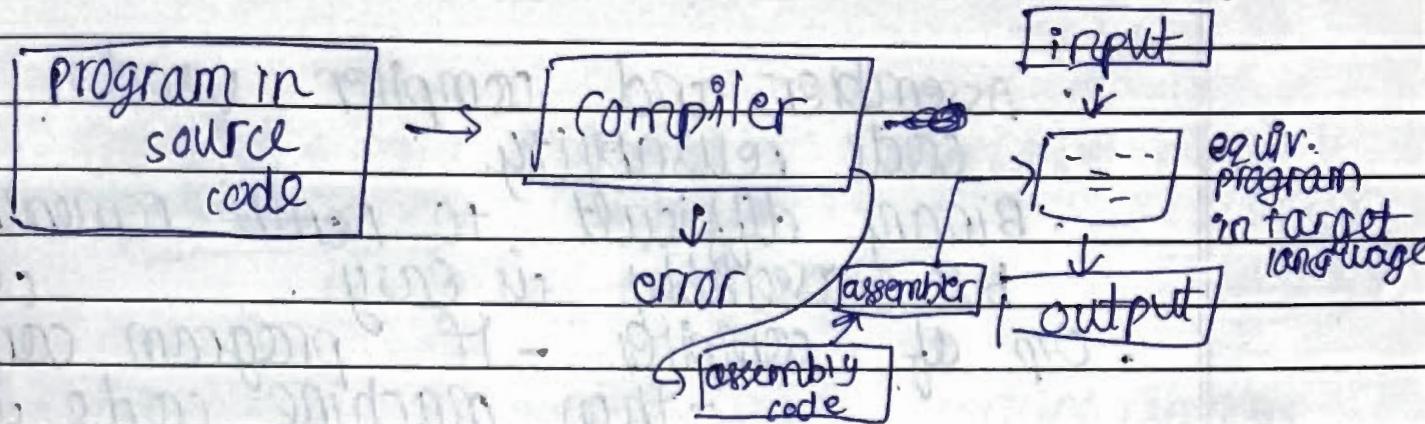


CCCC

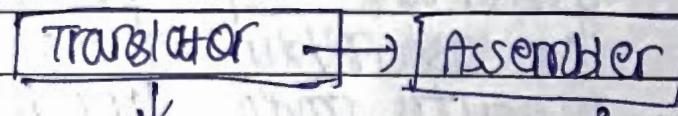
- Compiler - Translates language target - low level computer lang.
- compiler - Translating program written in source language into a semantically equivalent program written language. ↗ done by translator. ↗ eg - binary ↗ assembly language
- easily understood by computer - low level language.



- compiler is a special type of translator.
- Assembly code is diff. for diff. architecture.

↳ eg - ADD Al Bl
 instruction registers
 ↓

↳ still not very low for computer.



↳ binary language

↳ type of translator compiler.
 ↳ target - binary.

* Python has both
compiler & interpreter
↳ py file

Date : / / 20

Page No.:

also a translator

- For compiler, target lang. can be binary or assembly

[compiler] → Assembly code

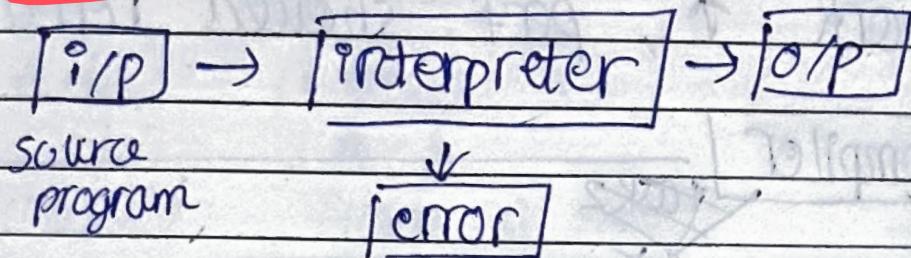
[compiler] → Binary code

- Assembler and compiler used code reusability.
- Binary difficult to remb remember but assembly is easy. syntactical
- O/p of compiler - If program correct - then machine code is generated.

→ Thus, compiler translates if i/p is correct, else list of error.

→ Interpreter, another translator, doesn't generate o/p program code in binary or assembly code but translates & executes code line by line
eg- translator for PM mode
translate line by line,
not all at once.
eg- command prompt

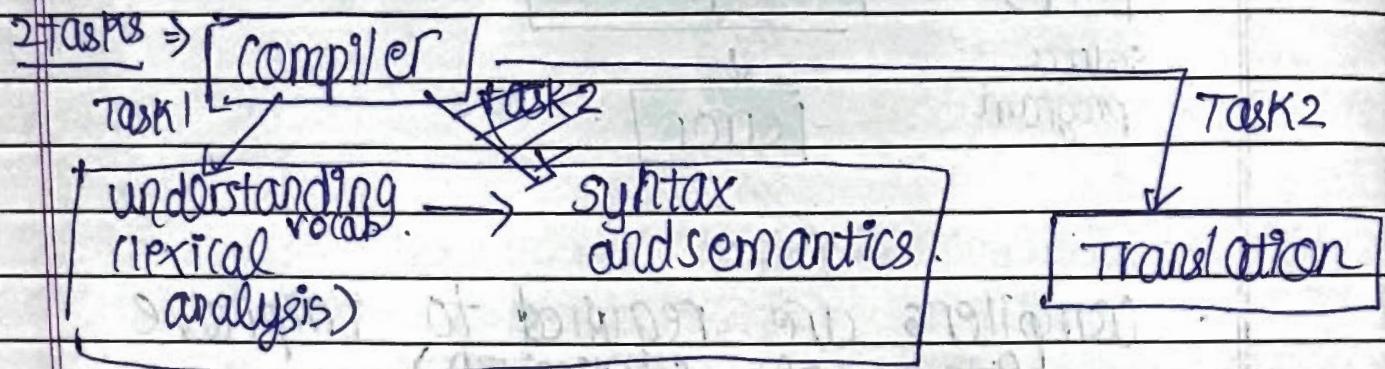
Note - Debuggers always require interpreter



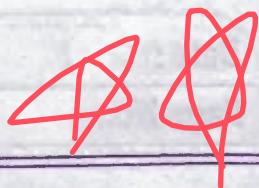
- Compilers are required to improve performance (optimiz'n)
↳ req. of every code.
 - compilers ↑ system perf., reduce system load & are needed for portability of compiled code.
- * 2 main purposes - trans'n
- optimiz'n of perf.

- A c exec file is machine dependant
 - may or may not run due to machine architecture.
- performance is needed in almost all domains.

- specific companies produce specific compilers.
- optimiz^n is a continuous req. Even if work ↑, perf. should remain same



- optimiz^n happens after correct transl.



Lab

Regular expressions

$a|b = a \text{ or } b$

$\rightarrow [a|b]$

$\hookrightarrow a \text{ or } b \text{ or } '1'$

$[a-c] = \text{from } a \text{ to } c$

$\rightarrow . = \text{anything}$

$ac = ac$

"ac" = ac

* $[\wedge a-c] = \text{except } a \text{ to } c$

* $^a a-c = \text{starts with } ac.$

$ac\$ = \text{end with } ac.$

$[a-z]^+ = \text{one or more occurrences of } a-z$

$[a-z]^* = 0 \text{ or more occurrences.}$

$[a-z]^? = 0 \text{ or } 1 \text{ occurrence.}$

* $[a-z]\{2,5\} = 2 \text{ to } 5 \text{ occurrences}$

$[a-z]^*\{2,3\} = 2 \text{ or more occurrences.}$

* $[a-z]^+^{\prime \prime} = \text{look ahead character (/)}$

$\rightarrow \text{To make a flex.c file, } \Rightarrow \underline{\text{Format}}$

% # { declarⁿ section
variable declarⁿ
func declarⁿ

% # ?

0/0 0#0 %.%

RE1 { Translⁿ Rule 3

RE2 { Translⁿ Rule 3

0/0 0#0 User function

first.l

{ dcl's \Rightarrow

2.5

二六

二、五

`[-a-zA-Z][a-zA-Z0-9]* {printf
 comment }
identifier`

"if" / "else"

% printf ("keyword")

↳ longest pattern match

$\Rightarrow \text{first.l} \rightarrow \boxed{\text{flex}} \rightarrow \text{lex.yy.c}$

↓
{ gcc }

a.out

HLL = high level language
LLL = Low level language

Date: 6 / 1 / 20

Page No.:

class

→ a/c produces assembly code.

Q: Identify assemblers & assemblers producing binary language.

Ans eg- NASM (netwide assembler), MASM (microsoft macro assembler), GAS (gnu assembler), TASM (turbo assembler)

→ compiler	Source	target
CC	C	Binary / machine
C++	C	Bin. / mach.
Javac	Java	Byte code

→ source to source compiler (transpiler / compiler)

- Produces o/p in another lang. rather than binary (1 programming lang i/p & another lang. o/p)

eg- C++ to C transpil.

- need legacy code to new code (lang.)

optimizⁿ

↳ parallel computat.

- need to check which operⁿ can be done parallelly

a=a+1 , b=b+1 X

a=a+1 , b=b+1 V

cross compiler

→ eg- compiler

runs on PC but generates that runs on android

used for cloud computing

- IAB, PAs, SAs

produce executable machine code

on another platform

(where other than one where compiler is running)

→ Bootstrap compiler

- written in programming lang. that they have to compile
- eg- compiler written in C language

→ decompiler

- translates exe file to high level source file which can be recompiled successfully

eg- source code not there, so reverse engineer the exe.

⇒ Java is both compiled & interpreted because source code is first compiled into a binary byte code for a machine which doesn't exist (Java virtual machine).

- This byte code runs on JVM which is software based interpreter.

~~P~~ Compiler is faster than interpreter as CPU can execute instruction faster in GHz when everything else is same. & I/O time is much faster than slower (CPU cycle).

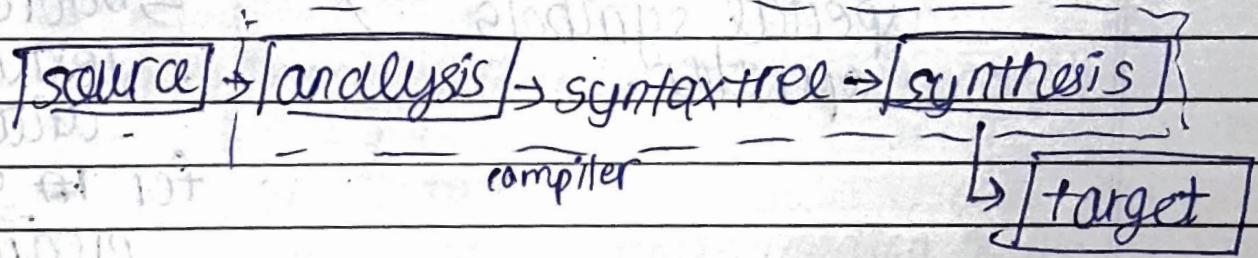
→ Analysis-synthesis model of compil^n.

- 2 parts to compil^n

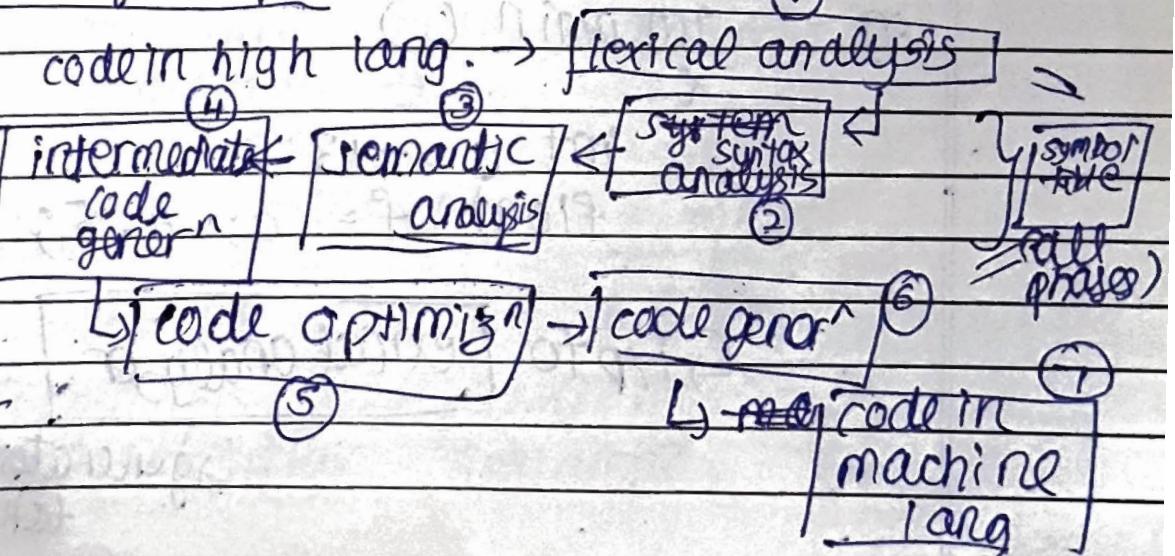
- ① - analysis → check if ip program correct or no
 eg-additn, -, etc
 - understand source program as per program lang used as source program
 - determines oper's implied by the source program which are recorded in tree structure

② - synthesis

- takes the tree structure & converts oper's in target language.



- Phases of compiler



tools using analysis-synthesis model :-

- editors, printers, interpreters,
- + text formatters etc.

Date: / / 120

Page No.:

eg - `int a;`
`float b, c;`
`a = b % c;`

binary ← ↗ semantic error
operand. C % needs int
not float).

⇒ Types of tokens

- comment
- identifier
- keywords
- constants
- special symbols
- operator

} diff. for
diff lang.
↳ token analysis
↳ is done by
lexical
analysis
for ~~the~~ source
program

eg - `int main ()`

{

`int a = 2 + 3;`

`float f = a * 4.5; //hello }`

↳ I/P to lexical analyzer.

Generates
tokens

$\Rightarrow \text{datatype(DT)}$
 $\Rightarrow \text{keyword, ID}()$ ✓

lexical o/p \Rightarrow DT ID() {
 \quad DT ID = const. op. const.; ~~const.~~ } \Rightarrow opto next phase.
 \quad DT ID = ID op. ID; 3
 \quad const. const.

newline & comments are discarded

\rightarrow Now syntax analysis checks this o/p.
 It only understands tokens
 eg- id=id+id ✓ understands
 ↳ not id=a+id
 ↳ X

a=b+c+d can be checked using

$\Rightarrow E + T$
 express' term

↳ binary oper. needing 2 ips

\rightarrow Now we do semantic analysis.

symbol table stores name of id,
 its datatype, etc.

eg- DT ID = const. op. const

↳ details in symbol table

* * : So 3 analysis took place

- lexical, syntax, semantic.

→ maintained

for this eg.

