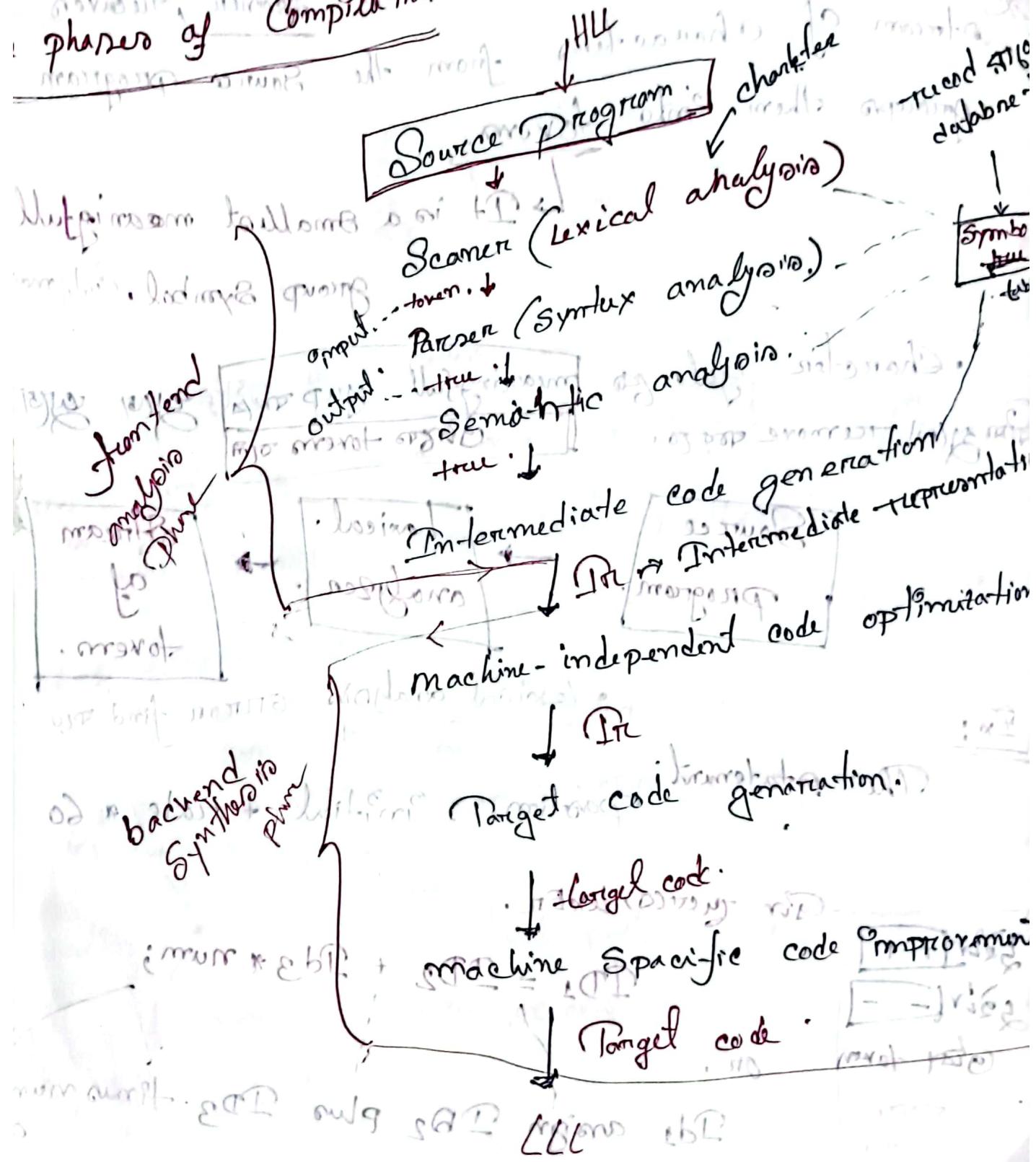


Bismillah hir Rahimahin ir Rahim.

implies \Rightarrow Source program \rightarrow Compiler

\rightarrow Target program ($HLL \rightarrow LLI$)

phases of compilation:



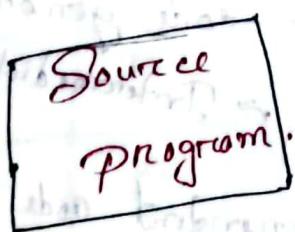
Scanner (Lexical analysis)

" "
 { } ; ,
 / \ , comment, whitespace
 for all token

- * It also called a Lexan or a Scanner, receives a stream of characters from the source program and groups them onto tokens.
- * It is a smallest meaningful group symbol. (int, main.)

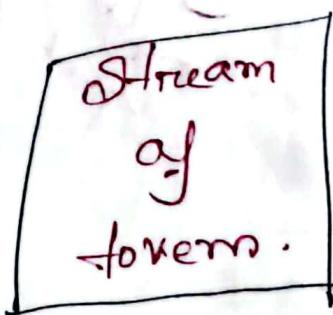
Character rule

Stream of tokens



meaningful group tokens

group tokens

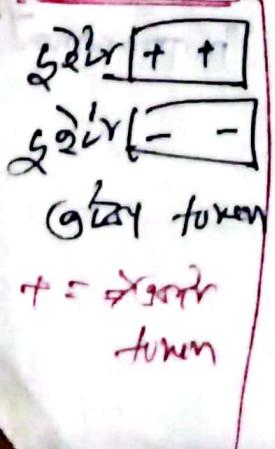


Lexical analysis error finding

Ex:

The statement position = initial + late, a 60

- Get generic lexer.



