.H.S



Here + is left associative

$E \rightarrow E + T / id$

leftmost symbol of RHS is equal to L.H.S.

(A3) Convert the following ambiguous grammars into unambiguous using precedence and associativity.

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

precedence: $* > +$

both are left associative

sg:

'+' w.r.t E & T:

Since '+' has least precedence, it should appear on the top of tree

i.e., '+' should be derived first.

$$\therefore E \rightarrow E + E$$

since '+' is left associative

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

'*' w.r.t T & F:

precedence of '*' is higher than '+'.

So now we write production for deriving '*'.

$$\therefore T \rightarrow T * T$$

But '*' is left associative

$$E \rightarrow \therefore T \rightarrow T * F \mid F$$

'id' w.r.t F:

E is finally changed to T

T is finally changed to F

\therefore we add production

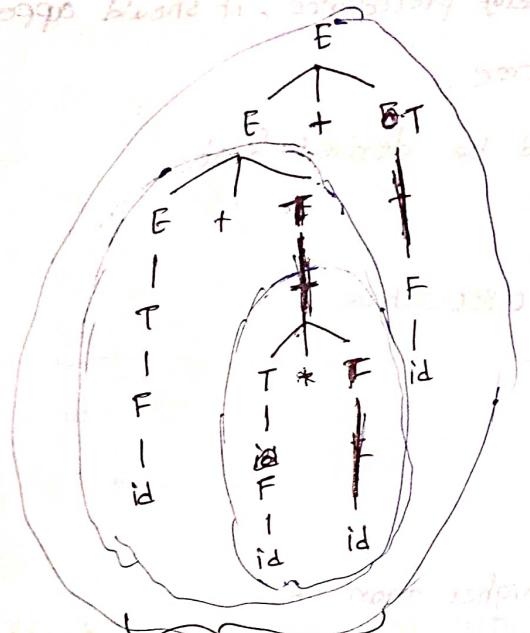
$$F \rightarrow id$$

∴ The modified grammar is

$$\begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T^* F \mid \text{if } F \\ F \rightarrow \text{id} \end{array}$$

Now this grammar obtained by considering precedence and associativity is unambiguous.

Consider the string $id + id * id + id$



Here, the expression evaluates to

$$\begin{array}{c}
 \text{id} + (\text{id} * \text{id}) + \text{id} \\
 \hline
 \text{---} \\
 \text{1st} \\
 \hline
 \text{---} \\
 \text{2nd} \\
 \hline
 \text{---} \\
 \text{3rd}
 \end{array}
 \quad \text{1st} \quad \text{2nd} \quad \text{3rd}$$

$$(ii) \quad E \rightarrow E+E \mid E * E \mid E - E \mid E \div E \mid E \wedge E \mid id$$

precedence: $\wedge > (* = \div) > (+ = -)$

Associativity: \wedge — right associative

$t_c - r \approx \frac{1}{2}$ ~~right~~ to associative
left

$+,-\ast$ should be on top of tree

$$E \rightarrow E+E \mid E-E$$

Associativity is left

$$\therefore E \rightarrow E+T \mid E-T \mid T$$

\ast, \div should be next

$$T \rightarrow T^{\ast}T \mid T \div T$$

Considering associativity

$$T \rightarrow T^{\ast}F \mid T \div F \mid F$$

\wedge should be next

$$F \rightarrow F^{\wedge}F$$

associativity is right

$$F \rightarrow G^{\wedge}F \mid G$$

Next we have id

$$\therefore G \rightarrow id$$

$$id \mid id + id \mid id \ast id \leftarrow id$$

$$id \mid id - id \mid id \div id \leftarrow id$$

So the required grammar is

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

$$T \rightarrow T^{\ast}F \mid T \div F \mid F$$

$$F \rightarrow G^{\wedge}F \mid G$$

$$G \rightarrow id$$

- Q14) Write precedence and associativity of the operators for the following grammar