→ Symbol Table

eg - int main()

{ int a = 2 + 3;

symbol name	Data type	mem loc	value type	float f = a / 4.5;
main				func
a	int			ID → DT ID () {
f	float			ID DT ID = const. op. const. DT ID = ID op. const.

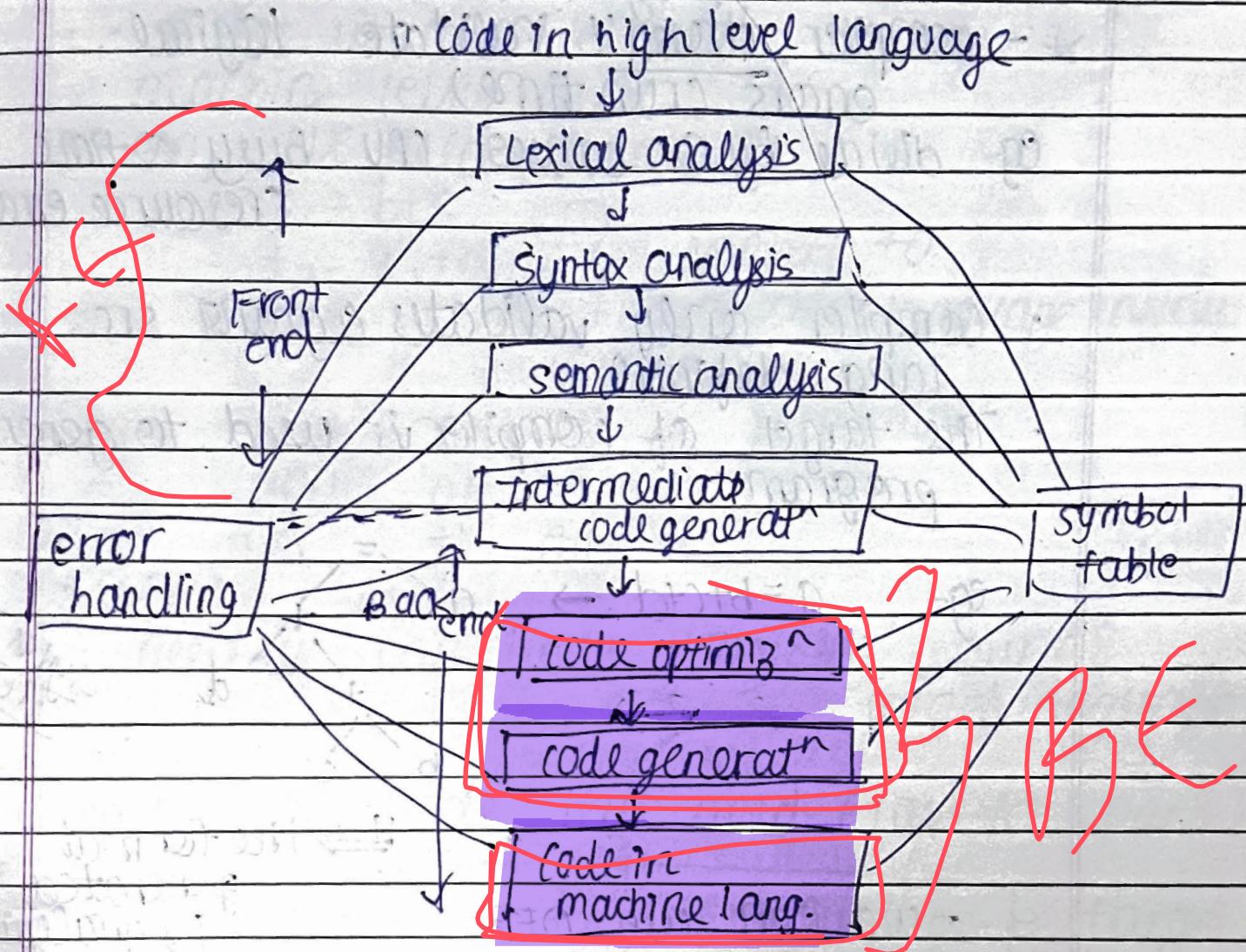
(DT+ID)

↳ through
multiple
datatype
end
already
integer

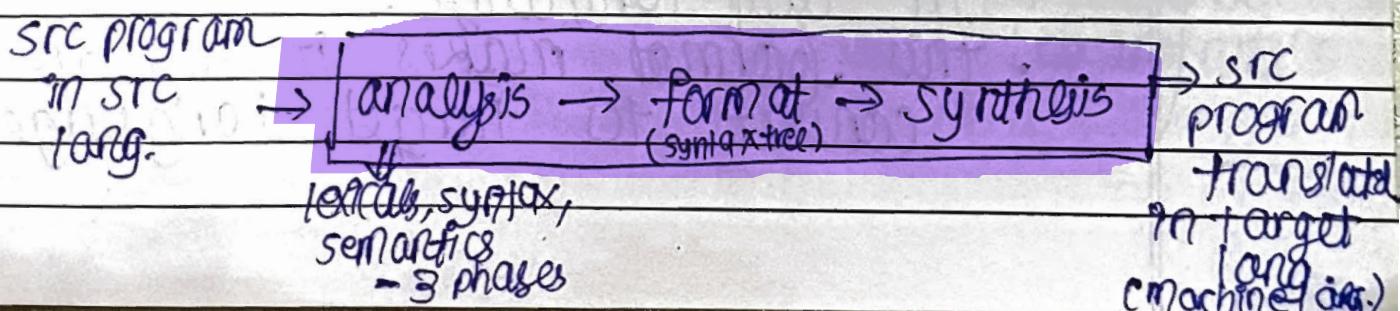
→ Now, for next phases, we can check datatypes in symbol table.

redrawn

Phases of compiler



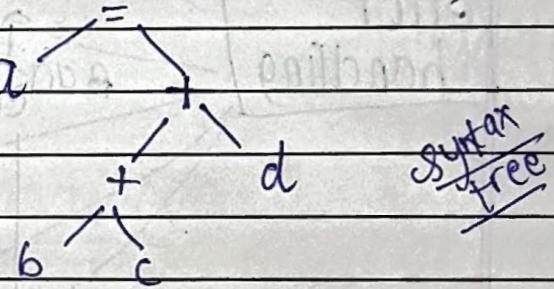
Compiler process



- Analysis - understand & validate source program
 - * - compiler doesn't validate logical errors (run time)
 - eg- divide by 0 causes CPU busy @ time (resource error)
- compiler only validates against src. lang. defn?
- The target of compiler is used to generate program

eg-

$$a = b + c + d \rightarrow$$



\Rightarrow the formal generated by analysis phase.

- we can use DFS for this.

so first of all, $t_1 = b + c$ is processed.
 then $t_2 = t_1 + d$, then $a = t_2$

- Tree makes it easy to write code in any language.

- this format makes it easy to convert to target language.

synthesis

- synthesis means gen?

eg- $a = t_2$ is the last step of synthesis

- But it is not enough, as we need machine level code.

so, we need to convert -

$$t_1 = b + c$$

$$t_2 = t_1 + d$$

$$a = t_2$$

convert to

machine language instruc?

eg-

MOV AH, b

ADD AH, C

Mov TI, AH

optimiz^n

- need of compiler - knowing source & target language

so, we can divide this need (task to 2 parts)

src → [analysis → format → synthesis] → target

1st need

2nd need.

- compiler has 2 jobs - transl & optimiz
↳ making less instruc's



eg- $a = a + 1$ can be done using \Rightarrow

```
mov AH, a
Add AH, #1
mov al, AH
```

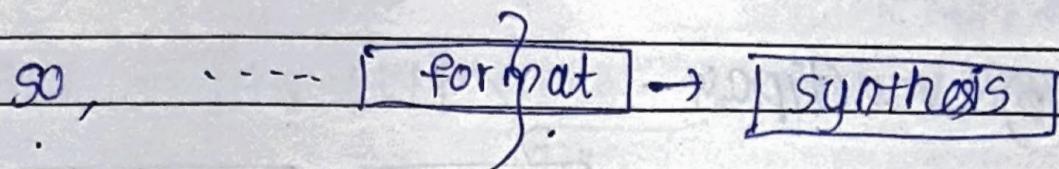
\hookrightarrow optimized \Rightarrow mov AH, a \Rightarrow can only be implemented if machine architecture has the function (current available)

\hookrightarrow increment
 \Rightarrow This is machine dependant optimizn.

2 optimizn- before machine & after machine code generated

machine indep.
use less instructions.
 \downarrow
machine dependant.

eg- $a = a + 1$
 $b = a \rightarrow$ useless, so remove before
 $c = a + 2$ transn
into machine code.



intermediate → optimizⁿ → code → op
code general?

* → optimizⁿ is an optional phase of compiler

• phases of compiler can be seen as functions.

eg - let's say 2 approaches

1) Functⁿ for each analysis phase

- pipeline - which while lexical generates each token, it keeps on passing it to next phases.

- single pass

↳ 1 place where functⁿ defined

mycompiler.c

lexer

syntax()

main()

code

optimizⁿ(c)

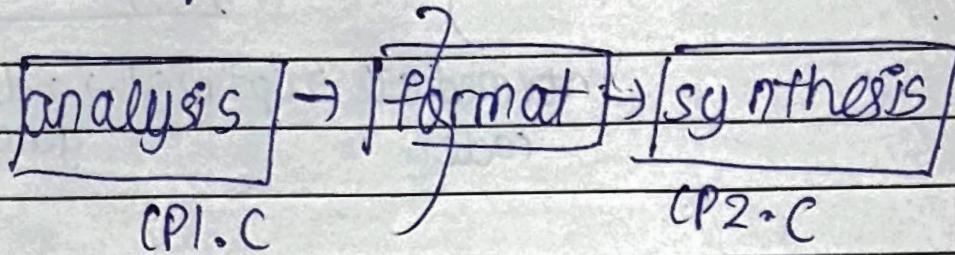
↳ --

- Phase is functionality (sequence)
- phase clubbing all functionality in I
- pass

Date : / / 20

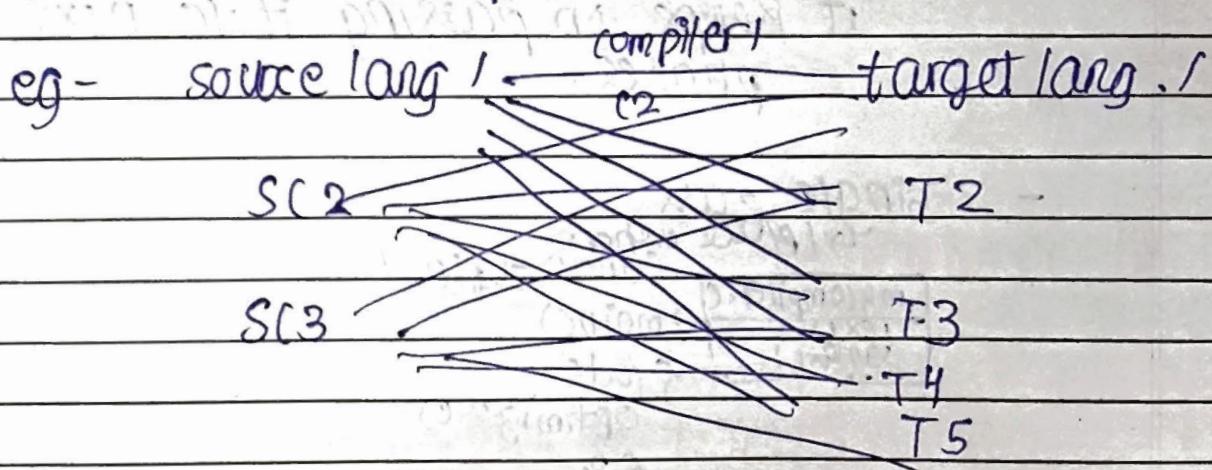
Page No.:

2) multipass -



more

- multiple time (scanning) needed for multipass
- more RAM in single pass
- more ROM in multipass to store.
- multipass doesn't generate code if program wrong, single pass can give intermediate o/p, as o/p of previous phase may get passed ahead.

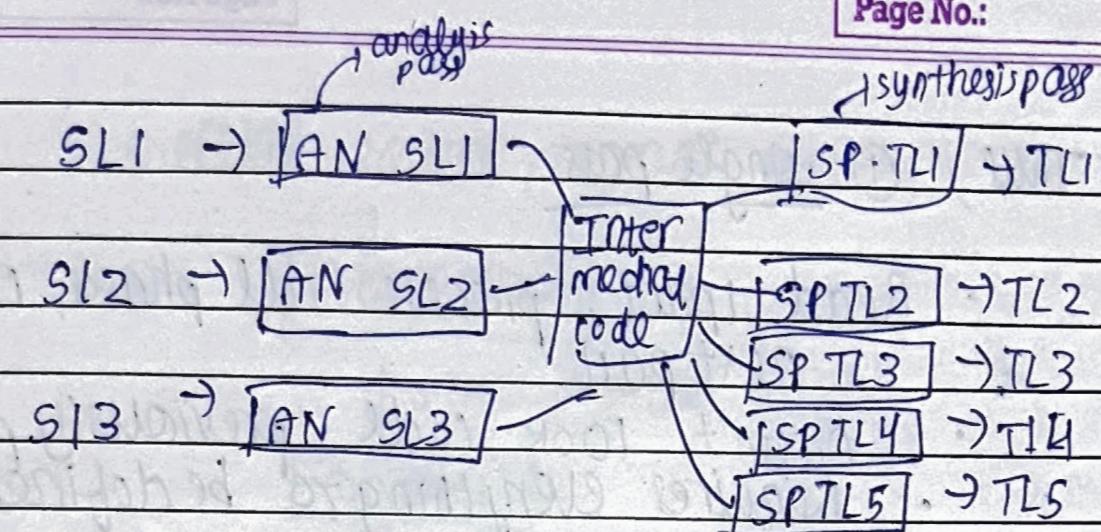


(single pass compiler)

$$\Rightarrow 3 \times 5 = 15$$

compilers
needed.

- For every compiler, we need 1 analysis, 1 synthesis
- ↳ total 2 for 1 compiler
- ↳ total 4 for 2 compilers
- ↳ total 6 for 3 compilers



(multipass compiler)

$3+5 = 8$ pass

reduct
using reusability

→ Frontend - who interacts with the user
 ↳ eg - ~~whether how~~ is ~~scanf~~ designed

→ Backend - functionality logic (not interacting with user)
 eg - C is both front & backend.