$$ID_1 = ID_2 + ID_3 * num;$$

or,

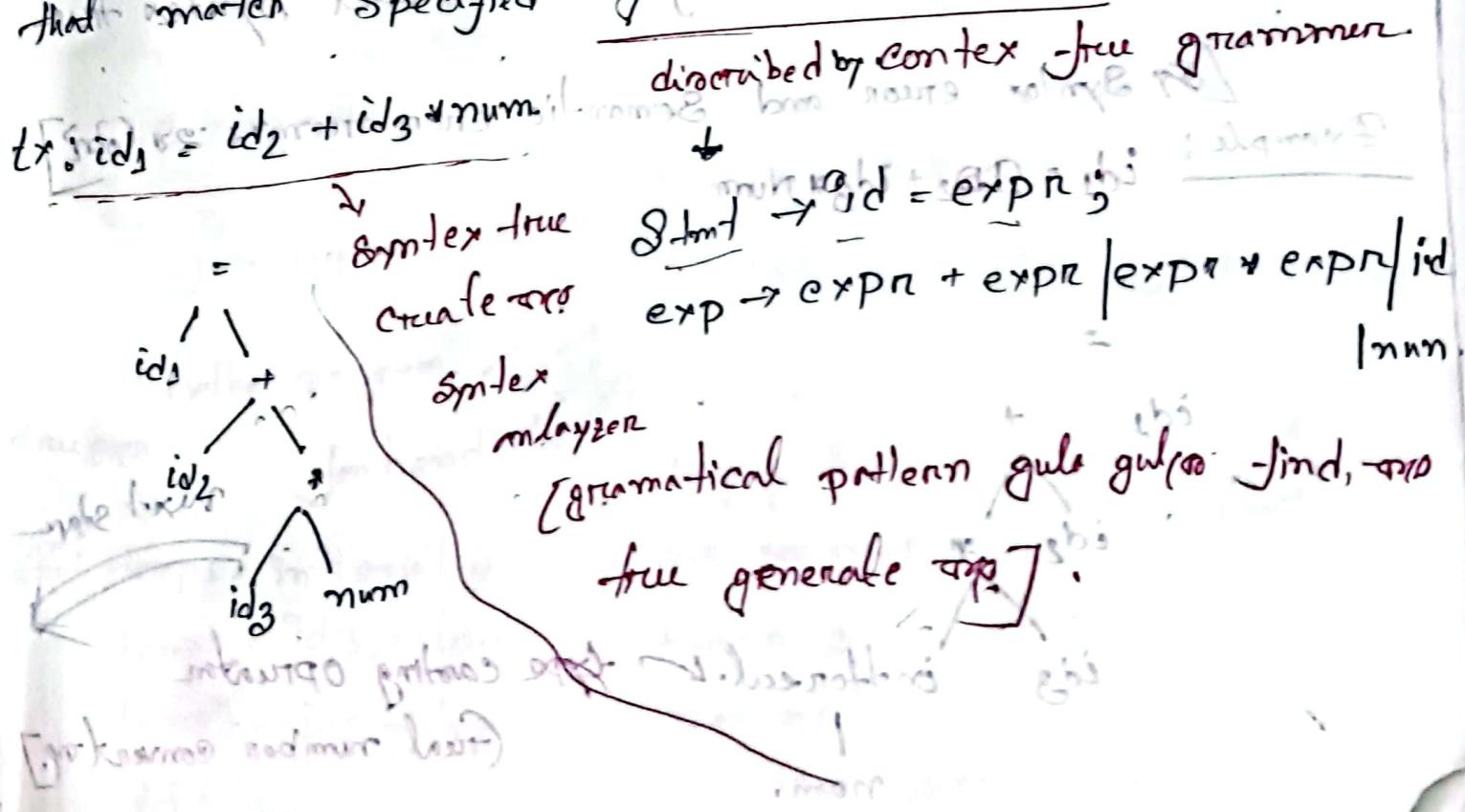
ID_1 assign ID_2 plus ID_3 times num

Lexical analysis Tools

- lex → C source code. (UNIX)
- flex → " "
- Jflex → Produces Java code (Java)
- Jflex → " " " " . (Java source code)

Syntax analyzer (parser)

A syntax analyzer also called parser, receives a stream of tokens from the lexer and groups them into phrases that match the specified grammatical pattern.



Syntax analyzer Tools.

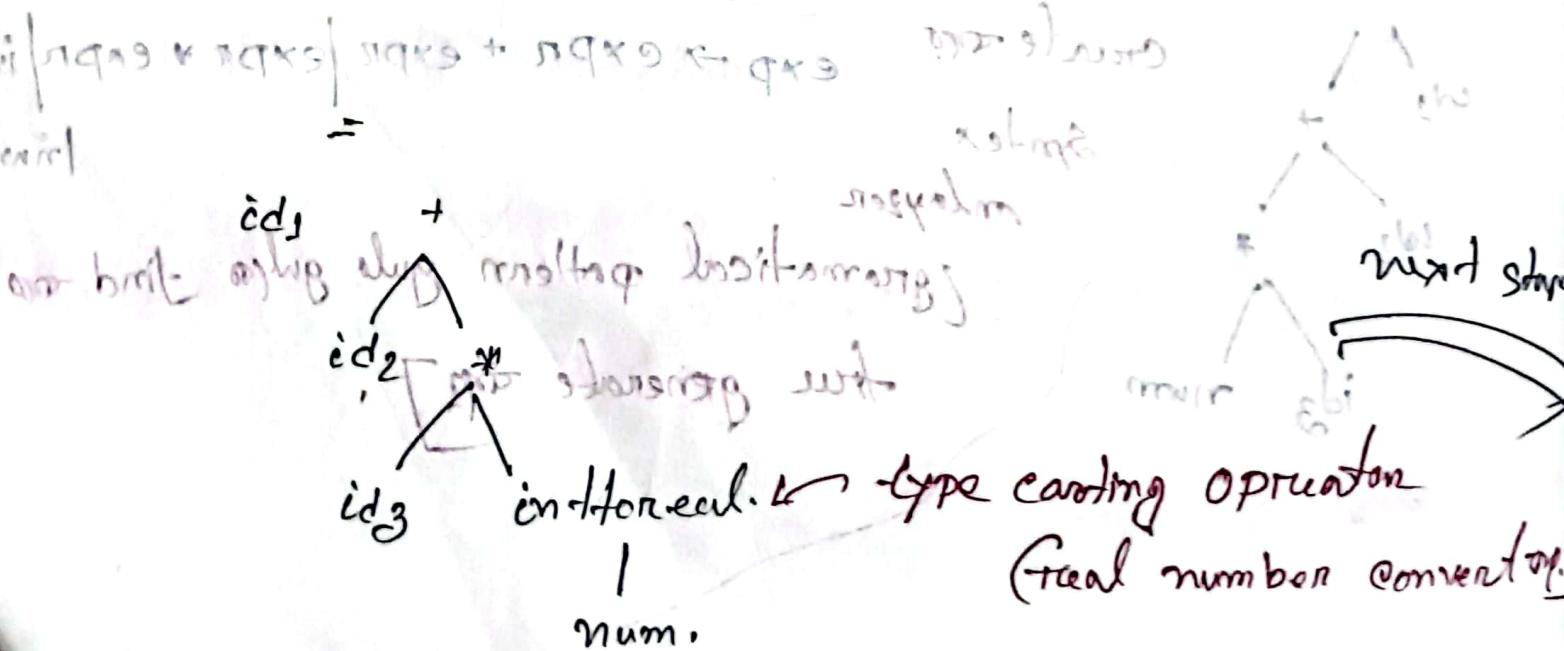
- Yacc \rightarrow C Source code (UNIX)
- Bison \rightarrow C " " (gnu)
- CUP \rightarrow Java Source code

(non-deterministic) \rightarrow Semantic analysis
Semantic Analysis.