$$i) E \rightarrow E^{\ast}F \mid F$$

$$F \rightarrow F^{\wedge}T \mid T$$

$$T \rightarrow P+T \mid P$$

$$P = id$$

Precedence: $(+ > (\cdot) > (*)$

Associativity: \cdot - left

\wedge - left

$+$ - right

$$(d) A \rightarrow @ A @ B | B$$

$$B \rightarrow C \# B | C ^ B | C$$

$$C \rightarrow C \$ D | C + D | D$$

$$D \rightarrow E - D | E * D | E$$

$$E \rightarrow id$$

Precedence: $(- > *) > (\$, +) > (\#, \wedge) > (@)$

Associativity: $@$ - left

$\#$ - ~~left~~ right

\wedge - ~~left~~ right

$\$, +$ - left

$-$, $*$ - right

Q15
G-

$$E \rightarrow E * F | F + E | F$$

$$F \rightarrow \cancel{F} E | F - F | id$$

choose the correct answer based on precedence

- a) $* > +$ b) $- > *$ c) $+ = -$ d) $+ > *$

Sol:

$$\text{opt } \cancel{E} \quad T | T - S | T + S < T$$

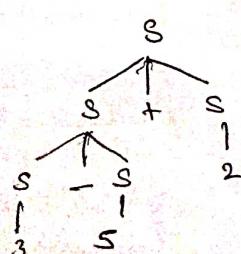
$$P(*) = P(+)$$

$$P(-) > P(*)$$

$\therefore \text{opt } b$

★
Q16
G-

Consider the following parse tree



which of the following based on precedence

- a) $(+) = (-)$ b) $(+) > -$
 c) $(-) > (+)$ d) none

Here question is according to the parse tree
and it is clear that '-' is being evaluated 1st

$$\therefore (-) > (*)$$

\therefore opt C

Consider the following grammar

(Q17)

$$S \rightarrow S + S \mid S * S \mid \text{id}$$

$$S^* F = \underline{S + S}$$

a) Choose correct option on precedence

- a) $(+) > *$
- b) $* > +$
- c) $+ = *$
- d) none

p

Note:

→ A production which is both, left and right associative is said

to be ambiguous grammar

Eg : $S \rightarrow S + S$
left associative
right associative

Consider the following left associative operators in decreasing order of precedence

(Q18)

what is the result of following expression

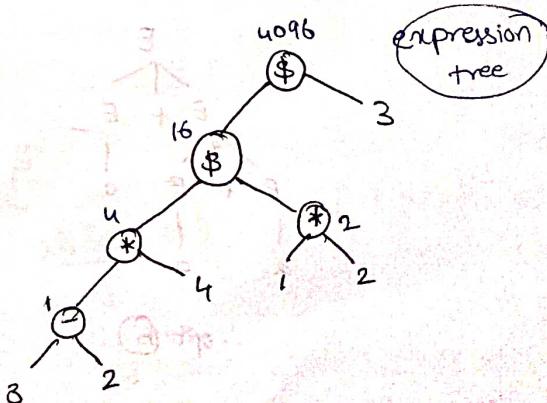
$$3 - 2 * 4 \$ 1 * 2 \$ 3$$

Sol :

$$\underline{3 - 2 * 4 \$ 1 * 2 \$ 3}$$

$$\begin{array}{c} 1 * 4 \$ 1 * 2 \$ 3 \\ \hline 4 \$ 1 * 2 \$ 3 \\ \hline 4 \$ 2 \$ 3 \end{array}$$

$$16 \$ 3 = 5096$$



operators

The operators in some programming language are given below

Q19
Q-15

2M

+	high precedence	left associative
*	medium precedence	Right associative
*	low precedence	left associative

the value of the expression $2 - 5 + 1 - 7 * 3$ in this language is _____

sol:

$$2 - 5 + 1 - 7 * 3 \quad (\text{from left to right})$$

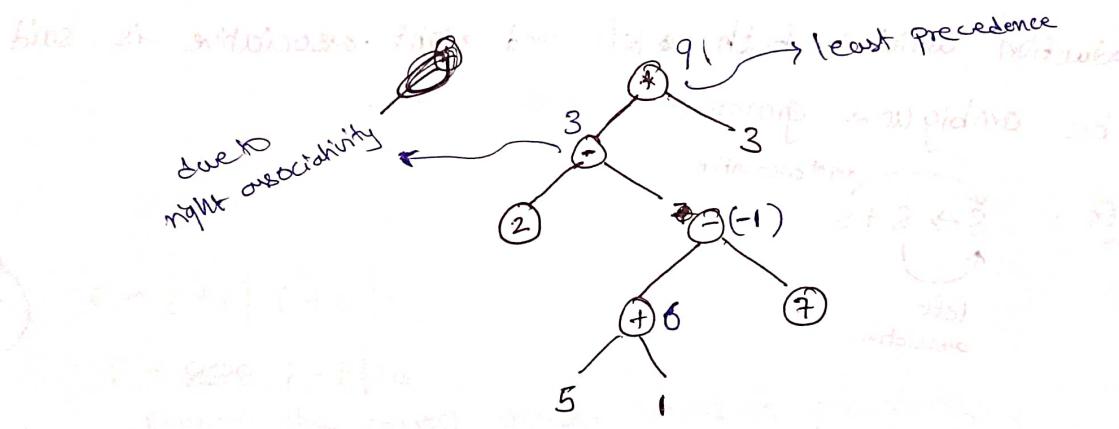
$$\underline{2 - 6 - 7 * 3} \quad (\text{from left to right})$$