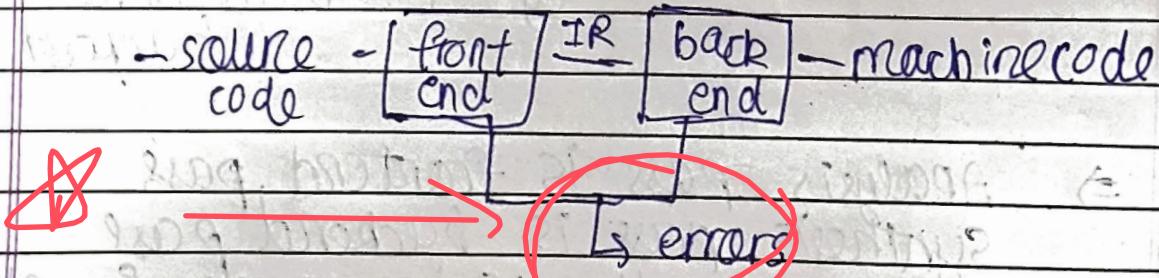
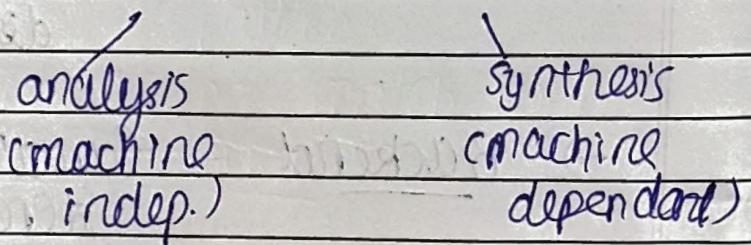
⇒ Analysis pass is frontend pass
 synthesis pass is backend pass
 • It is preferable to separate frontend & backend pass.

Thus, ① single pass

- Read 1 part, process all phases, read next part
- Doesn't look code previously processed
- Requires everything to be defined before
- Require large memory.

② multipass

- every pass results new representation & rip to next pass
- compiler -> frontend - backend



- similar to compilers -> preprocessors, assembly, linker (phases of compiler)

skeletal source program

preprocessor

↓ source program

compiler

↓ target assembly program

assembler

↓ relocatable object code

linker

← libraries & reloc. obj. files

↓
absolute
machine
code

→ applic's of compiler techniques

- Lexical analyzer - text editors, IRS, pattern recognit'

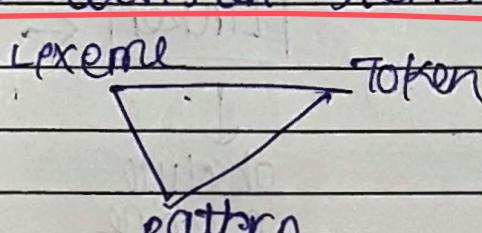
- Syntax analyzer - SQL query processing, Kmap for circuit design

- ~~Syntax + semantic~~ - equation solver analyzer

- NLP systems

- - Lexical analyzer doesn't give comments & errors - only gives other tokens.
 - It tries to match patterns of the substring (considering entire code as 1 string)
 - It generates tokens - also called classes
id, const., operator etc.
- ⇒ Lexeme - the matched part of the pattern.
Token - the class.

→ So, Lexical analysis works on 3 terminologies -

- 1) Lexeme
 - 2) Token
 - 9 3) pattern
- 

• string / substr. that doesn't match with any pattern - lexical error

* Lexical analyzer role

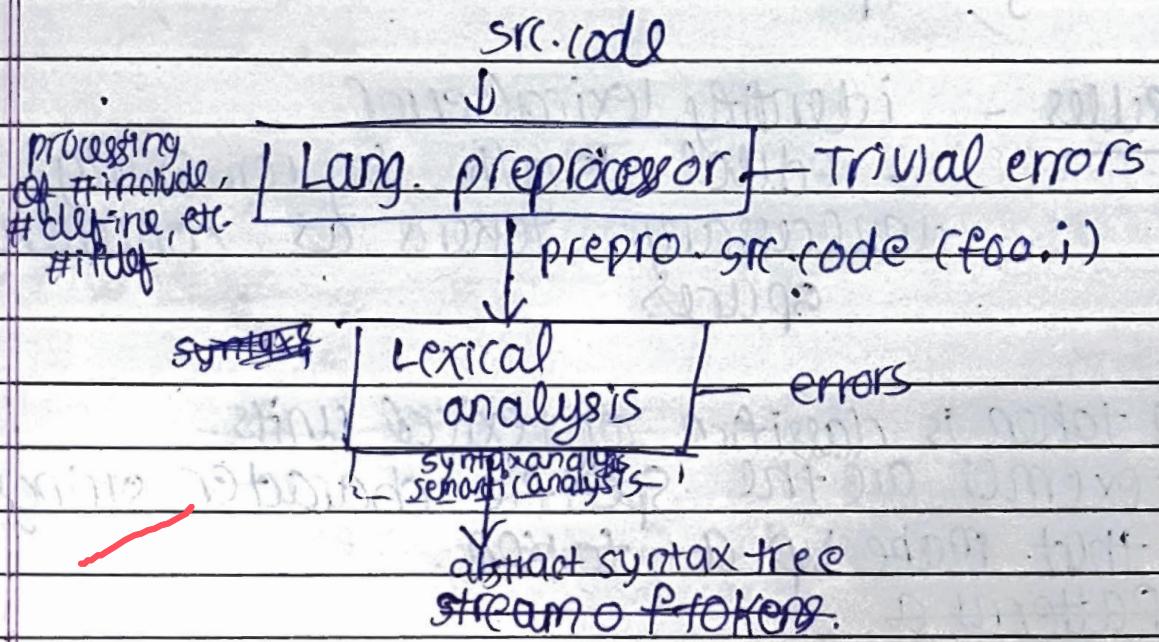
- scan src. program
- translates char stream to token stream
- eliminate unnecessary tokens (eg - comments, spaces)
- error symbols along with locn, type into symbol table (all symbols)
- lexical error identification
↳ but doesn't classify them.

⇒ symbol is a label
Identifiers

Date : / /20

Page No.:

- Lexical analyzer works only at 1 token at a time
 - To classify token we need previous & next token.

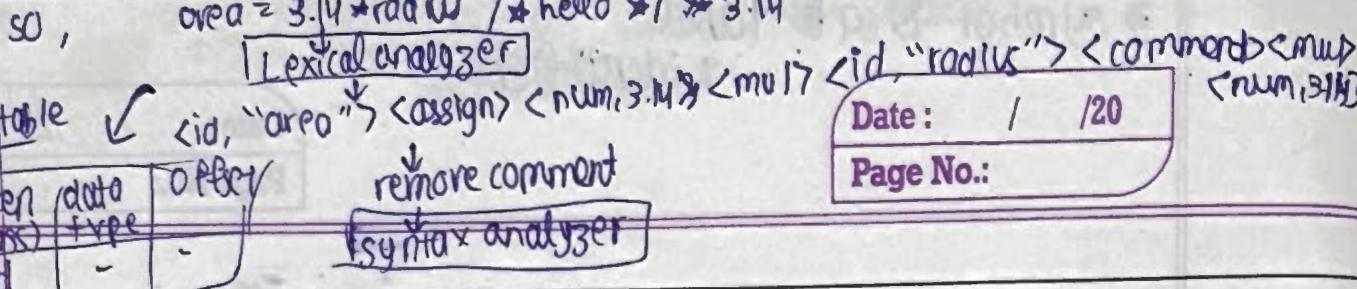


eg - " area = 3.14 * radius * radius ; @ 1 * eq?*/if..

\Rightarrow area = token cid). -symboltable ✓

radius = $\pi d - \checkmark$
→ formula → take attribute like π example

OP <id, "area"> <assign> <num>
 ^c taken <multiply> <id, "radius">
comment class <mul><id, "radius"><;>
ignored
ignores
@ ⇒ error as invalid



- If we execute this as pass, then we get o/p. of intermed. successful tokens.
- If taken as 1 phase, then it doesn't give o/p.

Rules -

- identify lexical error
- reduce length by removing unnecessary tokens as comments, spaces.

- A token is classifier of lexical units.
- Lexemes are the specific character strings that make up a token.
- Patterns a

⇒

test.c → comp → .test.exe
(src. program)

abc.l → Flex → .c file

eg. float const. RegEx

↳ float "[0-9]*[.][0-9]+

→ 4.25 ✓

[0-9]*"."[0-9]+

• 25 ✓

A. X

eg → Preprocessor statements

"#" "<include>"

^ "#" • *

→ multiline comment

• "/*" • * " */"

"/" /* ~~*/~~ . *

/* . * " \n " . * */

"/" ((1*)|"\n")* (^/) /* */

eg - if as a keyword \Rightarrow if()

if as a variable \Rightarrow int if;

\Rightarrow "if" / "C" { keyword }

"if" { ID }



→ Rules to solve matching conflicts

- Longest matching pattern wins.
- Ties in length resolved by priorities
→ token specific order often defines priority.

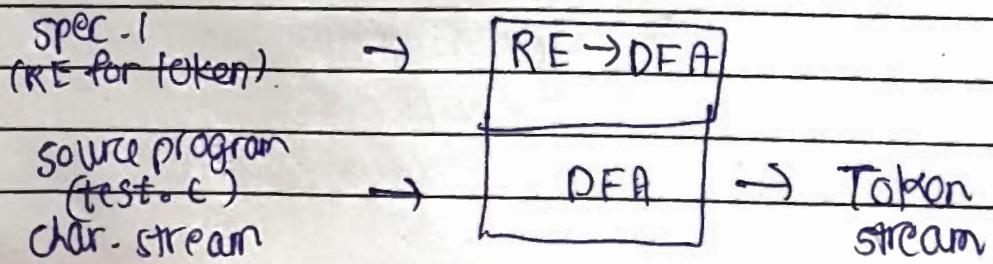
* → Design of lexical analyzer generator

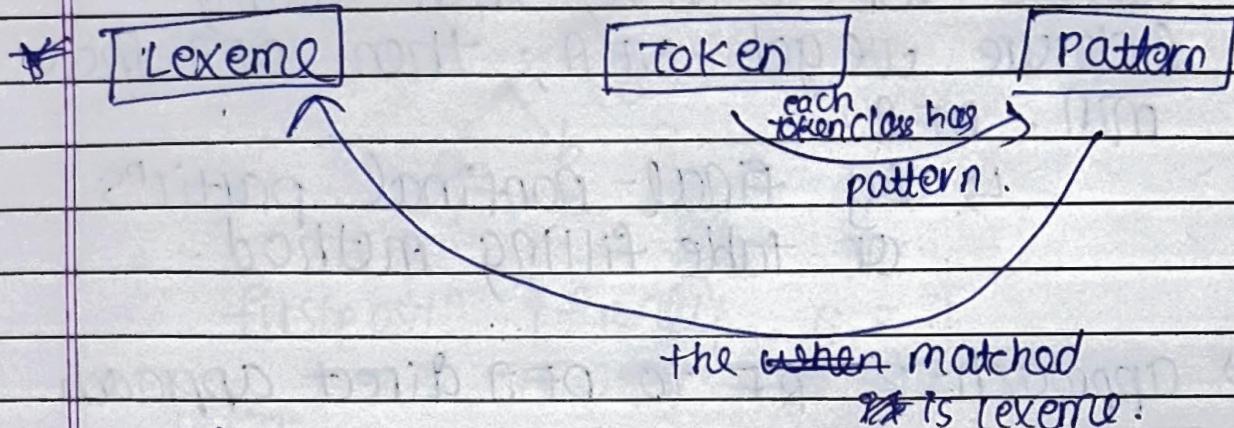
- 1) Approach to operate & recognize token using DFA
- 2) Input buffering
- 3) Optimizn using sentinel.

Lab

- The string token is stored in `yytext`.
- To store an id in symbol table, we can add an + we a self defined function and

→ Lexical analyzer gener. (Lex/Flex) working :-

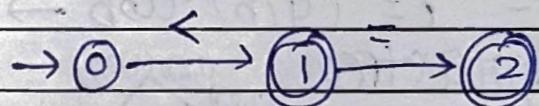


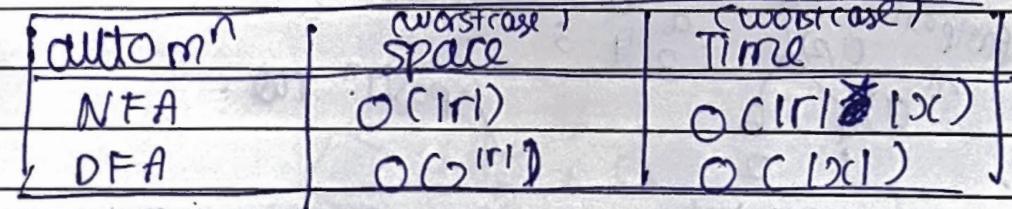
class

we can

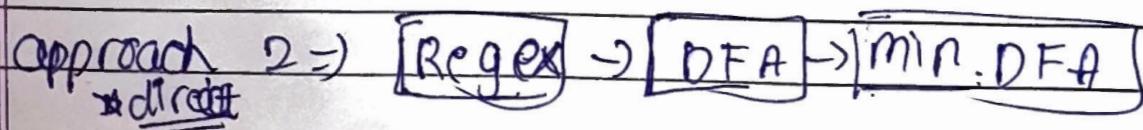
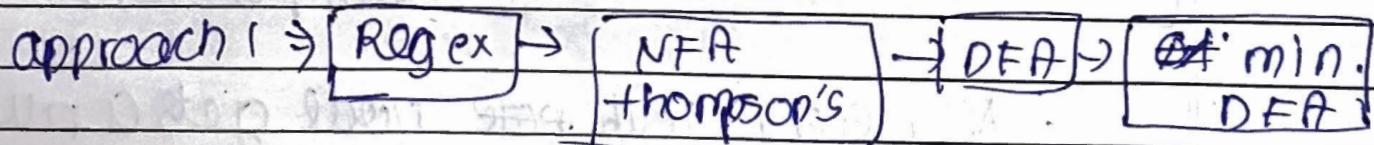
- use Finite automata to check regex match in a loop

eg - check <=



* 

autom^n	worstcase space	worstcase Time
NFA	$O(1^n)$	$O(1^n \cdot 1^n)$
DFA	$O(2^n)$	$O(1^n \cdot 1^n)$

→ Regex to DFA

→ For appr. 1, we can use thompson construct closure for e-NFA, then using closure we get NFA, then DFA then min-DFA.

↳ by final-nonfinal partitions or table filling method.

