

Language Translator

- (i) Assembler
- (ii) Interpreter
- (iii) Compiler

Source code / HLL code



Preprocessor

Compiler

Assembler

Linker / Loader



Machine code

High Level Language

↓
Pre processor

↓ pure HLL

↓ Compiler

↓ Assembly Lang.

↓ Assembler

↓ m/c code (relocatable)

↓ Loader/Linker

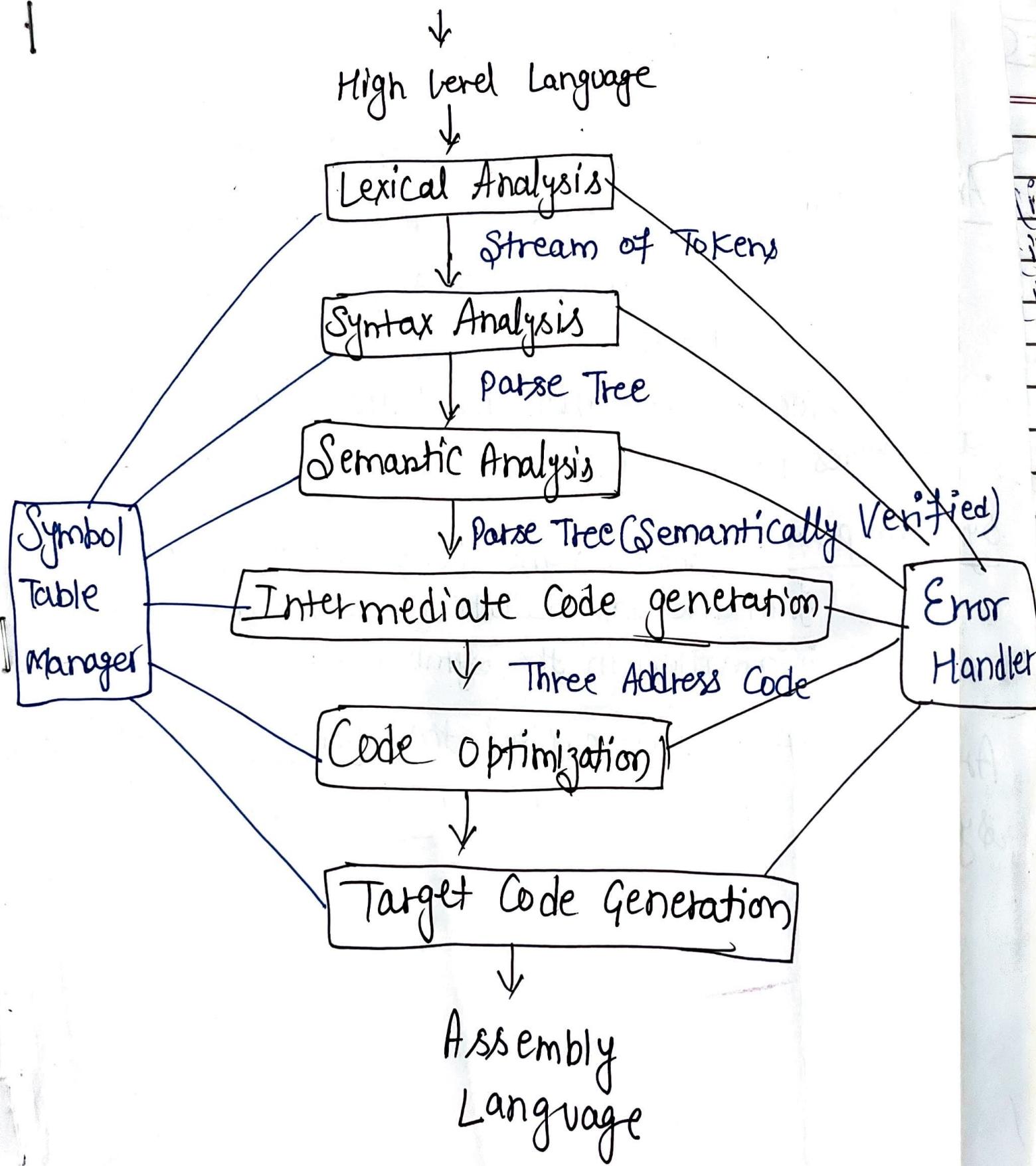
↓ Executable code/ absolute m/c cycle

↑ Preprocessor directive

HLL :- C, C++, Java, etc with Header files

↓ Include
 Preprocessor :- converts these

header files into the main program
and remove these header file inclusion.
and convert it into a pure file



The structure of a compiler

Two parts :- Analysis and Synthesis Part

Analysis Part I - Breaks up the source program into constituent pieces and imposes a grammatical structure on them.

- Also collects information about the source program and stores it in a data structure called a Symbol Table

Synthesis part:- • Constructs the desired target program from the intermediate representation and the information in the symbol table.