$$2 - (-1) * 3 \quad (\text{right associative})$$

$$3 * 3$$

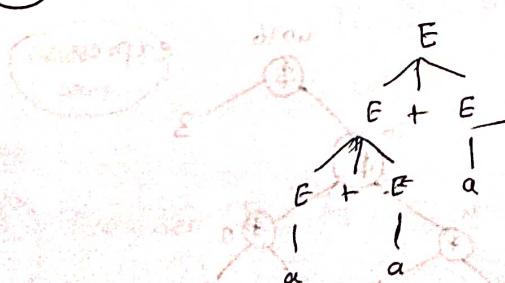
9

Solving by expression tree



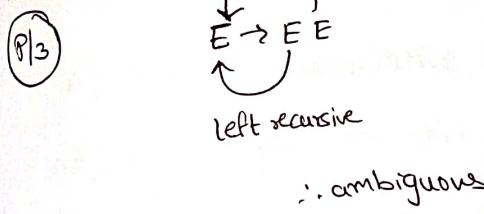
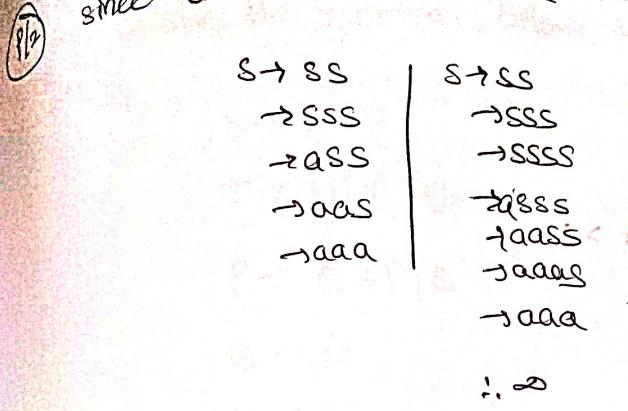
Material problems on parsing techniques:

(P/I) Since it is given that it is left associative



∴ opt ⑥

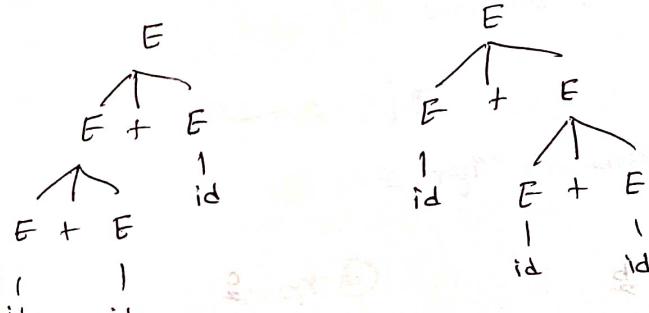
Since we have ' ϵ ' we have infinite derivations



(P14) observing grammar we can say for 'a' at right end we have 'b' at left end and viceversa. And in middle we have b.

This established by opt (d)

(P15)

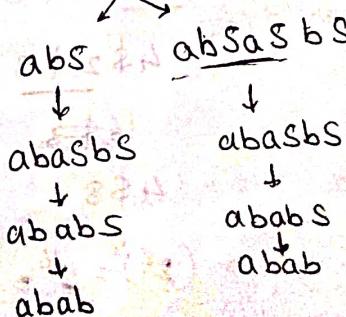


\therefore 2 trees.