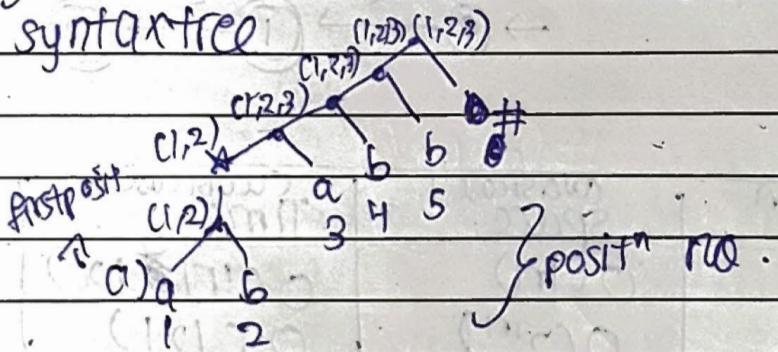
→ approach-2 RE to DFA direct approach

steps

1) append RE r with # to make accepting state

$$\text{eg} - r = (a/b)^*abb \rightarrow (a/b)^*abb\#$$

2) syntax tree



eg-
class $(a/b)^*a$

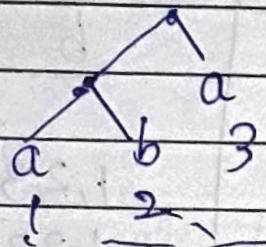
1 2 3

⇒ first positⁿ can
 be 1, 2, 3.
 first pos (1,2,3)

Now every leaf node gets a number.

3) For every node calculate firstposiⁿ.

eg -

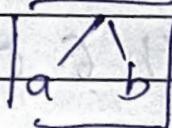


firstposiⁿ if only $a = 1$

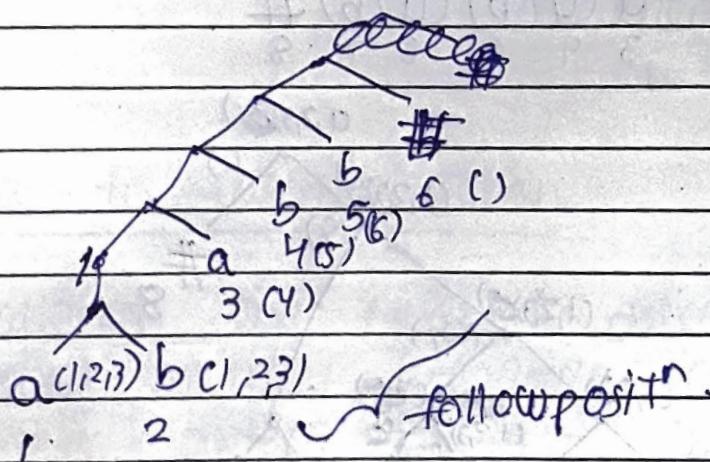
if only $-b = 2$

if

$\xrightarrow{\quad}$ firstposiⁿ = (1, 2)



- In our main example, we see the positiⁿs.
- similarly, followposiⁿ is calculated.



∴ we now make a transition table :

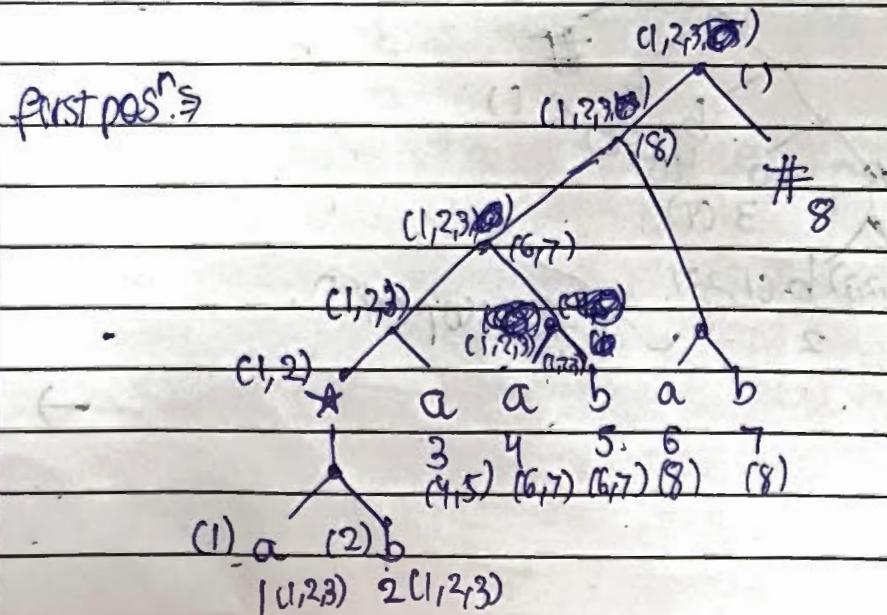
start	a	b	⇒ states having # are final states -
state name ↪ A = {1, 2, 3, 4}	1, 2, 3, 4	1, 2, 3	
B = {1, 2, 3, 4}	B	1, 2, 3, 5 ↳ (C)	
C = {1, 2, 3, 5}	B	1, 2, 3, 6 ↳ (D)	
D = {1, 2, 3, 6}	B	A	

• # ⇒ 6 → states with 6 are final states

• Now we can minimize.

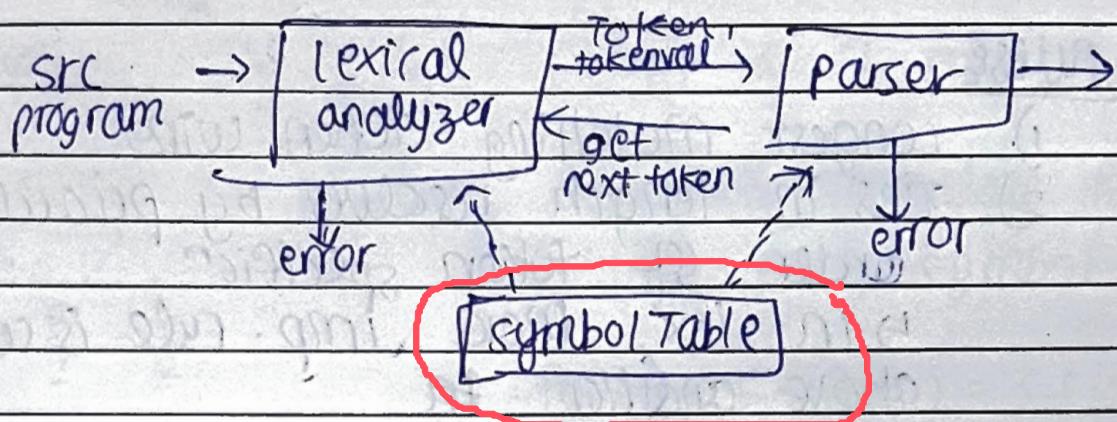
$$\Rightarrow \text{② } (q/b)^* a c a/b) (q/b)$$

$$\begin{matrix} \Rightarrow (q/b)^* & a & (q/b) & (a/b) & \# \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \end{matrix}$$



Self Notes

→ Lexical analysis



- Token is a classification of lexical units
 ↳ a class
- Lexemes - specific char. string that make up a token.
 ↳ the matched part becomes a lexeme

→ How to describe Token

- Patterns are 'rules to describe set of lexemes belonging to a token'
 eg - pattern for id :- "letter followed by digits"



• To resolve matching conflicts like :-

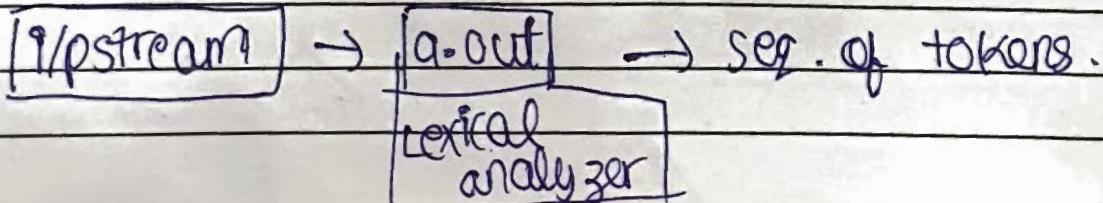
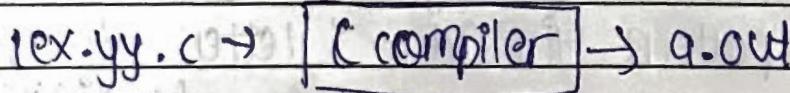
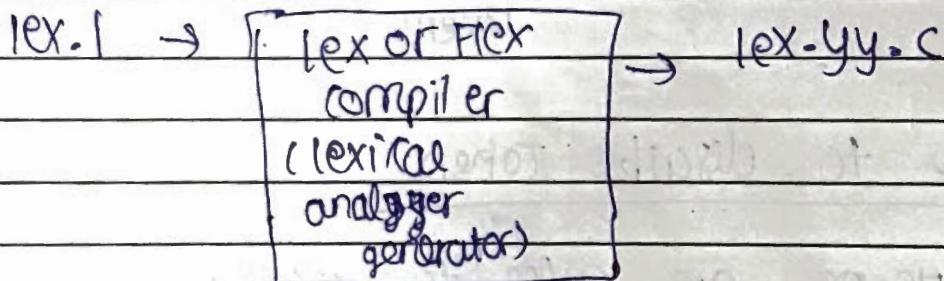
eg- `elsex = 0 ;`

- ① `else, x, =, 0, ;`
- ② `elsex, =, 0, ;`

Rule

- 1) longest matching token wins
- 2) Ties in length resolved by priority order of token specification
↳ in lex, more imp. rule is written above another
- 3) REGEX + priorities + longest matching token rule = definition of Lex

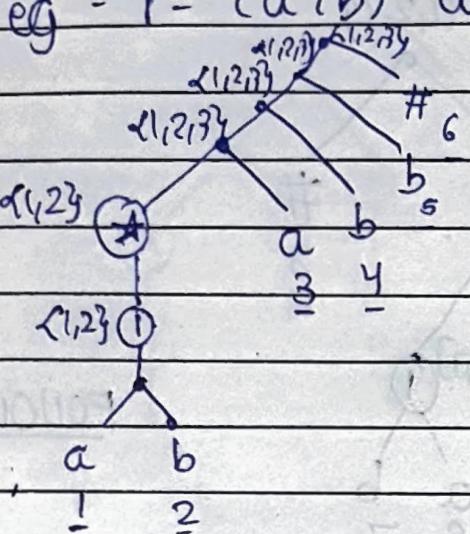
→ Lex working



* for \emptyset or (1)
 $\text{firstpost}^r = \text{firstpost}^n(L) \cup \text{firstpost}^n(R)$
 if nullable
 if yes
 $\text{firstpost}^r(L) \cup \text{firstpost}^n(R)$
 $\text{for } \star, \text{ firstpost}^r$
 firstpost^r
 of child
 if not nullable
 $\text{firstpost}^n(R)$
 of child

eg - RE to DFA

eg - $r = ca(b)^*abb \Rightarrow (a/b)^* \cdot a \cdot b \cdot b \cdot \#$



The {1...3} are firstpost^r

Now calculate followpostⁿ of each symbol.

followpostⁿ - $fp(1) = \{1, 2, 3\}$

$fp(2) = \{1, 2, 3\}$

$fp(3) = \{4\}$

$fp(4) = \{5\}$

$fp(5) = \{6\}$

$fp(6) = \emptyset$

Now table \rightarrow take firstpost^r of root as 1st state

State	\leftarrow symbol \rightarrow	a	b
$A = \{1, 2, 3\}$	$\{1, 2, 3, 4\}$	$\{1, 2, 3\} \rightarrow A$	
$B = \{1, 2, 3, 4\}$	B	$\{1, 2, 3, 5\} \rightarrow C$	
$C = \{1, 2, 3, 5\}$	B	$\{1, 2, 3, 6\} \rightarrow D$	
$D = \{1, 2, 3, 6\}$	B	A	

Now, we use the ~~FP~~ of each non-terminal state

herein $A = \{1, 2, 3\}$

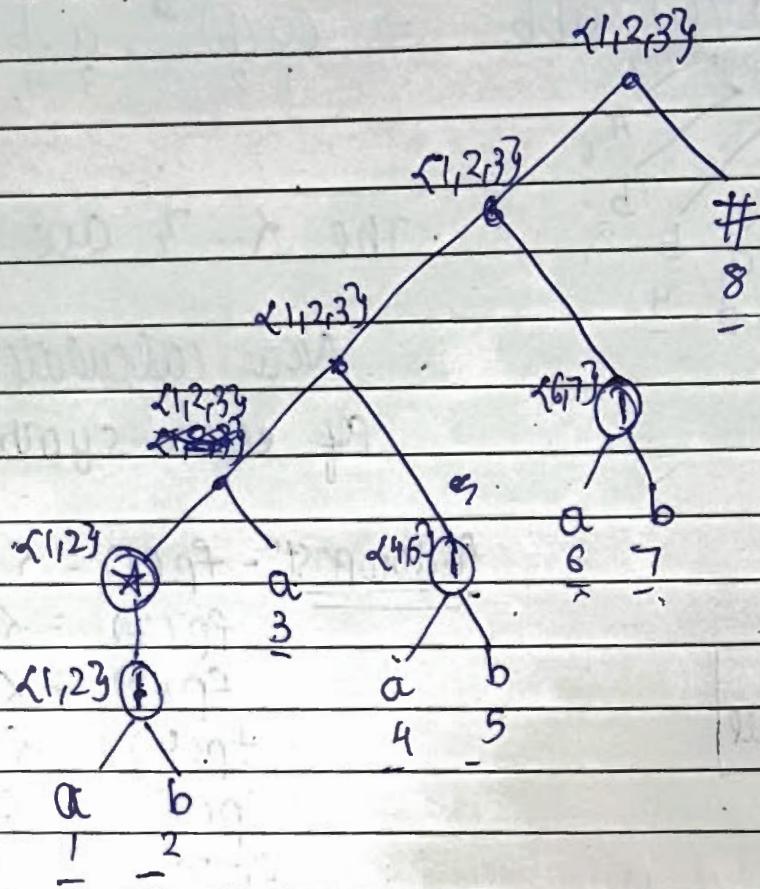
1 & 3 have same symbol

\hookrightarrow so union them

\rightarrow states having # (6) are final
so, D is final

Now we can minimize this DFA

$$\text{eg} - (a/b)^* a (a/b) (a/b) \Rightarrow \underset{1}{(a/b)}^* \cdot \underset{2}{a} \cdot \underset{3}{(a/b)} \cdot \underset{4}{(a/b)} \cdot \underset{5}{(a/b)} \cdot \underset{6}{a} \cdot \underset{7}{(a/b)} \cdot \underset{8}{(a/b)} \#$$