Semantic analyzer traverses the abstract tree, Checking that each node is appropriate for context. It checks for semantic errors. It outputs refined abstract syntax tree.

[By Syntax error and Semantic error handling 20 Gm]

Example: $id_1 = id_2 + id_3 \# Num$



Intermediate code generation.

Intermediate code

Intermediate code is code that represents the semantics of a program but is machine-independent.

(কোড ব্যবহার করে প্রোগ্ৰাম সমূহের অন্তর্ভুক্ত কোড হ'ল ইন্টেরমিডিএট কোড)।

Discusses

Abstract syntax \rightarrow Intermediate code \rightarrow Intermediate code

Program generator \rightarrow machine code

স্থির ও বেশি স্থির বেশি

স্থির ও বেশি স্থির বেশি

স্থির ও বেশি স্থির বেশি

⇒ x86
code.

C program

Intermediate code

Java program.

Code + SB

Python program.

MIPS86

code.

Front end

Back end

(machine)

language specific

language

target machine

common source

target code

Gen. generate code.

Ex: temp1 = 9 * 10 + 1 (6)

temp2 = 9 * 3 + temp1.

temp3 = id2 + temp2

id1 = temp3

• multi-dimensional array has intermediate code

Code optimizer

• optimize by reducing complexity of code, memory requirement

→ reduce memory requirement by using 2-line code
[2-line code is good]

- An optimizer reviews the code, looking for ways to reduce the number of operations and memory requirement
- optimized for speed or size.

Example: the intermediate code in may be optimized as

$$\text{temp1} = \underline{\text{id3}} * \underline{60.0}$$

$$3889\text{id1} = \underline{\text{id2}} + \underline{\text{temp1}}$$

• ibus

(bus word)
(bus width)

read rd

int regd

float regf

word mem

float mem

• inter storage rd

Machine Code generation.

- The code generator receives the (optimized) intermediate code.
- It produces:
 - Machine code for specific machine
 - Assembly code for a specific assembler

If it produces assembly code, then assemble to machine code.

Ex: the inter

$$\begin{aligned} G &= ID_3 * 60.0 \\ id_1 &= cd_2 + t_1 \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} DR$$

\downarrow
assembly code:
 $\begin{aligned} &\text{movf } id_3, R_2 \\ &\text{multf } \#60.0, R_2 \\ &\text{movf } id_2, R_1 \\ &\text{addf } R_2, R_1 \end{aligned}$

Intermediate code
→ Code generation

Assembly → Assembly
prioritizing transitions
Machine code

Machine code < assembly code

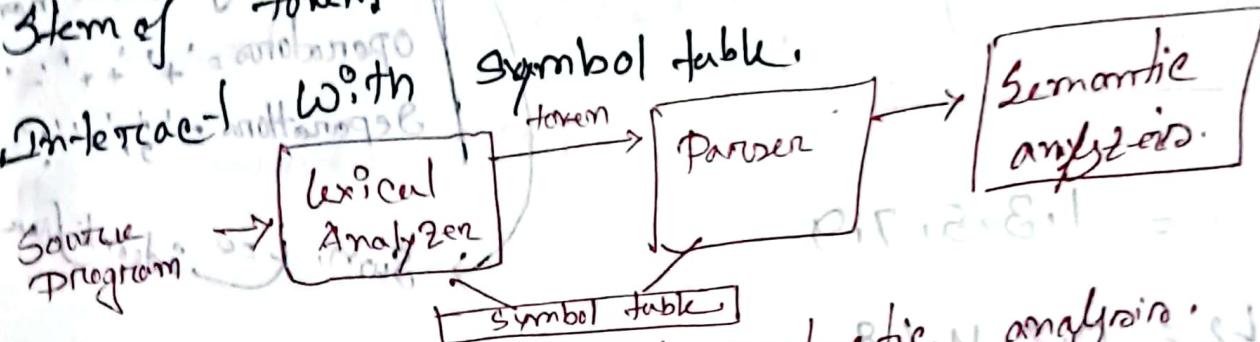
Symbol table

- Records <variable name, attributes>
- Lexical, syntactic elements
- store
- base of address
- functions
- variables
- values
- arguments
- parameters
- return type
- procedure names
- scope
- type
- storage information
- Attributes:

Design principle: find, store, delete, search quickly.

Tasks of Lexical Scanner / Tokenizer (TTR)

- Read the Input Char of the Source Programs.
- group them into tokens
- produce as output a sequence of tokens for each lexeme in the Source Program.
- Some of tokens sent to the parser.



why Separating lexical and syntactic analysis is important consideration.

- Simplicity of design is the most important consideration.
- efficiency is improved.
- portability is enhanced.

Lexeme : Lexeme to pattern Match token.

Scanning : comment white space deleted (UN).

Lexical analysis : It produces tokens (tokens from scanner)



Position = $0 \times 60 + id_1 * 60$ // Lexeme

id_1

$$= id_2 + id_3 * 60 // Token$$

variables to compare o buffering
memory error art. no. of times
new up art of time

Operations of language: Identifiers

$$L_1 = 1, 3, 5, 7, 9$$

start Token

$$L_2 = 0, 2, 4, 6, 8$$

$$L_1 \cup L_2 = 0, 1, 2, 3, 4, 5, 6, 7, 8, 9$$

$$L_1 \cap L_2 = 10, 12, 14, 16, 18, 30, 32, 34, 36, 38$$

concatenation L_1 followed by L_2 (carat)

$$(L_1 \cup L_2)^3 = \{ (00), (012), (034), (123) \}^3 \text{ length 3.}$$

$$L_1^3 = 111, 333, 555, 133, \{ \text{length 3 at } n \}$$

$$(L_1 \cup L_2)^* = L_2^* = \{ 0, 1, 2, 4, 6, 8, 00, 02, 04, 06, 08, \dots \}$$

$$n = (0, 1, 2, 3, 4, 5, 6, 7, 8, 9) \cdot (0, 1, 2, 4, 6, 8) \cdot L_2^* \text{ count } 0, 1, 4, 8, 6 \text{ etc.}$$

Symbol Table Store 24 Oct 2017.

key words = for, while, if

Identifier = Variable name

function name

Operators = '+', '++', '-'

Separators = ',', '.', '{', '}', '(', ')'

Sum, i = 0, ; identifier

Regular Expression:

• π represents $L(\pi)$, $L(s)$ ^{Geno.}, $\pi \mid s$ ^{represents union.}

• π^* represents concatenation. $(\pi) \cdot L(s)$

• π^* represents more replaced in $L(\pi)^*$. $[e-\emptyset] \leftarrow \text{digib}$
 $+ \text{digib} \leftarrow \text{digib}$