(P16)

~~abab~~

$S \rightarrow aSbS$



\therefore 2 derivations
 \therefore 2 parse trees

P/9

$$\begin{array}{l} B \rightarrow E \uparrow T / T \\ T \rightarrow F + T / F \\ F \rightarrow i \end{array}$$

precedence
↑ - left
+ - right
↓ high

precedence : $+ > \uparrow$

$\therefore \text{opt } d$

142

P/10

~~P(+)~~ $P(+)$ $> P(\uparrow)$

$\therefore \text{opt } c$

143

P/11

$[P(*) = P(+)] < P(-)$

$\therefore \text{opt } b$

144

P/12

$3 - 2 * 4 \$ 2 * 3 \$ 2$

$1 * 4 \$ 2 * 3 \$ 2$

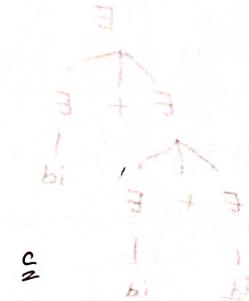
(b) \uparrow p. 144

$1 * 16 * 3 \$ 2$ ($\because \$$ is left associative)

$1 * 16 * 9$

$16 * 9$

144



P/13

a) $3 - 2 * 4 \$ 1 * 2 \$ 3$

$1 * 4 \$ 1 * 2 \$ 3$

$4 \$ 1 * 2 \$ 3$

$4 \$ 2 \$ 3$

$16 \$ 3$

$a = 4096$

$= 2^{12}$

$3 - 2 * 4 \$ 1 * 2 \$ 3$

$1 * 4 \$ 1 * 2 \$ 3$

$1 * 4 \$ 2 \$ 3$

$16 \$ 3$

$4 \$ 8$

256

$(2^8)^8 = 2^{16} =$

$3 - 2 * 4 \$ 1 * 2 \$ 3$

$3 - 2 * 4 * 2 \$ 3$

$3 - 2 * 4 * 8$

$1 * 4 * 8$

$4 * 8$

$c = 32$

$\therefore c < a < b$

145

146

Q20
6-07

Consider the grammar below with two operations * and +.

$$S \rightarrow T * P$$

$$T \rightarrow U | T * U$$

$$P \rightarrow Q + P | Q$$

$$Q \rightarrow id$$

$$U \rightarrow id$$

which one of the following is TRUE?

A) + is left associative, while * is right associative

B) + is right associative, while * is left associative

C) Both + and * are right associative

D) Both + and * are left associative

Sol:

$$P(+)>P(*)$$

$$T \rightarrow T * U | U$$

$\therefore *$ is left associative

$$P \rightarrow Q + P | Q$$

$\therefore +$ is right associative

\therefore opt B

If G is a grammar with productions

$S \rightarrow Sas | asb | bsa | ss | e$

Where S is the start symbol, then which one of the following strings

is not generated by G?

- a) abab b) acbab c) abbaa d) babba

Q21
6-07

\$ is at bottom

i. \$ has higher precedence than #

also in $b\$c\d

$b\$c$ is evaluated first

$\therefore \$$ is left associative

from expression it is clear the # is right associative

20/06/20

CFG based on number of strings:

1) Non-Recursive grammar:

A grammar which generates finite no of strings is said to be non-recursive grammar.

Eg: $S \rightarrow Aa \mid Bb$

$A \rightarrow c \mid e \mid b$

$B \rightarrow d \mid f$

ca, ea, db, fb are only string generated by the grammar.

\therefore It is non-recursive grammar.

2) Recursive Grammar:

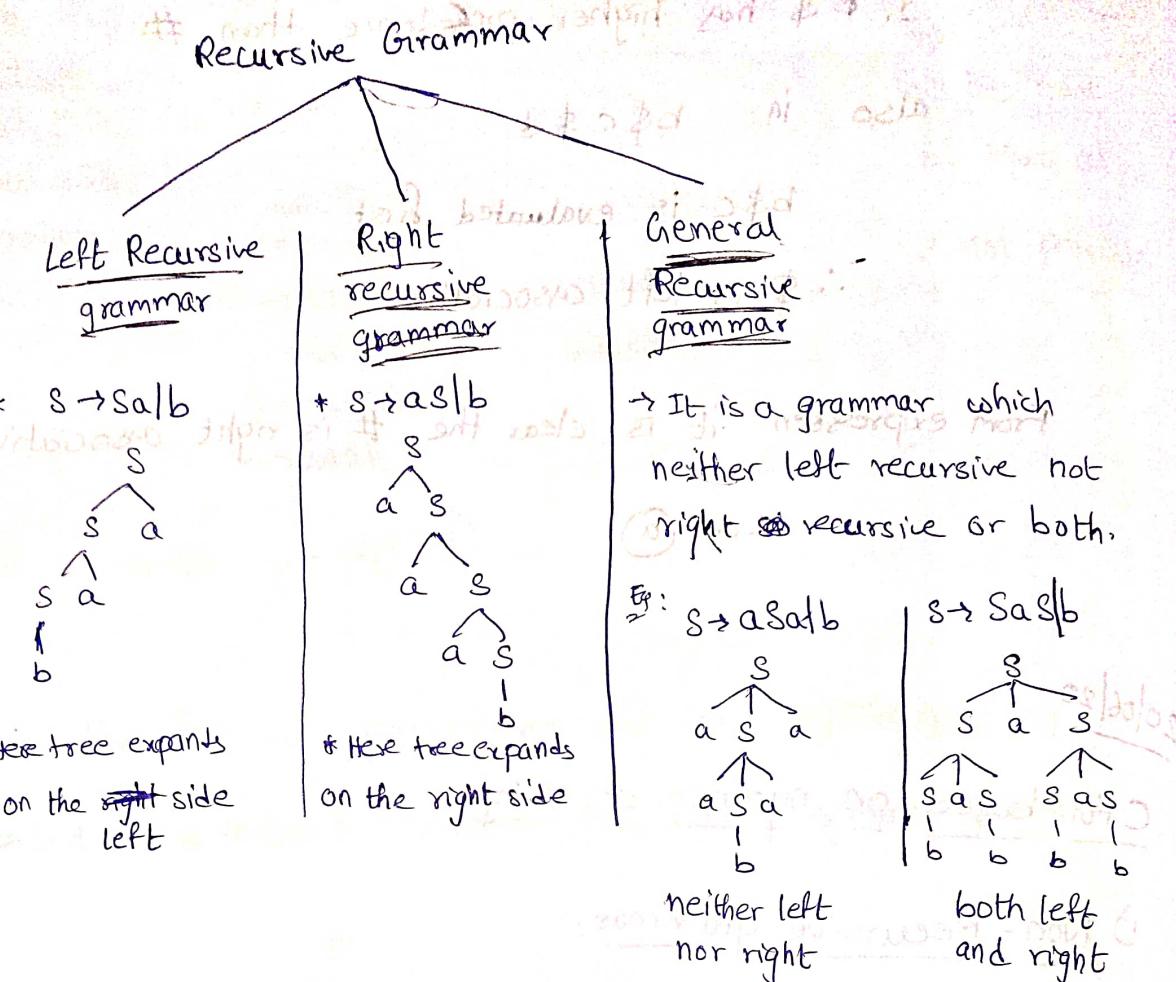
A grammar which generates infinite no of strings is called recursive grammar.

Eg: $S \rightarrow Sa \mid b$

The grammar generates b, ba, baa, baaa...

\therefore Recursive grammar.

Classification of Recursive grammars :



Note: A grammar which is both left & right recursive is said to be ambiguous grammar.

- Non-Recursive grammars are not used in the programming languages because we need to generate infinite string in programming language.
- Recursive languages are suitable for programming language.
- ~~Therefore~~ it is not that all recursive languages are suitable for every parsing technique.
- If a recursive ~~language~~ grammar is not suitable we need to make appropriate modifications on the grammar.

Modifying the grammar:

① Left recursive:

A grammar is left recursive, if the leftmost symbol in the RHS is same non-terminal as that in the LHS.

(Q1)

A grammar is said to be left recursive if it has a non-terminal 'A' with production $A \rightarrow A\alpha/\beta$ for $\alpha = (VUT)^*$

* Consider the production

$$A \rightarrow A\alpha/\beta$$

$A \rightarrow A\alpha$
 $\rightarrow A\alpha\alpha$
 $\rightarrow A\alpha\alpha\alpha$
 $\rightarrow A\alpha\alpha\dots\alpha$
 $\rightarrow \beta\alpha\alpha\dots\alpha$
 $\rightarrow \beta A'$ Here A' generates $\alpha^m n_2 \alpha$

The left recursion is removed by adding new non-terminal A' as follows

$$A \rightarrow A\alpha/\beta = A \rightarrow \beta A' \\ A' \rightarrow \alpha A'/\epsilon$$

A grammar is said to be left recursive if it has at least one non-terminal from whose derivation we can obtain a sentential form such that the leftmost symbol of the sentential form is same as that