Followuposit's (FPJ)

$$f_D(1) = \{1, 2, 3\}$$

$$fp(2) = \{1, 2, 3\}$$

$$f_D(3) = \{4, 5\}$$

$$f_D(4) = 4.674$$

$$f(0.5) = 4.6 \cdot 10^{-4}$$

$f'(6) = 607$

$$\text{EP}(7) = 84$$

$$fp(8) = \emptyset$$

Table

	← symbols →	
STATE	a	b
A = {1, 2, 3}	{1, 2, 3, 4, 5}	{1, 2, 3}
B = {1, 2, 3, 4, 5}	{1, 2, 3, 4, 5, 6}	{1, 2, 3, 6}
C = {1, 2, 3, 4, 5, 6, 7}	{1, 2, 3, 4, 5, 6, 7, 8}	{1, 2, 3, 6, 7, 8}
D = {1, 2, 3, 6, 7}	{1, 2, 3, 4, 5, 8}	{1, 2, 3, 5}
E = {1, 2, 3, 4, 5, 6, 7, 8}	{1, 2, 3, 4, 5, 6, 8}	{1, 2, 3, 6, 7, 8}
F = {1, 2, 3, 6, 7, 8}	{1, 2, 3, 4, 5, 8}	{1, 2, 3, 4, 8}
G = {1, 2, 3, 4, 5, 8}	{1, 2, 3, 4, 5, 6, 9}	{1, 2, 3, 6, 7}
H = {1, 2, 3, 8}	{1, 2, 3, 4, 5}	{1, 2, 3, 7}
=		

} now complete.

states having 8 (#)
are final

eg - $(0/1)^* \cdot 10^* 10^* \rightarrow (0/1)^* 10^* 10^* \#$

follow posn

$fp(1) = \{1, 2, 3\}$
 $fp(2) = \{1, 2, 3, 4\}$
 $fp(3) = \{4, 5, 6\}$
 $fp(4) = \{5, 6\}$
 $fp(5) = \{6, 7\}$
 $fp(6) = \{6, 7\}$
 $fp(7) = \emptyset$

Table

	0	1
$A = \{1, 2, 3\}$	$\{1, 2, 3, 7\}$	$\{1, 2, 3, 4, 5\}$
$B = \{1, 2, 3, 4, 5\}$	$\{1, 2, 3, 4, 5\}$	$\{1, 2, 3, 6, 7\}$
$C = \{1, 2, 3, 6, 7\}$	$\{1, 2, 3, 6, 7\}$	$\{1, 2, 3, 4, 5\}$
<u>Final</u>		

- Ques -
- when we compile a C program, .obj file is generated
 - compiler never generates .exe

- lexeme = a valid token
 ↳ collect of character according to lang of the

→ Why are only identifiers put in symbol table

- only value of ids change

eg - int a,b,c; ,

↳
 ↳ <kw1, int> → int
 <idl, ptr. to symbol table> → a
 <id2, ptr>
 <id3, ptr>

ptr → id1	type	mem-lcm	scope

lexical error → [int a;]

↳ not a lexical error
 it will be identified as an identifier

- Similarly, $f_i(a>b)$ \Rightarrow not an lexical errors

\rightarrow also called scanner - re-scanning L to R

=* Lexical analyzer is the only phase reading full source program

↳ Lexical analyzer will attach the line no. with code.

eg - Line 9 - invalid token.

~~Cousins of Computer~~ → # include <stdio.h> \rightarrow to use inbuilt funcn.

define PI 3.14 \rightarrow inline / lambda
void main()

{
 printf();
 scanf();
}

- #include allows funcn's to be used.
- Lambda / inline / macro - replaces the value

↳ eg - PI 3.14

↳ replace it (replace PI by 3.14)

↳ speed up process

↳ this is expansion of source program

- This is role of preprocessor.

* ~~Studio~~ - only has found ~~name &~~, not the code.

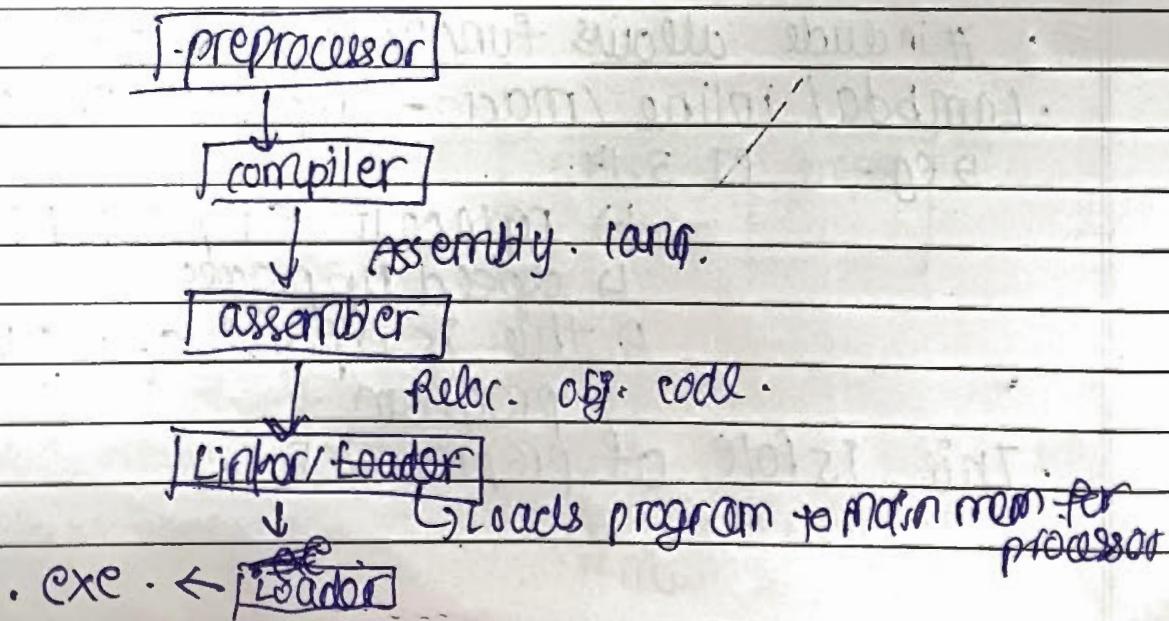
Date: / /20

Page No.:

- Now compiler will get this expanded program.
- Compiler generates assembly language.
 - eg - ADD STA BUN
- Then assembler generates relocatable machine obj. code.
- Then, linker links the fundⁿ declared ^{for} loader

(eg - scanf & printf) , & then also
Keyboard monitor

links the relocatable obj. code of our program, & obj. code of monitor-printer keyboard.

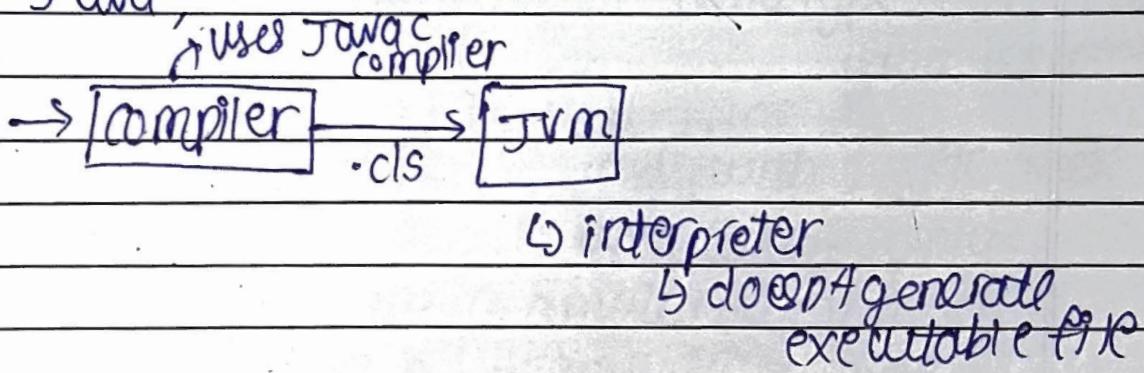


* Java & Python use both compiler & interpreter - so, they don't have .exe file.

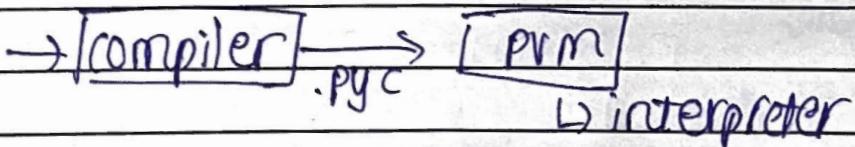
→ Thus, compiler can't function without linker, assembler, loader, preprocessor - so, they are called 'cousins of compiler'.

- * If we compile & run, entire all the phases are run
- If we only run, then only Linker & Loader will run

→ In Java,



→ In Python,



· So, Java & Python use both compiler & interpreter.

LabPractical-2

- For a grammar, check if string present or not.

eg - $S \rightarrow ASI / 0A / 0A$
 $A \rightarrow 011$

0011 \Rightarrow ✓

S has 3 fuctions.

- Take 3 variable grammar.

class

class → Syntax analyzer

eg- void main()

{ int a,b,c;

 int int;

}

)

not a lexical

error.

→ smallest unit of program

= word

↳ verified by
lexical analyzer

→ verify word to check
program

- once verified, words generate sentences
- syntax analyzer checks syntactic str. of the statement
- So, in C, we have \Rightarrow datatype variable;
 every line ends with ; etc. these are rules of syntax.

\Rightarrow These rules are grammar of the lang.

- We have 4 types of grammar - type 0, 1, 2, 3

free context context
unrestricted sensitive free

regular
grammar

• Here, we only use context-free grammar

\hookrightarrow used for by
syntax analyzer.

so, syntax analyzer works with CFG
as CFG works with grammar which can be verified.

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

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

variables ↗ start symbol
 ↗ production rule
 ↘ terminals

The grammar uses process of derivation

* Derivⁿ - is process of deriving a string from starting symbol.

Parse tree - A graphical representation of derivⁿ process using tree structure

- ↳ root node at start (start sym.)
- ↳ intermediate vertex - nonterminals
- ↳ leaf node - terminals.

$S \rightarrow DL$
 ↳ ~~product~~

$S \rightarrow DL$
 ↳ derivation

* use ' | ' for separation

$\Rightarrow D \rightarrow \text{int} | \text{float}$

$$\cdot S \rightarrow DL;$$

grammar $\Rightarrow D \rightarrow int \mid float$

$$L \rightarrow L, L \mid R$$

$$R \rightarrow a \mid b \dots$$

Derivⁿ \Rightarrow

$$\cdot S \rightarrow DL$$

$$\Rightarrow int \ L, L;$$

$$\Rightarrow int \ a, L, L;$$

$$\Rightarrow int \ a, b, c;$$

sentential form

Derivⁿ

sentence

\rightarrow sentence - It is a collectⁿ of terminal symbols during the derivⁿ process.

so, if $S \xrightarrow{*} G$ meaning if using
($*$ means 0 or more)

0 or more steps if from S we generate a string G only of ~~variables~~^{terminals}, then we get a sentence.

\rightarrow sentential form - Derivⁿ steps that may contain nonterminals is known as sentential form.

So, $S \Rightarrow DL;$

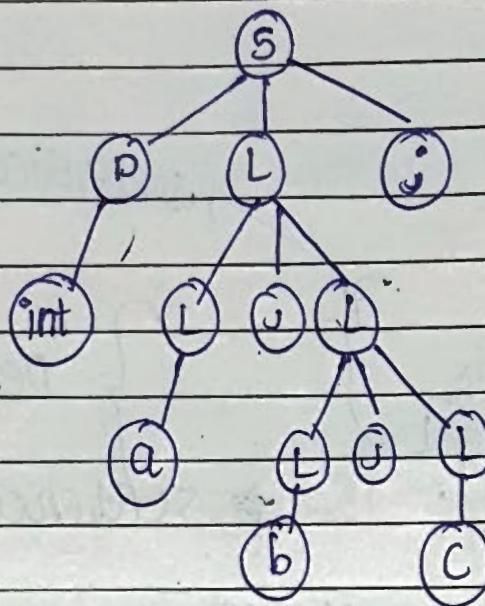
$$\Rightarrow int \ L, L;$$

$$\Rightarrow int \ a, L, L;$$

sentential

sentential form





⇒ This is a parser tree

now, if $\Rightarrow \text{int } a, b, c \& \text{ no semicolon}$,
1st product itself won't execute

If $\boxed{a, a ;}$ then DL ,
↳ not present.

⇒ 1st type of error generated is called
misspelled keyword error

or
⇒ If no $;$ or $,$ then we get
missing / invalid punctuation marks.

⇒ operator missing error can also be generated

3 errors
gen.
bit



~~Note - Only lexical & syntax analyze phase can update symbol table~~

Date : / / 20

Page No.:

• if \Rightarrow int int ;

Lexical analyzer sends them to syntax analyzer, but syntax error - invalid identifier.

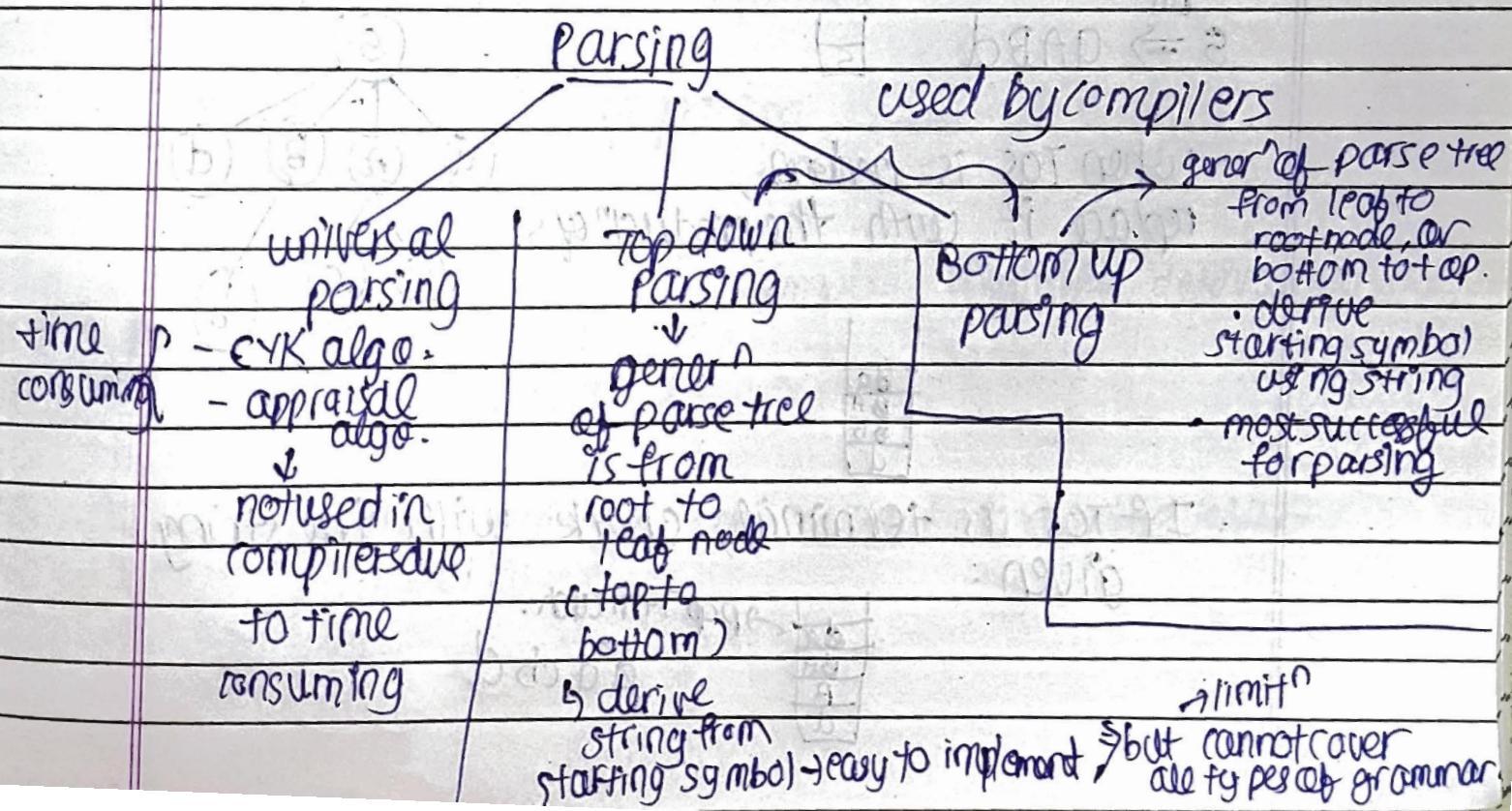
• so, if \Rightarrow int a, b, c ;

lexical enters \Rightarrow a, b, & c in symbol table
syntax enters \Rightarrow int datatype of a, b, c

• syntax analyzer also scans data LTR.