• π represents $L(\pi)$

$S(\text{atigob } e[-,+]) \cup S(\text{atigob } .)$ $\text{digib} \leftarrow \text{random}$

Example:

$\pi = A \mid B \cdots 1z \text{ a1b} \cdots 1z$ ^{/letter}

$\alpha = 0 \mid 1 \cdots 19$ ^{/digit}

$f = \pi (\pi \mid s)^*$ ^{replaced concatenation}
 π concatenation (π union s)

id \rightarrow letter (letter (digit))^*

Cf + identification follows them.

digit^+ $\cap_{\text{un}} 0 \in \Sigma^*$ $\cap_{\text{dig}} \text{dig}^+$

- $\pi^+ = \pi \cdot \pi^*$ ↗ zero on one machine (greater than 10^{23})
- $\pi? = \pi/E$ ↗ zero on one machine (greater than 10^{23})
- # floating Point number

digit $\rightarrow [0-9]$

digit $\rightarrow \text{digit} +$

Number \rightarrow digit (.⁰ digits)? $E^{[+,-]} 9$ digits)?

$$(0.123456789) \cdot 10^{-5} = 1.23456789 \times 10^{-5}$$

$$(21\pi) \approx 65.4$$

With overflow and underflow +1)

trap point are very close to min. trap

4.1) class tokens

digit - 0-9

digit₀ → digit⁺

number → digits (· digit₀) ? (E [+ -] ? digit₀) ?

letter → [A-zA-Z].

id → letter (letter/digit)^{*}.

if → if.

then → then

else → else

tuple → (|), <= | > | { } | []

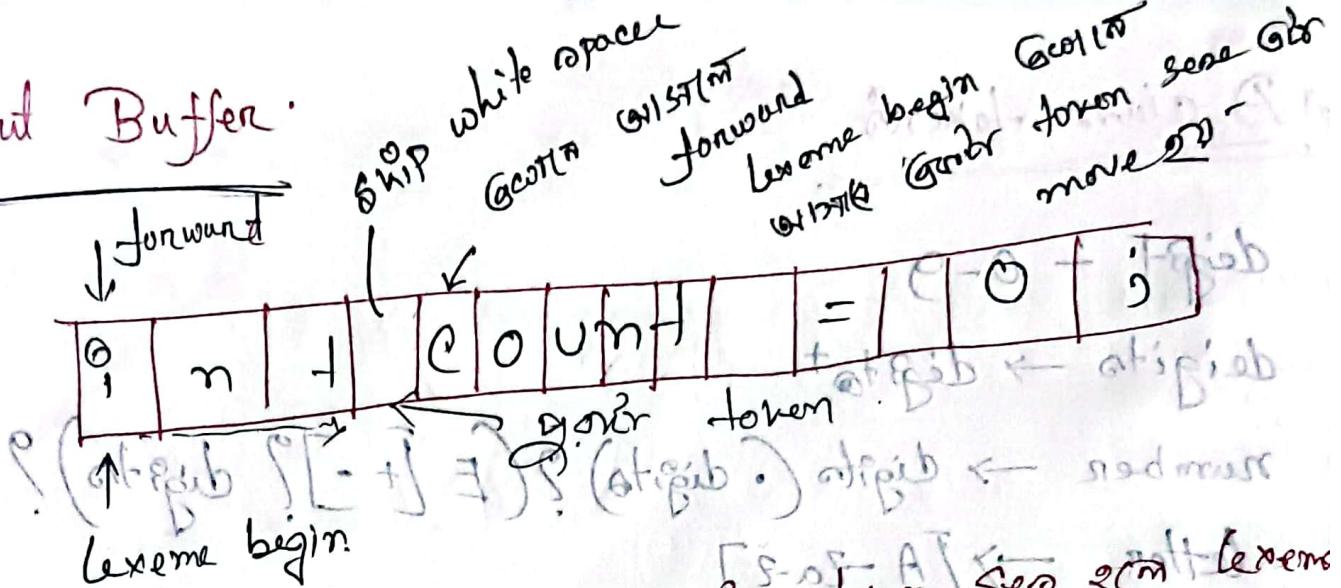
white space : WS → (blank / tab / newline)⁺.

• whitespace is not a token, but it must be recognized by the lexer to add it to the output

• lexer remove whitespace from the input

• lexer knows what to skip

#Input Buffer



Buffer pointer: a two buffer scheme that handles ends
 Large aheads safely $\xrightarrow{\text{checked}}$ 2 buf Buffer

Scan time: save time checking for the ends

+ of Buffers. (Scan from left off the right end of the current buffer).

~~Plane two position of Buffer~~

- Lexeme begin : holds the starting position of the current token.

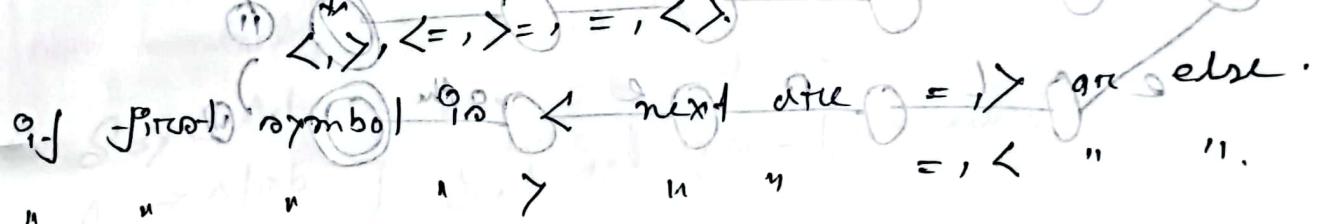
forward pointer : points to the current symbol.

backward pointer : points to the previous symbol.