Eg: Eliminate left recursion from below grammar of L.H.S

$$E \rightarrow E + T / T$$

$$T \rightarrow T^* F / F$$

$$F \rightarrow id$$

Sol:

$$\underline{E \rightarrow E + T / T}$$

$$\Rightarrow E \rightarrow T E'$$

$$E' \rightarrow + T E' / \epsilon$$

$$\underline{T \rightarrow T^* F / F}$$

$$T \rightarrow F T'$$

$$T' \rightarrow * F T' / \epsilon$$

$$A \xrightarrow{*} A\alpha$$

α is $(VUT)^*$

$\xrightarrow{*}$ mean series of steps in derivation

$F \rightarrow id$ is not recursive

\therefore Required grammar is

$$E \rightarrow TE'$$

$$E' \rightarrow +TE'/G$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT'/G$$

$$F \rightarrow id$$

Q23) Eliminate left recursion from the below grammar

$$S \rightarrow Sa/b/d$$

$$S \rightarrow bS'/dS'$$

$$S' \rightarrow aS'/E$$

$$A \rightarrow A\alpha_1 | \cancel{A\alpha_2} | \cancel{A\alpha_3} | \dots | \cancel{A\alpha_m} B_1 | B_2 | B_3 | \dots | B_m$$

↓ Eliminating left recursion

$$A \rightarrow B_1 A' | B_2 A' | \dots | B_m A'$$

$$A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_n A' | E$$

~~A₁, A₂, A₃~~

$$A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_n | B_1 | B_2 | \dots | B_m$$

leftmost not final \downarrow

most significant position

$$A \rightarrow B_1 A' | B_2 A' | \dots | B_m A'$$

$$A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_n A' | E$$

Q24) Eliminate left recursion from below grammar

$$A \rightarrow Aa | Ab | c | d$$

↓ eliminating left recursion

$$A \rightarrow cA' | dA'$$

$$A' \rightarrow aA' | bA' | E$$

Indirect Left recursion:

Consider the grammar

$$S \rightarrow Aa/b$$

$$\& A \rightarrow ad/e$$

Consider derivation

$$S \rightarrow Aa \rightarrow Sd/e$$

Indirect left recursion

$$\begin{aligned}
 S &\rightarrow AaBb \\
 B &\rightarrow CAd/e \\
 C &\rightarrow d/e \\
 \end{aligned}$$

A is not indirect
left recursion.
 $\xrightarrow{A \rightarrow ad} \xrightarrow{C \rightarrow d} \xrightarrow{A \rightarrow ad} Ada$

In this case try to draw
parse tree marking parser
for string "ada" and see
how left recursion effects
immediate left recursion.

- Direct left recursion is also called immediate left recursion.
- Indirect left recursion should also be eliminated if needed.
able to do down parsers.

Eliminating indirect left recursion

→ we first convert indirect left recursion to direct left recursion.

Then we eliminate it.

Eg: Consider the below grammar

$$S \rightarrow Aa/b$$

$$A \rightarrow Ac/Sd/e$$

Here S has indirect left recursion
A has direct & indirect left recursion

Eliminating indirect left recursion of S

$$S \rightarrow Aa/b$$

$$A \rightarrow Ac \mid Aad \mid bd \mid e$$

↳ elimination of indirect left recursion of S.

with this indirect left recursion of A is also eliminated.

Now we remove direct left recursion of A

$S \rightarrow Aa|b$

$A \rightarrow BdA' | eA'$

$A' \rightarrow cA' | adA' | e$

Thus the above grammar is now free from both direct and indirect left recursion.

(Q23) Eliminate left recursion from below grammar.

$S \rightarrow Aa|b$

$A \rightarrow Bd|e$

$B \rightarrow Sf|g$

Sol:

$S \rightarrow Aa \rightarrow BdA \rightarrow SfdA$

$\therefore S$ has indirect left recursion

$A \rightarrow Bd \rightarrow Sfd \rightarrow Aafd$

$\therefore A$ has indirect left recursion

$B \rightarrow Sf \rightarrow Aaf \rightarrow Bdaf$

$\therefore B$ has indirect left recursion

Eliminating indirect left recursion of S & A .

$S \rightarrow Aa|b$

$A \rightarrow Bd|e$

$B \rightarrow Aaf | bf | g$

$A \rightarrow Bd \rightarrow Aafd$

Eliminating indirect left recursion of A .