→ Parsing - deriving the string using the parse tree.

• There are 3 types of parsing



- Left most derivⁿ = The leftmost variable is derived first
- Right most = vice versa.

~~→ Top down parsing uses left most derivⁿ, & bottom up parsing uses right most derivⁿ~~

~~in reverse.~~

use of stack →

$$\begin{aligned} S &\rightarrow QABd \\ A &\rightarrow QC \\ B &\rightarrow b \end{aligned}$$

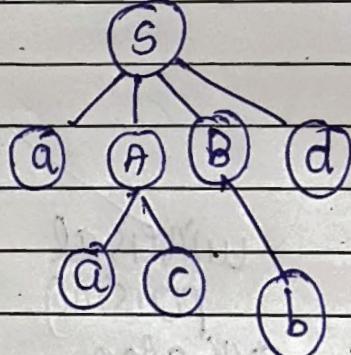
string = aacb d

For top down parsing,

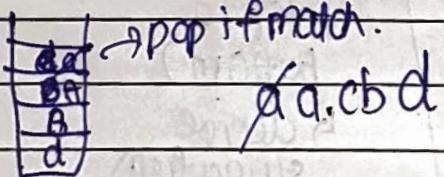
push S in stack

$$S \xrightarrow{^m} QABd \quad \boxed{S}$$

When TOS is nonterm.,
replace it with the products of S



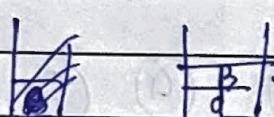
If TOS is terminal, check with the string given.



now, as TOS = A \Rightarrow nonterm., expand it.
 replace it with product = abc
 $A \xrightarrow{lm} abc$



abc'd



replace B by b

$$B \xrightarrow{lm} b$$



bcd

\Rightarrow stack & string over.
 • So, string is a member

\Rightarrow For bottom up parsing,

acbcd

$$\begin{array}{l} S \xrightarrow{rm} aABd \\ \xrightarrow{rm} aAbd \\ \xrightarrow{rm} aaacbd \end{array} \quad \left. \begin{array}{l} aABd \\ aAbd \\ aaacbd \end{array} \right\} \text{simple rightmost derivation}$$

- While parsing, we want to derive start symbol of string





a not there
on right side of
product's

(a)

I have got
a mother
on right
side is a handle
it is called a
Handle



a not there
on right side
of product's

(a)

(a)



ac present of
right side -

(a)

(a)

(c)

we are reducing
so :-

(a)

(a)

(c)

→ handle



a A not
on right side

(a)

(a)

(c)

(c)



a A b
b on right
side

(a)

(A)

(a)

(c)

(b)

↳

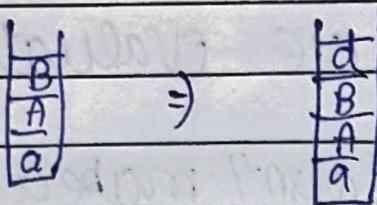
(a)

(A)

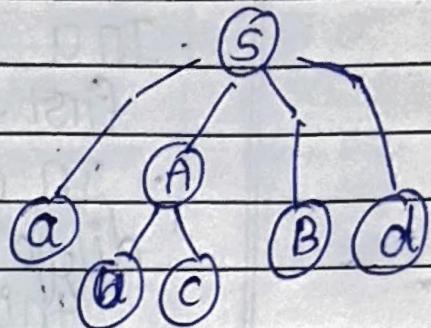
(a)

(c)

(B)



aABd present
on right side



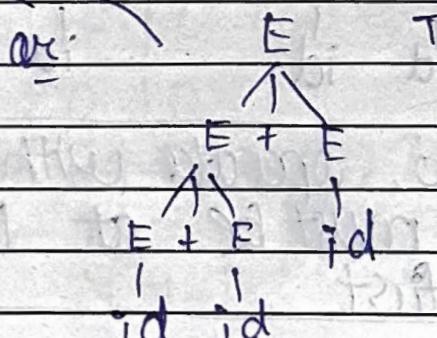
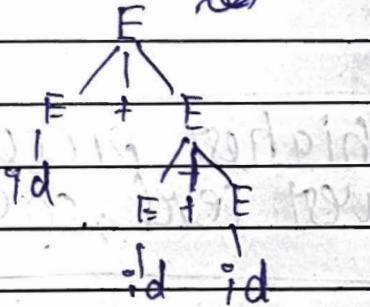
$\boxed{s} \Rightarrow$ we reduced entire
thing by S.

We used this as \Rightarrow reverse of rightmost
derivation - as S was derived in the end,

L) eg- in ~~left~~ rightmost,
B was derived
first, but here
B was derived last.

class $\rightarrow E \rightarrow E+E|id$
derive id fid+id \Rightarrow eg- 2+3ts

Here
associativity
is the problem



This is correct as '+'
is left associative
- id+id+id
solve first

$$E \Rightarrow E+E$$

$$E \Rightarrow id + E+E$$

$$E \Rightarrow id+id+id$$

$$E \Rightarrow E+E$$

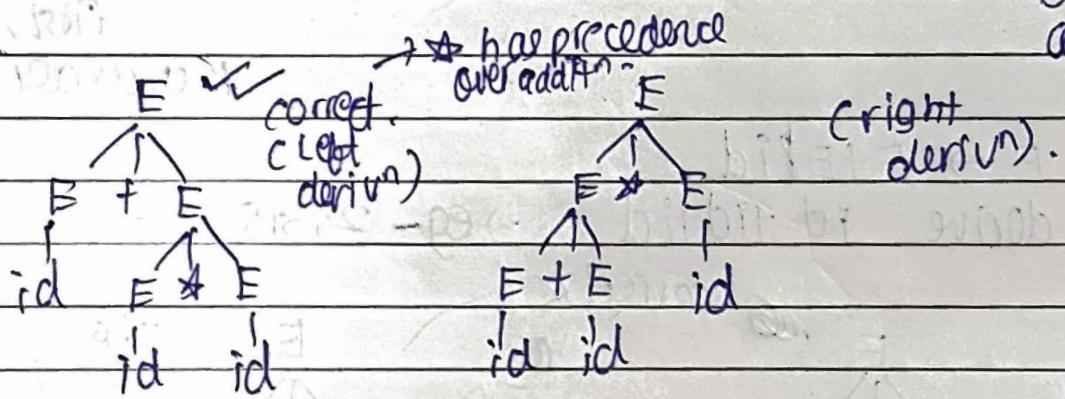
$$E \Rightarrow E+E+id$$

$$E \Rightarrow id+id+id$$

- In a tree, the leaf nodes are evaluated first.
- In addition to this doesn't make a difference, but if 2 power 3 powers, this makes diff grammar answers for both trees.
- This grammar is ambiguous
 - more than 1 tree.

e.g. $E \rightarrow E+E \mid E * E \mid id$
 $id \mid id * id$

here, precedence is not taken care of



So, operator with highest precedence must be at lowest level; so evaluated first.

→ How to make an ambiguous grammar an unambiguous

Steps to make unambiguous grammar

so, if we take care of associativity, precedence, we can make the grammar

If we make grammar ambiguous, if operator is left associative, introduce a new nonterminal or right side, if right asso., then new terminal on left side.

Date: / / 20

Page No.:

→ write first as lowest precedence.

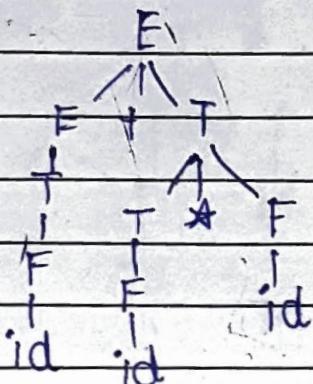
$$\rightarrow : E \rightarrow E + T \quad T \rightarrow T * F \quad F \rightarrow id$$

we can take care of precedence by this rule

eg - id + id * id

we can't write $T * T$ as again ambiguous.

id has highest precedence, hence at lowest level.



left side to be evaluated first.

left associative
a+b+c
↳ evaluate 1st - a-b-c
↳ evaluate first.

Q - $E \rightarrow E+E \quad | \quad E * E \quad | \quad E-E \quad | \quad E/E \quad | \quad E \uparrow E \quad | \quad id.$

~~$E \rightarrow E+E \quad | \quad E * E \quad | \quad E-E \quad | \quad E/E \quad | \quad E \uparrow E \quad | \quad id.$~~

$E \rightarrow E \uparrow A \quad E \rightarrow E \uparrow B \quad E \rightarrow E-B \quad | \quad E+B \quad | \quad B$

$A \rightarrow A \quad A \rightarrow B \quad B \rightarrow A+B \quad | \quad B$

$B \rightarrow C \star \quad B \rightarrow B * C \quad | \quad C$

$C \rightarrow C/D \quad | \quad D$

$D \rightarrow id \quad | \quad D \uparrow E \quad | \quad id$

$E \rightarrow id \quad | \quad N \rightarrow C \quad | \quad N$

$N \rightarrow id$

Q - above Q. with | (E)

brackets



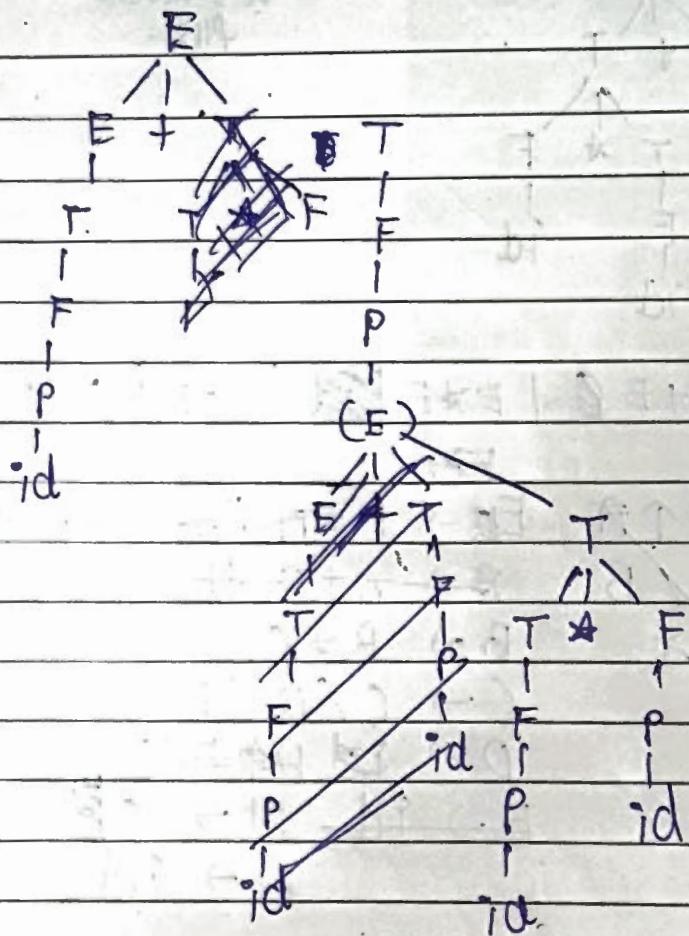
Q, $E \rightarrow E + T \mid T \mid E - T$

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

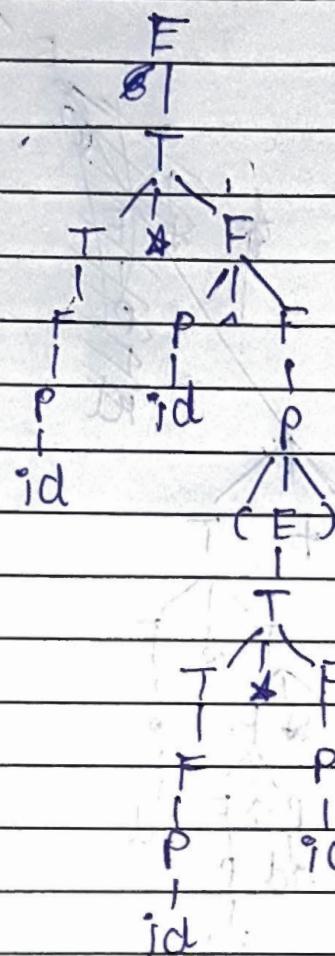
$F \rightarrow P \uparrow F \mid P$

$P \rightarrow (E) \mid Id$

string = id + (id * id)