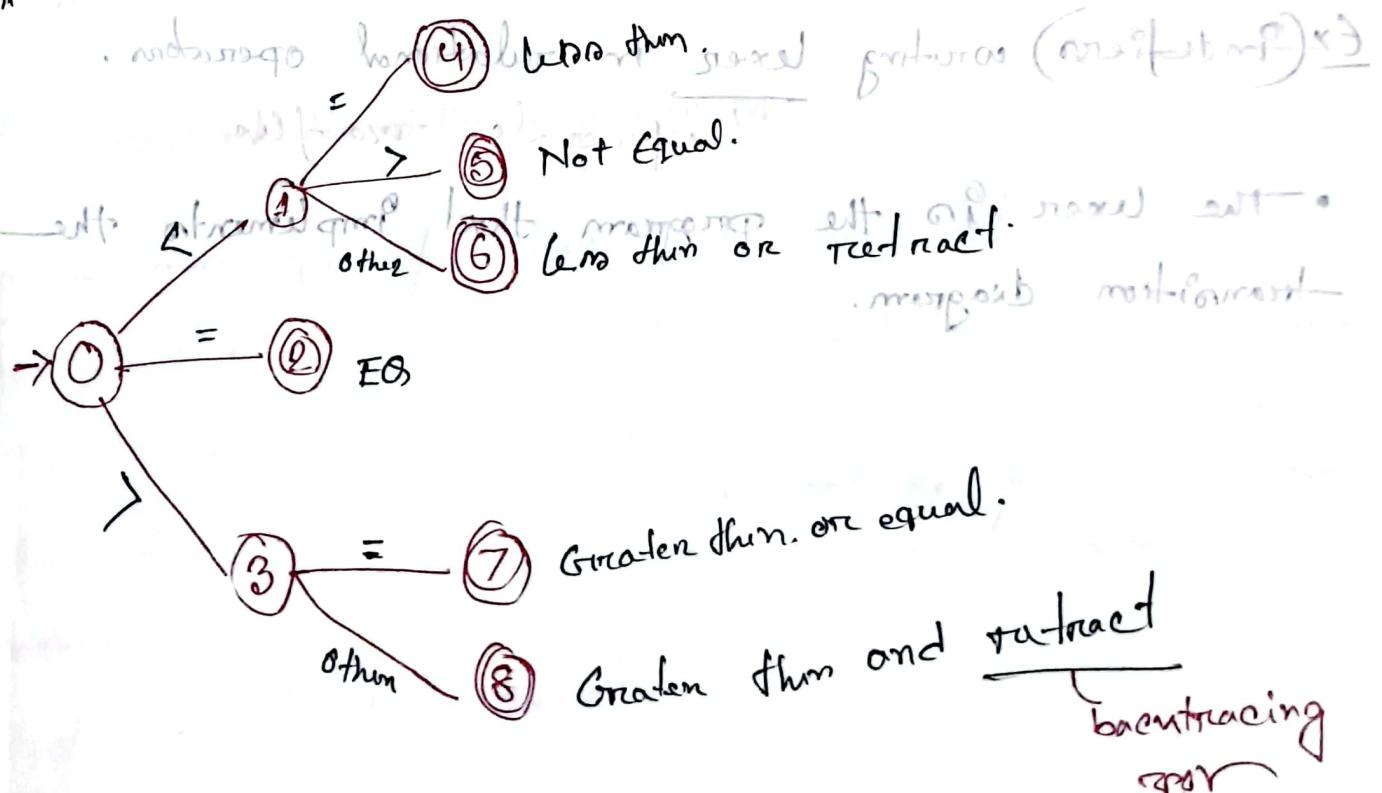
Transition Diagrams

- Transition Diagram is a directed graph.
- It consists of a finite set of nodes called states. One of them start state.
- Set of the states are called accepting states. The remaining states are rejecting states.
- The directed edges between states are transitions.

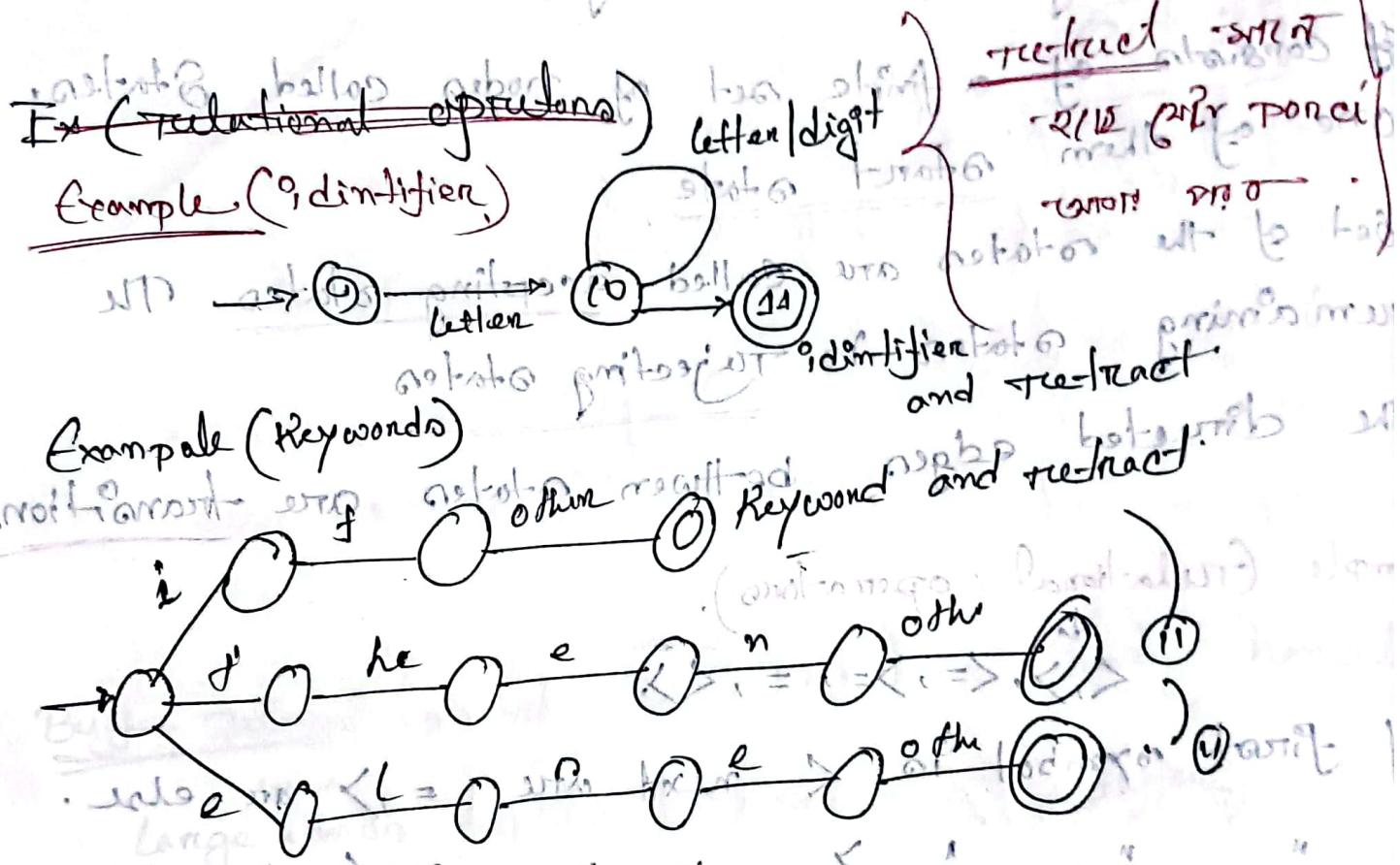
Example (Relational operators),



if first symbol is next after =, > or < else =, <, >, " ", "



Transition Diagram



Ex (Identifier) writing lexer ~~using relational operators.~~
by generate ~~no~~ flex

- The lexer is the program that implements the transition diagram.