$S \rightarrow Aa|b$

$A \rightarrow Bd|e$

$B \rightarrow Bdaf | eaf | bf | g$

Now we need to eliminate direct left recursion of B

$s \mapsto Aa/b$

$A \rightarrow Bd/e$

$$B \rightarrow eaf B' \mid bf B' \mid g B'$$

$$B' \rightarrow \text{def } B' | e$$

Note: Here all we did is ~~making~~ making B direct recursive
and eliminating it.

Similarly we can do the same for S or A.
 Thus we don't have unique answer for this kind of question.

which of the following grammars is free from left recursion.

$$\begin{array}{l} \text{a) } S \rightarrow A B \\ \quad A \rightarrow A a b \\ \quad B \rightarrow C \end{array}$$

$$\begin{array}{c}
 b) S \rightarrow Ab \mid Bb \mid c \\
 A \rightarrow Bd \mid \epsilon \\
 B \rightarrow \epsilon
 \end{array}$$

$$\begin{array}{l} \text{c) } S \rightarrow Aa \mid B \\ \quad A \rightarrow Bb \mid Sc \mid e \\ \quad B \rightarrow d \end{array}$$

d) $S \rightarrow Aa | Bb | c$
~~A $\rightarrow Ba | c$~~
~~B $\rightarrow Ae | f$~~

SOL:

a) ~~S \ominus~~ A \rightarrow Aa

∴ left recursive

b) $\text{g} \rightarrow \text{Ab} \rightarrow \text{Bd b}$
 $\quad \quad \quad \downarrow$
 $\quad \quad \quad \text{Bb} \rightarrow \text{eb}$

\therefore no left recursion

$$c) S \rightarrow Aa \xrightarrow{\text{K}} Sca$$

- left recursive

$$(d) \quad \text{Aa} \rightarrow \text{Aa} \rightarrow \text{Ae}$$

$$A \rightarrow Bd \rightarrow Aed$$

∴ left recursive

Q27
G1-17

Consider the following expression grammar G1:

$$E \rightarrow E - T \mid T$$

$$T \rightarrow T + F \mid F$$

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

which of the following grammars is not left recursive, but is equivalent to G1?

a) $E \rightarrow E - T \mid T$

$T \rightarrow T + F \mid F$

$F \rightarrow (E) \mid id$

b) $E \rightarrow TE' \mid$

$E' \rightarrow -TE' \mid E$

$T \rightarrow T + F \mid F$

$F \rightarrow (E) \mid id$

c) $E \rightarrow TX$

$X \rightarrow -TX \mid E$

$T \rightarrow FY$

$Y \rightarrow +FY \mid E$

$F \rightarrow (E) \mid id$

d) $E \rightarrow TX(TX)$

$X \rightarrow -TX + TX \mid E$

$T \rightarrow id$

sol:

opt ④ & ⑤ has direct left recursion

so they can be eliminated.

e) $E \rightarrow TE'$

$E' \rightarrow -TE' \mid E$

$T \rightarrow FT'$

$T' \rightarrow +FT' \mid E$

$F \rightarrow (E) \mid id$

But in question X is used as E'

Y is used as T'

\therefore opt ⑥ is correct

d) Though the grammar has no left recursion, it produces different grammar.

Material problems

(P/7)

$$A \rightarrow Sa/b$$

$$S \rightarrow Sc/Ad$$

↓

$$A \rightarrow Sa/b$$

$$S \rightarrow Sc/Sad/bd$$

↓

$$A \rightarrow Sa/b$$

$$S \rightarrow bdS'$$

$$S' \rightarrow CS'/ads'/\epsilon$$

opt (a)

(P/8)

same as (Q/7)

Left Factoring :

It is a process which isolates the common parts of two or more productions into a single production.

Any production of form

$$A \rightarrow \alpha\beta_1/\alpha\beta_2/\dots/\alpha\beta_n$$

can be modified as

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1/\beta_2/\dots/\beta_n$$

→ Here if string starts with ' α ' the parser doesn't know which of derivations to choose. So we take ' α ' as common and create new non-terminal. This process is called factoring.

→ Left factoring helps parser make accurate decisions and

reduces ambiguity in taking decision.

→ Also left factoring reduces backtracking.