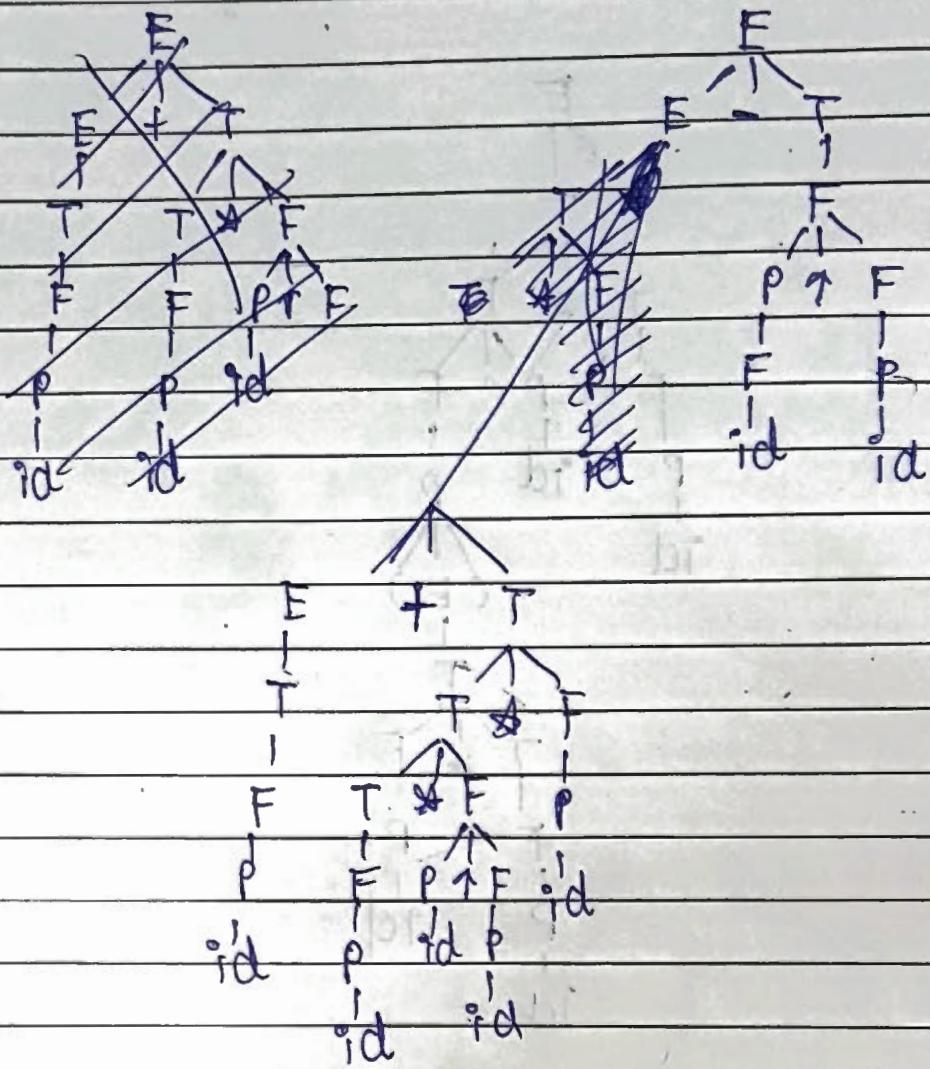


string \Rightarrow id * id^{c(id * id)}



~~Trick~~: Traverse till right, then
 ↓ while returning, split based on lowest precedence operator
(check end on return)

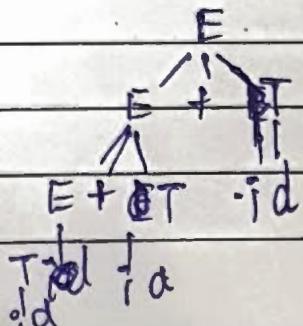
Q- id + id * id / id * id - id + id.



→ self notes

1+2+3 \Rightarrow left associative

so, we grow tree on left side \rightarrow left recursive grammar



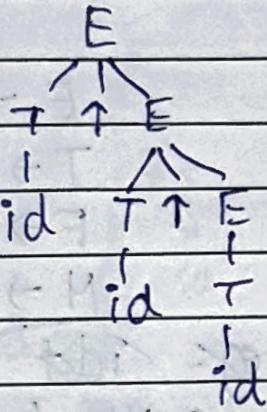
$$\begin{aligned} E &\rightarrow E + id \mid id \quad T \\ T &\rightarrow id \end{aligned}$$

grammar

⇒ $2 \uparrow 3 \uparrow 2$ - power is right associative
 $- 2^{3^2} = 2^9$ ↳ right recursive grammar

$$E \rightarrow T \uparrow E \mid T$$

$$T \rightarrow id$$

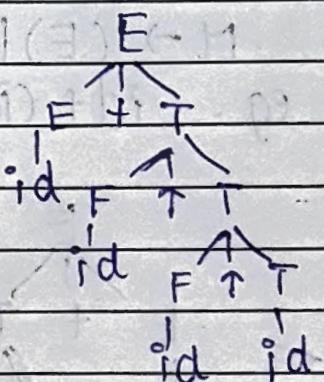


⇒ $2 + 3 \uparrow 2 \uparrow 4$ ~~end~~

$$E \rightarrow E + T \mid T \rightarrow \text{left recurs}$$

$$T \rightarrow F \uparrow T \mid F \rightarrow \text{right recurs}$$

$$F \rightarrow id$$



$$\text{eg} - A \rightarrow A \$ B \mid B$$

$$B \rightarrow C @ B \mid C$$

$$C \rightarrow C \# D \mid E$$

$$E \rightarrow id$$

⇒ here, associativity of
 $A = \text{left}, B = \text{right},$
 $C = \text{def}$

precedence

⇒ $\# > @ > \$$

class eg. resolved

eg. $E \rightarrow E+E \mid E \cdot E \mid E/E \mid E^{\star} E \mid id.$

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

$T \rightarrow T \cdot F \mid T/F \mid F$

$F \rightarrow M \uparrow F \mid M$

$M \rightarrow id$

eg. ~~1st~~

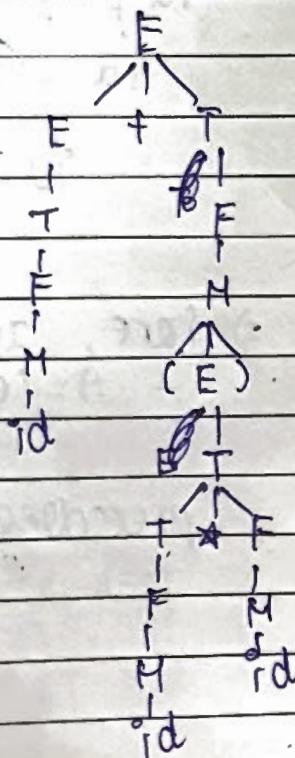
with brackets,

{ same }

$F \rightarrow M \uparrow F \mid M$

$M \rightarrow (E) \mid id$

eg. $id + (id \cdot id)$



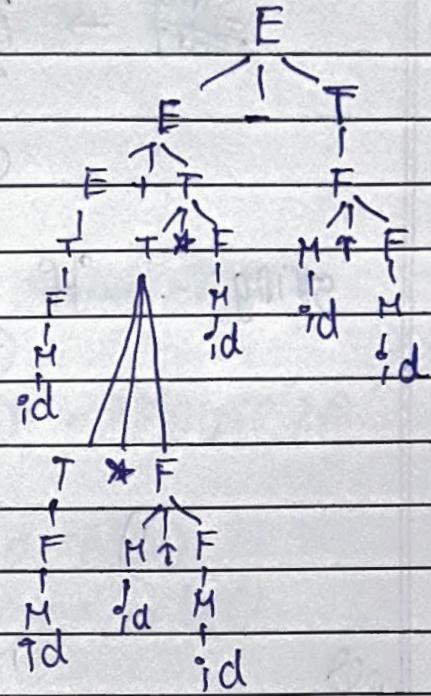
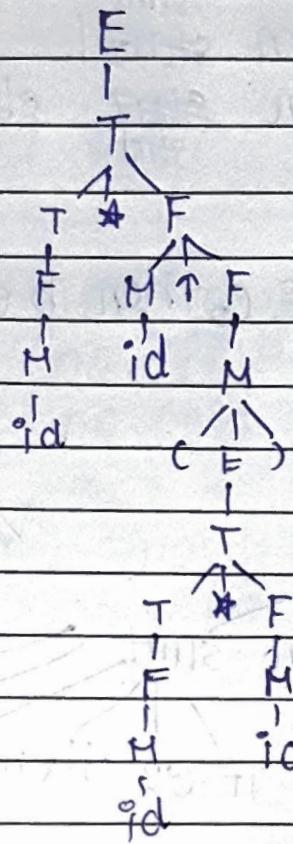
Date : / / 20

Page No.:

split based on layout
precedence from left to right

eg - $\text{cid} * \text{id}^n (\text{cid} * \text{id})^*$

$\text{pd} + \text{id} * \text{id} \text{ T id} * \text{id} - \text{id} \text{ T id}$

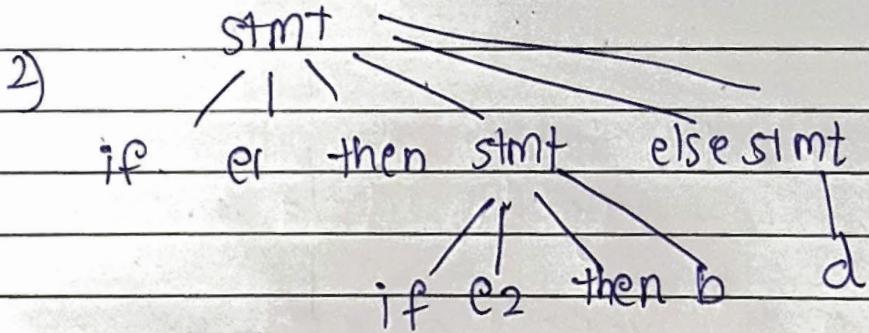
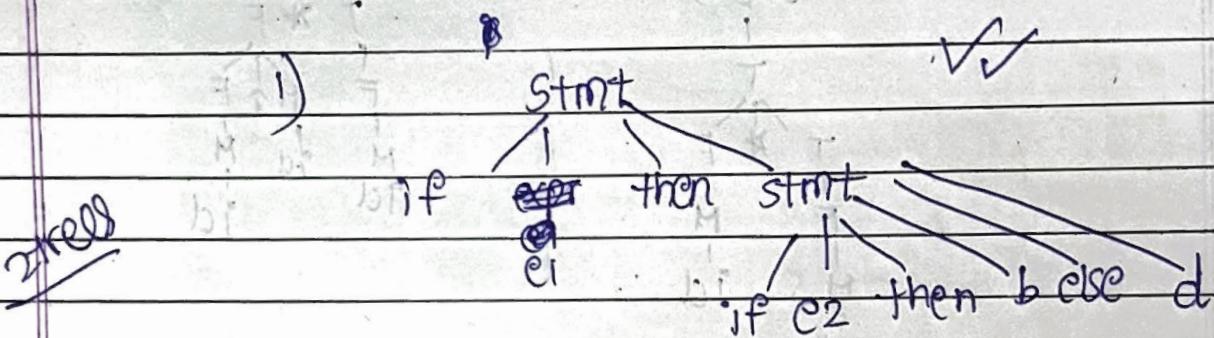


• we also have inherited ambiguous grammar \rightarrow $a^n b^n c^n d^n$
 • ambiguos \rightarrow convers \rightarrow unambiguous.
 \hookrightarrow not removable \rightarrow inherited ambiguous grammar
 Date: / /20
 Page No.:
 amb "c" d"
 inherited ambiguous grammar

Ariya \rightarrow Dangling else problem

~~stmt~~ \rightarrow if expr. then ~~stmt~~
 if expr. then ~~stmt~~ else ~~stmt~~
 other.

string: - if e, then if e, then b else d



\Rightarrow 1st is correct, as if is always attached with nearest if.

- This is dangling else problem, as we don't know which else is attached with which if.
- The grammar is ambiguous.

Solution to the problem

; There are 2 types of if statements

$\Rightarrow \text{no. (if)} = \text{no. (else)}$

$\Rightarrow \text{no. (if)} \neq \text{no. (else)} \rightarrow \text{simple if}$

$\begin{matrix} \\ \hookrightarrow m-s \end{matrix}$

match stmt $\Rightarrow \text{no. (if)} = \text{no. (else)}$

unmatch stmt $\Rightarrow \text{no. (if)} \neq \text{no. (else)}$

$\hookrightarrow \text{vn-s}$

$\Rightarrow \text{stmt} \rightarrow m-s \mid \text{vn-s}$ equal if & else

$m-s \rightarrow \text{if expr then } m-s \text{ else } m-s \text{ (other)}$

$\text{vn-s} \rightarrow \text{if expr then stmt.}$

$\text{if expr then } m-s \text{ else } v-s \text{ (other)}$

e.g. if a then if b then c else d .

stmt

vn-m

if

expr

then

stmt

3. vn-m (2)

ei

m-s

match

if

expr

then

m-s

else

m-s

e'2

op

op

op

b

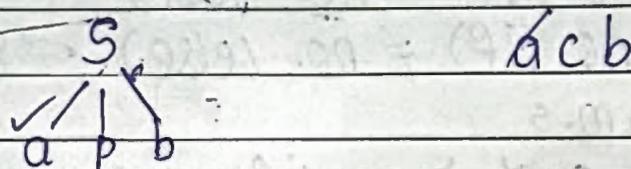
b

b

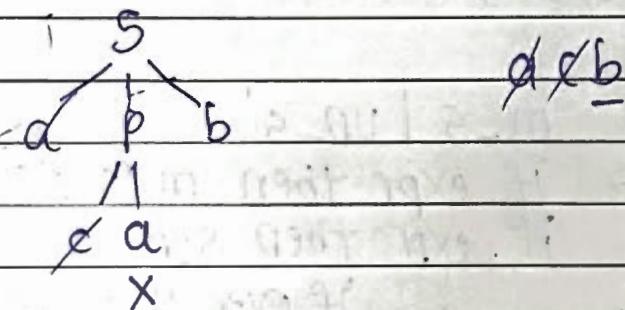
→ Recursive Descent Parsing → simplest form of top down parsing

eg $\Rightarrow S \rightarrow aPb$
 $P \rightarrow c \mid ca$
 $\text{string} \Rightarrow acb$

Start from starting symbol \Rightarrow



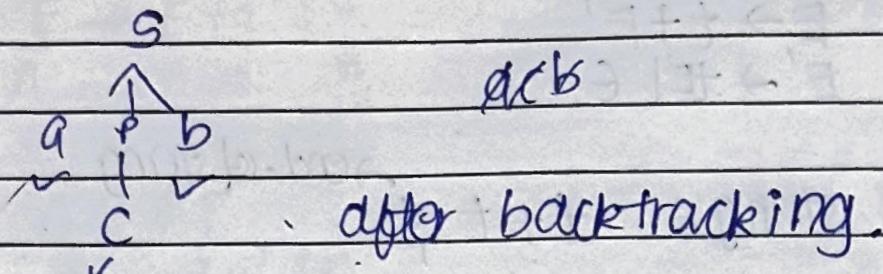
if we choose $\Rightarrow ca$



now ~~the~~ as we do not get a match, we do backtracking.

Problem of backtracking \Rightarrow let $P = calcb \mid cc$, we don't have c.

so, prog parser backtracks & tries all combin', but we don't get answer. $\rightarrow \infty$ loop \rightarrow parser crash.



- There are 2 more problems.