Grammar & Expression (N. Datta Sir)

✓

Expr \Rightarrow Expression

Smt \Rightarrow Statement

Such rules are called Production.

- In a production terminal elements like the keyword if or else in parentheses are called terminals.
- variables like expr and string sequences of terminals called non-terminals.

$S \Rightarrow C A d .$ } here a,b,c,d are terminals
 $A \rightarrow \lambda | ab$ } λ, d are (non-terminals)

Context Sensitive grammars (CSG.):

CSG is more powerful than CFG, and it has 4 tuples

$G = \{N, \Sigma, P, S\}$

$N =$ Set of non-terminals

$\Sigma =$ Terminal

$P =$ production

$S.$ = Start symbol of the production.

CSG is left sides and R.H.S hand of any production rule may be surrounded by context of terminal and non-terminal symbols.

d | S | A | S | A

CFG

4 components.

- c. A set of terminals symbol referred are tokens.
2. A set of Non-terminals "called" Syntactic variable
3. A set of Productions (P/R).
4. Start Symbol (S)

digit Sappared by (+, -)

digit (abbreviated) are b, g, d, s and
of Expressions.

• 9-5+2 into syntax of b, g, d, s

Word \rightarrow Word + digit

Word \rightarrow " - " non terminal

digit \rightarrow digit

digit \rightarrow 0/1/2/3/4/5/6/7/8/9

done 89 bro able best of 3020

all the production $A \rightarrow B\alpha$

$A \rightarrow A\alpha / B\alpha / b$

• doing terminal were bim

all the production $A \rightarrow \alpha B$ or $A \rightarrow \alpha$ then Right

$A \Rightarrow \alpha A / \alpha B / b$

Terminal \rightarrow Right
Left

Derivation Tree / Parse tree

graphically represents

the demanded information.

Properties
Parse tree

Root - start symbol.

Interior nodes - non terminal symbol.

leaves - terminal

Two approach to draw derivation tree.

Top down Approach.

starts with starting Symbol S.

Goes down to true leaves using Production.

Bottom up:

Starts with leaves.

upward to the root which is the Starting symbol S

Derivation of $9 - 5 + 2$

Unit → Unit + digit / Unit - digit / digit.

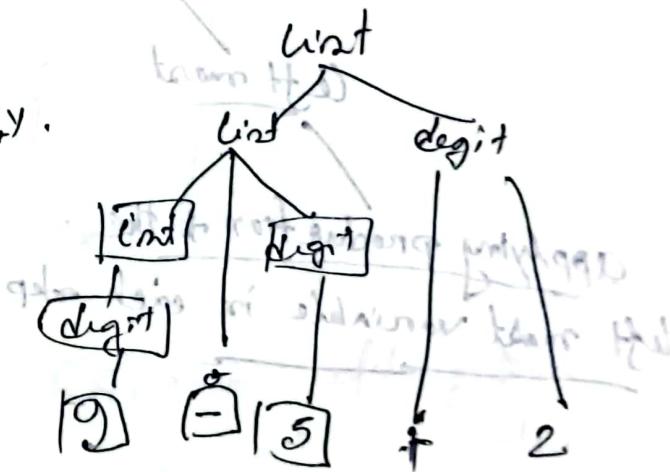
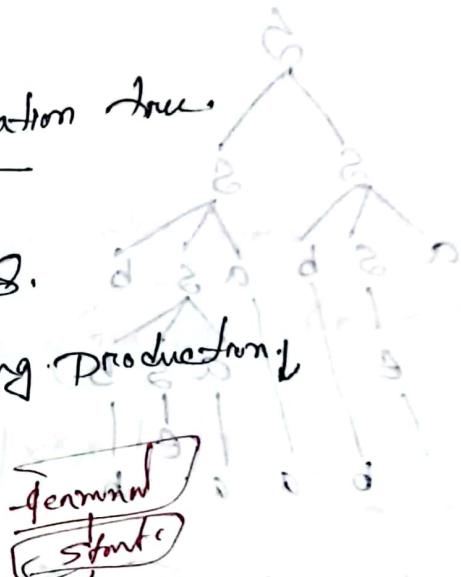
Unit → Unit + digit.

→ Unit + digit + 2

→ Unit - 5 + 2

→ digit + 5 + 2

→ 9 + 5 + 2.



Example

CFG (N, T, P, S)

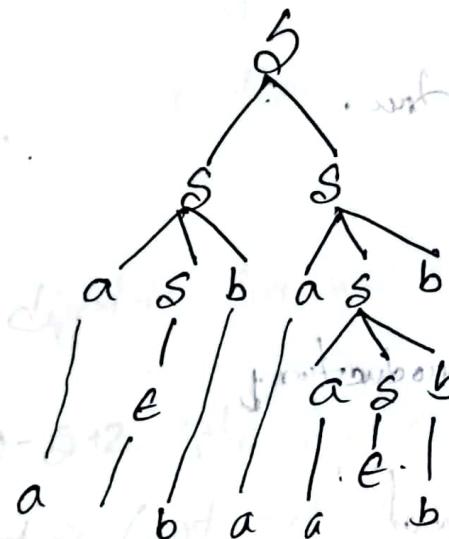
$N = \{S\}$, $T = \{a, b\}$

Non-terminal
Terminal

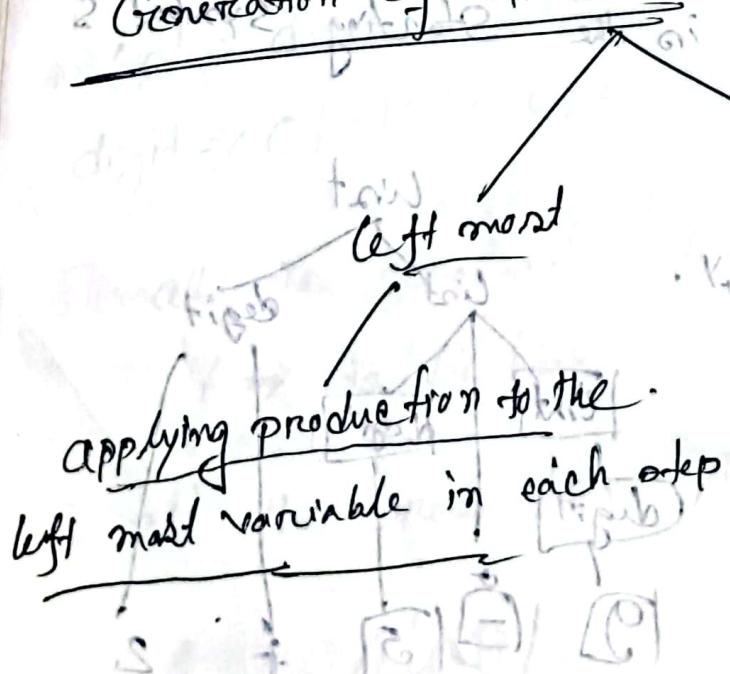
$P = S \rightarrow SS \mid asb \mid \epsilon$