Eg: Left factor the following grammar:

$$S \rightarrow iEts \mid iEtSe \mid a$$

$$E \rightarrow b$$

Sol:

iEts is common

$$S \rightarrow iEtsS' \mid a$$

$$S' \rightarrow eS \mid \epsilon$$

$$E \rightarrow b$$

$$a \mid ab \mid aa$$

$$BA \mid BB \mid B$$

$$d \mid bd \mid dd$$

$$ab \mid bba \mid bbb$$

$$d \mid bd \mid dd$$

$$12 \mid bd \mid dd$$

$$ab \mid bba \mid bbb$$

Eg: $A \rightarrow abA \mid abB \mid abc \mid ad$ (or)

Eg: //

$$A \rightarrow aA'$$

$$A' \rightarrow bB \mid bc \mid ad$$

//

$$A \rightarrow abA' \mid ad$$

$$A' \rightarrow B \mid C$$

//

$$A \rightarrow aA'$$

$$A' \rightarrow bA'' \mid d$$

$$A'' \rightarrow B \mid C$$

$$A \rightarrow aA''$$

$$A'' \rightarrow bA' \mid d$$

$$A' \rightarrow B \mid C$$

: partition first

After this point B and C
are not part of the
initial part of string

Parsing:

Parser: A program for a given CFG 'G' which
take a string w as input and output

(i) parse tree if $w \in L(G)$

(ii) error message if $w \notin L(G)$

This process is known as parsing.

Parsing techniques are of two types:

15B

Parsing

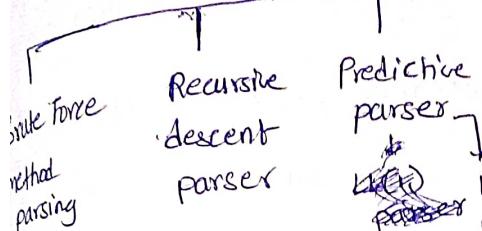
Top-Down parsing

* uses LMD

* tree is constructed

from root to leaves

* top-down parsing algorithms

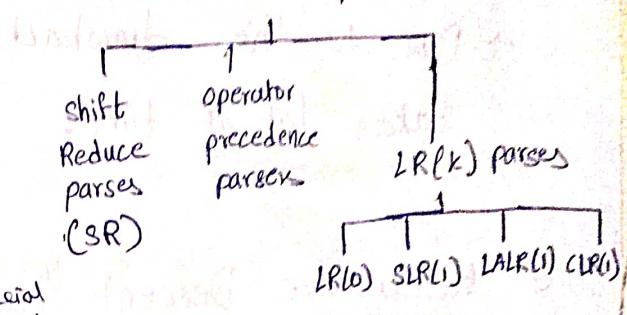


Bottom-up parsing

* uses PND in reverse

+ leaves to top

* Bottom-up parsers



Top-down parsing:

A top-down parser constructs a parse tree starting from starting non-terminal symbol (root) to input string (leaves)

→ During this process leftmost derivation is used.

Brute Force Method: It is a topdown parser with full backtracking

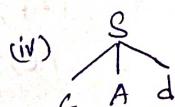
Eg: Consider the grammar

$$S \rightarrow C A d$$

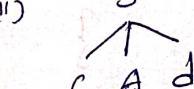
$$A \rightarrow a b / a$$

Below is topdown parsing of string "cad"

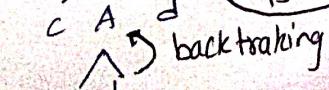
Step: (i) S



(ii) S

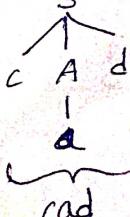


(iii) S



Always 1st alternative is tried

(iv) S



drawbacks: ① Backtracking

→ We use other parsing techniques to overcome this drawback of ~~back~~ backtracking.

→ The term "full backtracking" means if needed the brute force parse tracks back even to start symbol if needed.

→ Due to this drawback, the brute force technique takes lot of time.

② Recursive Descent Parser:

It is a top-down parser ~~with no backtracking~~.

* It is done by representing each production rule by a function call (with if else or switch case stmts)

Eg: For the below grammar

$$\begin{array}{l} S \rightarrow aA \\ A \rightarrow b \end{array}$$

We write functions as

S()

{ if(ip[i] == 'a')

{ i++;

A();

else

error();

A()

{ if(ip[i] == 'b')

accept();

else

error();

main()

{ S();

Drawbacks: Recursive Function Calls

- * Sometimes there is a chance for infinite loop due to recursive function calls.
- * Backtracking.
- Also Recursive descent parser and brute force parser can parse only a small class of ~~languages~~ grammars.
- A simple form of recursive descent parsing without backtracking is called predictive parsing.
- A left-recursive grammar can cause a recursive-descent parser, even one with backtracking, to go into an infinite loop.