- 1) Left recursion

$$\Rightarrow S \rightarrow \cancel{S} a$$

$$\Rightarrow \infty \text{ loop}$$

- 2) Left factoring

$$P \rightarrow c a | c b | c c$$

↓ ↓

same prefix 'c'.

so which to choose?

Because of problems with ambiguity,
left recursion & left factoring, CLR don't
use recursive descent parsing.

→ To solve left recursion,
problem $\Rightarrow A \rightarrow A\alpha / \beta$,

ans. $\Rightarrow A \rightarrow \beta A'$
 $A' \rightarrow \alpha A' / \epsilon$

$$\text{eg} - E \rightarrow t + E' \\ E' \rightarrow t \mid \epsilon$$

let string \Rightarrow $t + t \text{ } \overset{\text{end-of-string}}{\downarrow}$

\rightarrow Recursions

Recursion

Left

$$E \rightarrow E \text{id} id$$

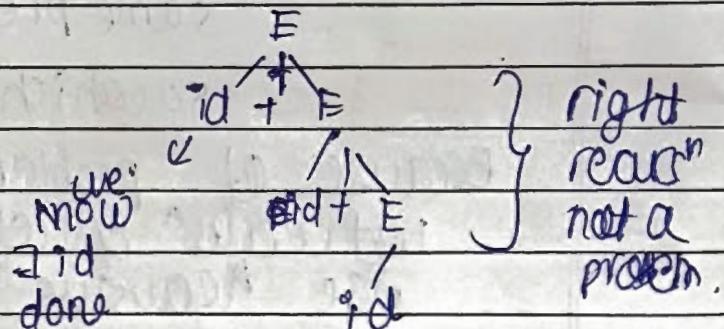
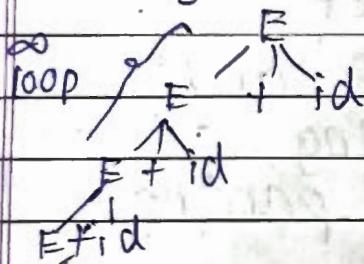
left
terminal
dot side

eg $id \text{id} id \text{id}$

Right

$$F \rightarrow id + E \text{id}$$

eg $id + id + id$



here parser

doesn't
know how
many id's
are generated

$$\text{Q9. } E \rightarrow \frac{E+t}{A} \mid \frac{t}{A \alpha} \mid \frac{t}{B}$$

$$A \rightarrow BA'$$

$$A' \rightarrow dA' \mid E$$

$$\therefore E \rightarrow tE'$$

$$E' \rightarrow t + E' \mid E$$

* Non deterministic grammar

- Here the parser gets confused when looking at the productions.

$$S \rightarrow qAr$$

$$A \rightarrow b \mid bc$$



more than 1 option - non determinism

\hookrightarrow confusion for parser

So, we need to apply left factoring

Left factoring - grammar technique to make production suitable for top down parsing.

Step - If there is a grammar with production
 $A \rightarrow \alpha \beta, \mid \alpha \beta_2$, then we can rewrite as:-

$$A \rightarrow \alpha A'$$

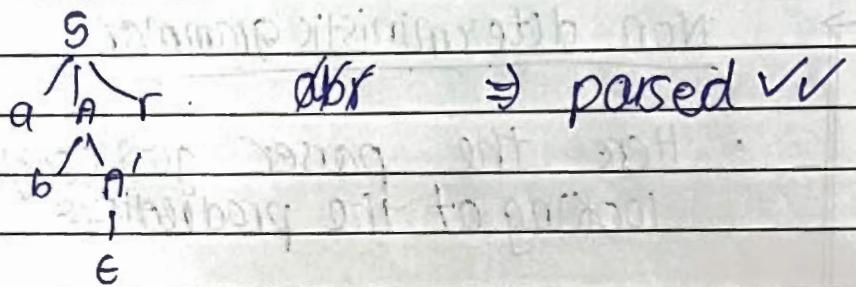
$$A' \rightarrow \beta, \mid \beta_2$$

\rightarrow keep common prefix & introduce new terminal.

$\therefore \text{in } S \rightarrow qAr$
 $A \rightarrow bA'c$

$\Rightarrow S \rightarrow qAr$
 $A \rightarrow bA'$
 $A' \rightarrow cIE$

now to get abr.



For top down parsing, we need 3 condit's -

- 1) grammar - unambiguous
- 2) grammar - no left recursion (no non del)
- 3) grammar - no non determinism
 (left factoring problem)

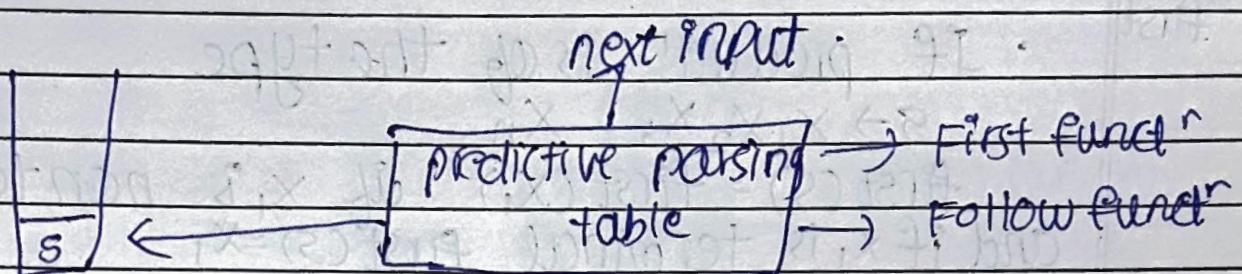
Predictive Parsing

- It is a parsing technique in which the select of the product is based on the input symbol.

$$S \rightarrow a \underbrace{lb}$$

parser predicts which is
suitable/correct
product?

- we have a stack & a predictive parsing table.



- This method uses TOS & the next input.
- This predictive parsing is of 2 types
 - 1) Recursive predictive parsing \rightarrow x not used.
 - 2) Non recursive predictive parsing
 - \hookrightarrow maintaining the stack implicitly.
 - \hookrightarrow we can move down (parsing)
- We use 2 functions to create this table
 - 1) FIRST function
 - 2) FOLLOW function

1) First funct'

- If the product' is of the type

if $S \rightarrow a$ then
 we expand S , we get a first
 and $\text{first}(S) = a$ → nonterm. → terminal

- If product' is of the type

$S \rightarrow X_1 X_2 X_3 \dots X_n$
 $\text{first}(S) = \text{first}(X_1)$ if X_1 is nonterminal.
 and if X_1 is terminal $\text{first}(S) = X_1$,

eg - $S \rightarrow A B C$ $\text{first}(S) = \text{first}(A) = a$

$$A \rightarrow a$$

and $\text{first}(S) = \begin{cases} \text{first}(X_1) \cup \emptyset & \text{if } X_1 \neq \emptyset \\ \text{first}(X_2) \cup \emptyset & \text{if } X_1 = \emptyset \\ (\text{if } X_1 \text{ is nullable}) \end{cases}$

eg - $S \rightarrow A B C$

$$A \rightarrow a \mid \emptyset$$

$$\Rightarrow abc \text{ or } bc \quad \therefore \text{first}(S) = a \text{ or } b$$

here $\text{first}(A) = a \text{ or } \emptyset$

2) Follow funct'

- If product' is of the type $S \rightarrow \alpha$ and S is the starting symbol,

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

- If product is of the type $S \rightarrow \alpha B \beta$

$\text{follow}(B) = \text{first}(\beta)$

we look
on right
side

- If B is nullable, $(B \xrightarrow{*} \epsilon)$

$\text{Follow}(B) = \text{follow}(S) \quad (\epsilon \in \text{Follow}(B))$

\Rightarrow eg - $S \rightarrow \alpha A B | A S B$
 $B \rightarrow b | \epsilon | S$

$\text{follow}(S) = \{\$, b\}$

$\text{follow}(B) = \text{follow}(S) = \{\$, b\}$

$\text{follow}(A) = \{b, \$\}$

- If product of the type -

$S \rightarrow \alpha B$

$\text{follow}(B) = \text{follow}(S)$

$\Rightarrow \text{follow}(A) = \{b, \$\}$



Q- $E \rightarrow ETT \mid T$
 $T \rightarrow T \star F \mid F$
 $F \rightarrow (E) \mid id$

$E \rightarrow TE'$
 $E' \rightarrow +TE' \mid \epsilon$
 $T \rightarrow FT'$
 $T' \rightarrow \star FT' \mid \epsilon$
 $F \rightarrow (E) \mid id$

\rightarrow	<u>First</u>	<u>Follow</u>
E	$\{ (, id \} \}$	$\{ \cdot \}, \$ \}$ $= \text{follow}(E) = \{ \epsilon,) \} \}$
E'	$\{ +, \epsilon \} \}$	$\Rightarrow \text{first}(E') = \{ (, id \}, \{ +, \epsilon \},) \} \}$
T	$\{ (, id \} \}$	$= \text{follow}(T) = \{ +, \epsilon,) \} \}$
T'	$\{ \star, \epsilon \} \}$	$= \{ \star, \epsilon, +,) \} \}$
F	$\{ (, id \} \}$	

eg-8 - $S \rightarrow aSA \mid SB$

$A \rightarrow \del{a} \mid H \mid E$

$B \rightarrow bcd \mid E$

$H \rightarrow h \mid E \mid b$

$S \rightarrow aSA \mid qSAS'$

$S' \rightarrow BS' \mid E$

$A \rightarrow aHIE \mid$

$B \rightarrow bcd \mid E$

$H \rightarrow h \mid E \mid b$

first

follow

$S \quad \{a\}$

~~{}\$~~ $\{s, q, b\}$

$S' \quad \{\epsilon, b\}$

$\{\$, q, b\}$

$A \quad \{a, \epsilon\}$

$\{\$, q, b\}$

$B \quad \{b, \epsilon\}$

$\{\$, q, b\}$

$H \quad \{h, b, \epsilon\}$

$\{\$, q, b\}$

→

class Q
negative

$$\begin{aligned} E &\rightarrow EFTT \\ T &\rightarrow \star T \rightarrow T \star F | F \\ F &\rightarrow (E) | id \end{aligned}$$

3 "left recurs"

$$\begin{aligned} \Rightarrow E &\rightarrow TE' \\ E' &\rightarrow +TE' | \epsilon \\ T &\rightarrow FT' \\ T' &\rightarrow \star FT' | \epsilon \\ F &\rightarrow (E) | id \end{aligned}$$

	<u>First</u>	<u>Follow</u>
E	{C, id}	{+, -, *, /, (,), \$}
E'	{+, \epsilon}	{+, -, *, /, (,), \$}
T	{C, id}	{+, -, *, /, (,), \$}
T'	{\star, \epsilon}	{+, -, *, /, (,), \$}
F	{C, id}	{\star, +, -, *, /, (,), \$}

Q- $S \rightarrow ASA | SB$
 $A \rightarrow AHIE$
 $B \rightarrow bcd | \epsilon$
 $H \rightarrow h | Elb$

3 "left recurs".

$$\begin{aligned} \Rightarrow S &\rightarrow ASAS' \\ S' &\rightarrow BS' \\ A &\rightarrow AHIE \\ B &\rightarrow bcd | \epsilon \\ H &\rightarrow h | Elb \end{aligned}$$

	: First	Follow
S	{a}	{a, b, \$}
S'	{b, ε}	{a, b, \$}
A	{a, ε}	{b, a, \$}
B	{b, ε}	{b, a, \$}
H	{h, ε, b}	{b, a, \$}

- class →
- Now to use predictive Parsing after calc. first & follow, we list down nonterm. row wise & terminals usually in order of increasing order of precedence) in rows :

	+	*	()	id	\$	→ predictive parsing table
E			E → TE'		E → TE		
E'	E' → +TE'			E' → ε		E' → E	
T			T → FT'		T → FT'		
T'	T' → ε	T → AFT'		T' → ε		T' → ε	
F			F → (E)		F → id		

in a producer

- 2). If we get ϵ , we put it in follow of that ~~use~~ term \rightarrow nonterm. $\rightarrow \epsilon$
- 3). we put the non-term.'s product in its follow, first

↳ eg - if we want to write for S, we check

First(S)

↳ then,

we write

product's.

eg - $E' \rightarrow +TE' | E$

↳ put $E \rightarrow E'$

in follow symbols of E'

- Now we can parse a string using the table.
- (as on right side)

stack	input	action
\$ E	id + id * id \$	
\$ E' T	match now see entry of T & id \rightarrow id + id * id \$	E \rightarrow TE'
\$ E' T' F	id + id * id \$	T \rightarrow FT'
now T' & F	\$ E' T' id \rightarrow pop	F \rightarrow id
\$ E'	+ id * id \$	T' \rightarrow E
\$ E' T' F	* id * id \$	E' \rightarrow + TE'
\$ E' T' F	id * id \$	T \rightarrow FT'
\$ E' T' id	id * id \$	F \rightarrow id
\$ E' T' F*	* id \$	T' \rightarrow FT'
\$ E' T' id	id \$	F \rightarrow id
\$ E'	\$	T' \rightarrow E
\$	\$	E' \rightarrow E

\Rightarrow [parsed.] \nwarrow accepted

- If we get do not get a prod - prod pair, then it is a error.
- Then error handler will recover.
 - ↳ panic mode error recovery.

\rightarrow Here, by looking at stack & i/p symbol, parser can predict what could be the next prod.

Q- $S \rightarrow A$

$A \rightarrow aB \mid Ad$ "def/ recurs"

$B \rightarrow bBC \mid f$

$C \rightarrow g$

string $\Rightarrow a b b f g g d$.

$\Rightarrow S \rightarrow A$

$A \rightarrow aBA'$

$A' \rightarrow dA' \mid \epsilon$

$B \rightarrow bBC \mid f$

$C \rightarrow g$

first

follow

S

$\{a\}$

$\{\$ \}$

A

$\{a\}$

$\{\$ \}$

A'

$\{d, \epsilon\}$

$\{\$ \}$

B

$\{b, f\}$

$\{g, d, \$\}$

C

$\{g\}$

$\{g, d, \$\}$

	a	b	d	f	g	\$
S	$S \rightarrow A$					SDA
A	$A \rightarrow aBA'$					ABAD
A'			$A' \rightarrow dA'$			$A' \rightarrow \epsilon$
B		$B \rightarrow bBC$		$B \rightarrow f$		
C					$\hookrightarrow g$	

stack

input

action

\$ S

abbfggd \$

-

S A

abbfggd \$

S → A

S A' B A'

abbfggd \$

A → A B A'

S A' C B C'

bbfggd \$

B → b B C

S A' C C B C'

bfggd \$

B → b B C

S A' C G F

f ggd \$

B → f

S A' C g

g d \$

C → g

S A' g

g d \$

C → g

S A' d

d \$

A' → d A'

\$

\$

A' e

string accepted