$S \Rightarrow S$

abaabb



2 Generation of tree



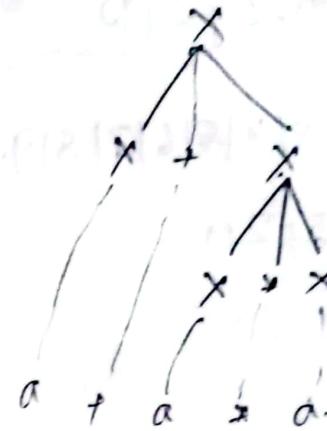
Right most

by applying production to the right most variable in each step.

$$\# \quad X \rightarrow X + X \mid X^* X \mid x/a.$$

left most (a+a*)

$$\begin{aligned} X &\rightarrow X + X \\ &\rightarrow a + X \\ &\rightarrow a + X^* X \\ &\rightarrow a + a^* a \end{aligned}$$

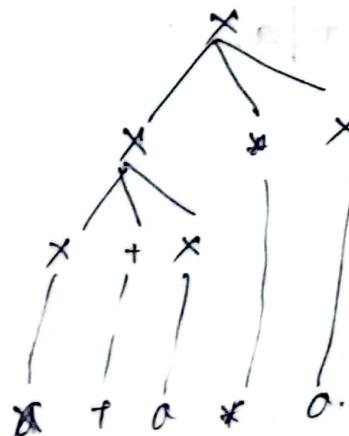


Ambiguity.

If CFG has more than one derivation tree for some string called Ambiguous grammar. (eg. left derivation etc.)

$$X \rightarrow X + X \mid X^* X \mid x/a$$

(a+a*)



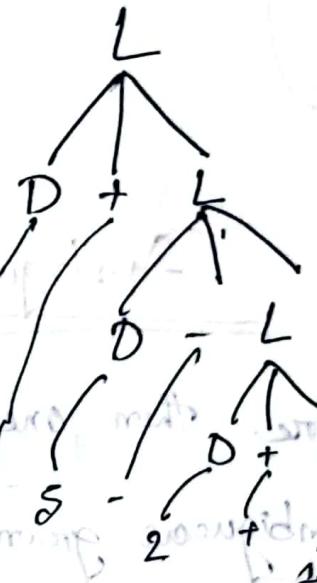
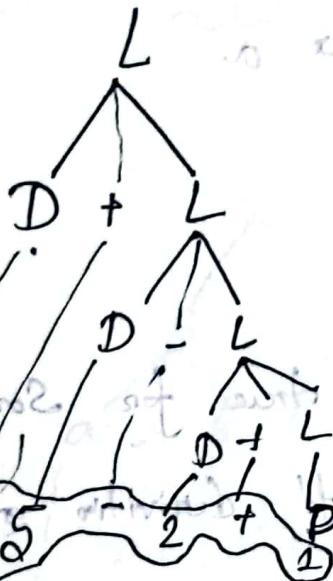
So, to define R_D, and left pattern tree
is ambiguous.

$L \rightarrow D+L | D-L | L|D$

$D \rightarrow 0|1|2|8|4|8|6|7|8|9.$

String $\rightarrow 9+5-2+1$

$0|x|x*x|x+x*$
 $(D+D) \text{ term}$
 $9+5-2+1$
 $x+D+$



a
 (a^*a+a)

$\rightarrow a|a-a|a$

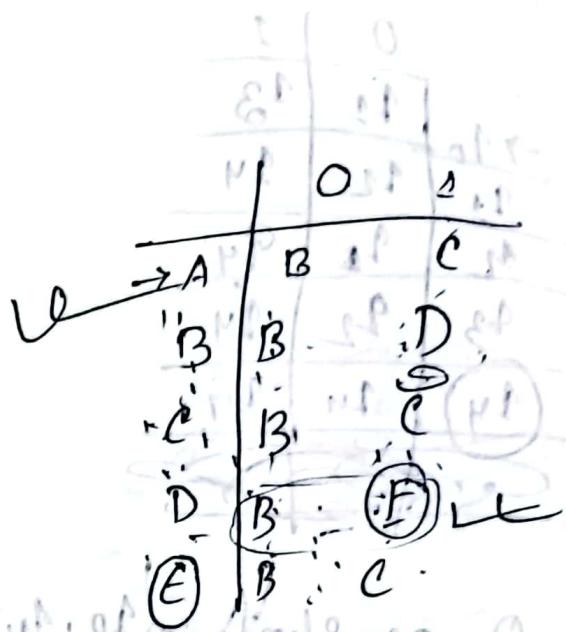
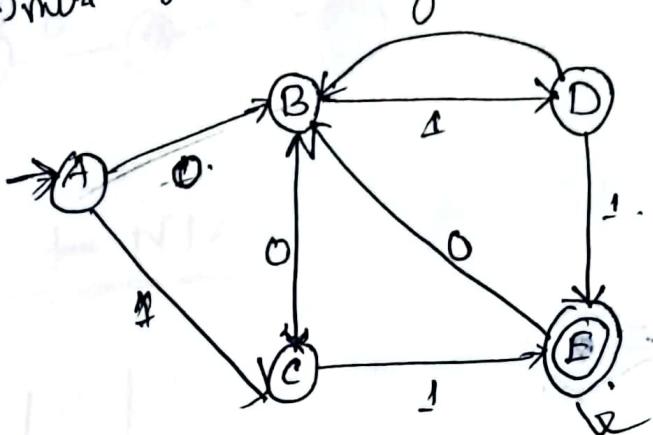


using HLD bmo.89 with of? 02
with

Seg-3

Minimization of DFA

Algorithm follows a step



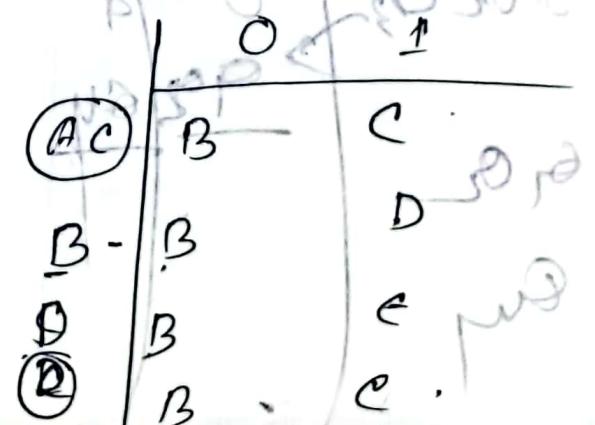
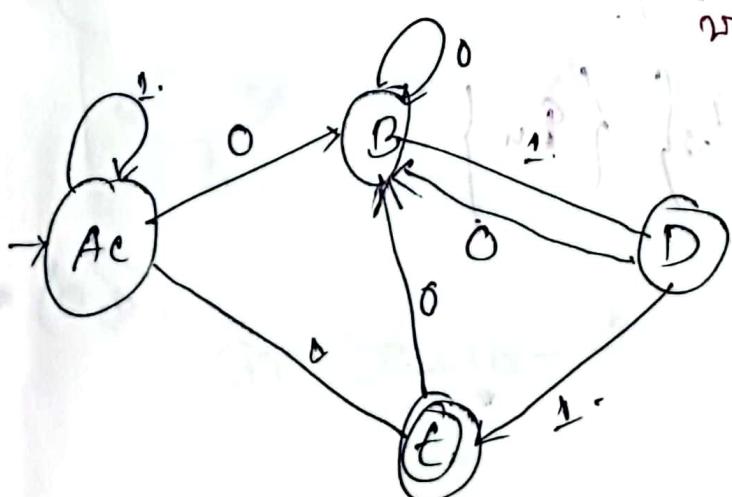
0 equivalent : $\{\underline{A, B, C, D}\}, \{E\}$ ← final stat CONTRARY to initial

1 equivalent : $\{\underline{A, B, C}\}, \{D\}, \{E\}$ ← A, B, C & D are B, B (pair)

2 equivalent : $\{\underline{A, C}\}, \{B\}, \{D\}, \{E\}$ Some contr. to 0 & C, D pair

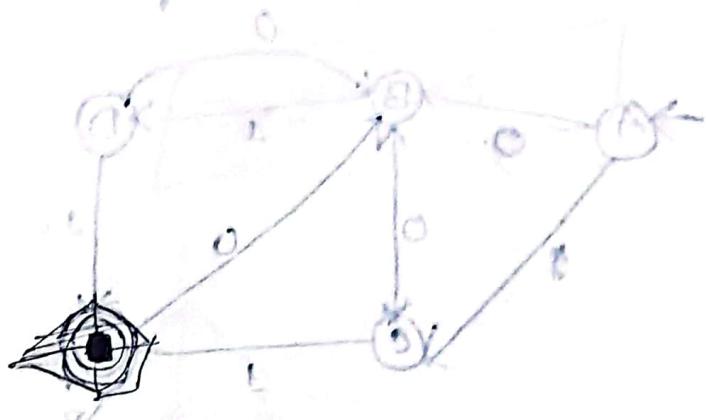
3 equivalent : $\{\underline{A, C}\}, \{\underline{B}\}, \{D\}, \{E\}$ C, D pair & 0 & 1 pair
also, C, D & D & E

E ∈ { } n² so close



~~AU 22~~ 3(b)

| | 0 | 1 |
|-------|-------|-------|
| 0 | q_1 | q_3 |
| 1 | q_2 | q_4 |
| q_1 | q_1 | q_4 |
| q_2 | q_2 | q_4 |
| q_3 | q_4 | q_4 |
| q_4 | q_4 | q_4 |



0 equivalent $\Rightarrow \{q_0, q_1, q_2, q_3\} \setminus \{q_4\}$

1 equivalent $\Rightarrow \{q_0, q_1, q_3\} \setminus \{q_2\} \setminus \{q_4\}$

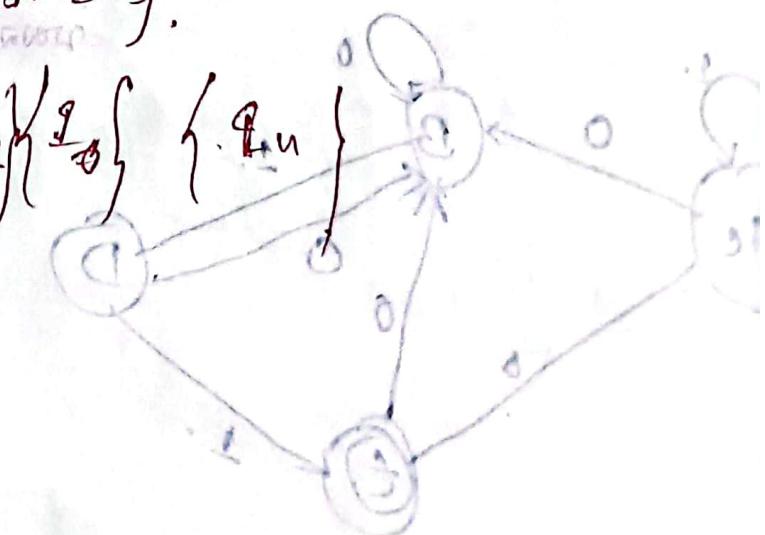
2 equivalent $\Rightarrow \{q_0, q_3\} \setminus \{q_2\} \setminus \{q_4\} \setminus \{q_1\}$

3 equivalent $\Rightarrow \{q_0, q_1, q_2, q_3\} \setminus \{q_4\}$

4 equivalent $\Rightarrow \{q_0\} \setminus \{q_0, q_1, q_3\}$.

$\Rightarrow \{q_0, q_1, q_3\} \setminus \{q_0\} \setminus \{q_4\}$

| | 0 | 1 |
|-------|-------|-------|
| 0 | q_0 | q_4 |
| 1 | q_0 | q_4 |
| q_0 | q_0 | q_4 |
| q_4 | q_4 | q_4 |



• Ignored number of states ~~NFA~~
 (Exact states which the machine moves can't be determined)

Phases of NFA and DFA

NFA

$$\begin{aligned}
 S &= \{a, b, c\} - \text{Non terminal} \\
 \Sigma &= \{0, 1\} - \text{terminal} \\
 Q_1 &= \{a\} \\
 \text{Initial} & \\
 F &= \{c\}.
 \end{aligned}$$

DFA

• The transition from state
 go to single particular next
 state for each input symbol.
 [Ex: for $a \rightarrow a, b, c$]

* empty string transition

Not shown

• Backtracking is allowed

• requires more space

• The transition
 multiple next state for each
 input symbol.

* ~~NFA~~ Empty string transition.

are shown here

• Backtracking not possible

• Requires less space

• Requires less space