① Lexical Analysis :

- first phase of a Compiler, also called scanning
- Reads the stream of characters making up the source program and groups the characters into meaningful sequences called lexemes.

token form :-

$\langle \text{token-name}, \text{attribute value} \rangle$

- for each lexeme, the lexical analyzer produces an output a token.
- Example:- position = initial + rate * 60

Characters of this statement, grouped into lexemes:-

① lexeme Tokens
position token $\Rightarrow \langle \text{id}, 1 \rangle$ points to the
 Abstract symbol for "identifier" symbol table entry.

② = token $\langle \text{op}, 2 \rangle$

no attribute-value is needed.

③ initial token $\langle \text{id}, 3 \rangle$

④ + mapped into the $\langle + \rangle$ token.

⑤ rate token $\langle \text{id}, 4 \rangle$

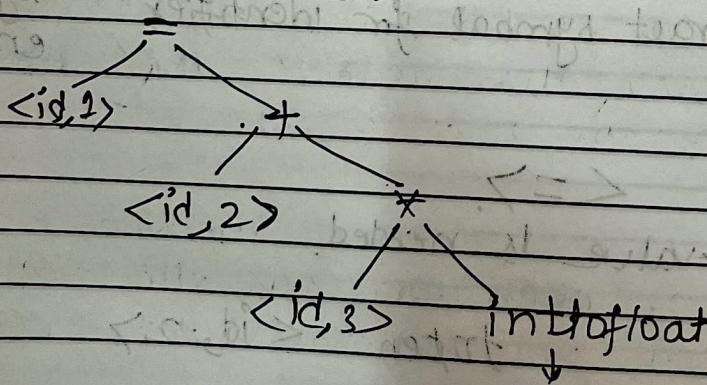
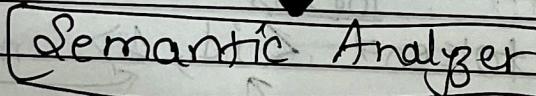
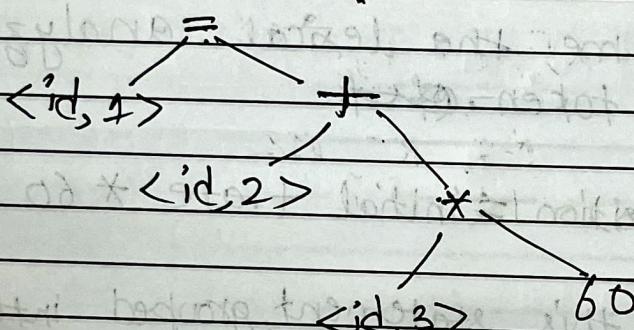
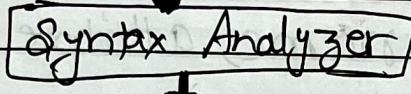
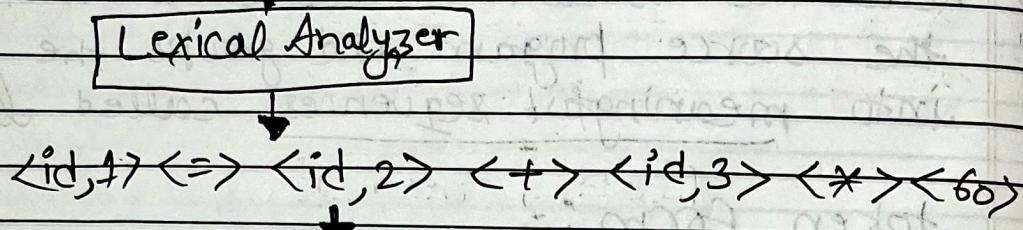
⑥ * token $\langle \ast \rangle$

⑦ 60 token $\langle 60 \rangle$

Blank space will be discarded by LA.

Translation of an Assignment Statement

position = initial + rate * 60



Three-Address
Code

$t1 = \text{inttofloat}(60)$

$t2 = id3 * t1$

$t3 = id2 + t2$

$s1 = t3$

1	position
2	initial
3	rate

Symbol
Table

- ② Syntax
- Second or Par
 - Uses represe
 - struc
 - Syntax represe
 - of th
 - the o
 - Tree shows
 - statem

↓

Code optimizer

1	position	..
2	initial	..
3	rate	..

symbol

Table

$t1 = id_3 * 60.0$
 $id_1 = id_2 + t1$
 ↓
Code Generator
 ↓
 LDF R2, id3
 MULF R2, R2 #60.0
 LDF R1, id2
 ADDF R1, R1, R2
 STF id1, R1

2) Syntax Analysis :

- Second phase of Compiler is Syntax Analysis or Parsing.
- Uses tokens to create a tree-like intermediate representation that depicts the grammatical structure of the token stream.
- Syntax Tree, in which each interior node represents an operation and the children of the node represent the arguments of the operation.
- Tree shows the order in which the operations in statement is given are to be performed.

(3) Semantic Analysis :

Semantic Analyzer uses the syntax tree and the information in the symbol table to check the source program for semantic consistency with the language definition.

- Also gathers type information and saves it in either the syntax tree or the symbol table, for subsequent use during intermediate code generation.
- Typechecking: Compiler checks that each operator has matching operands.

(4) Intermediate Code generation:

Many compiler generate an explicit low-level or machine-like intermediate representation after syntax and semantic analysis, which we can think of as a program for an abstract m/c.

- It should have 2 properties:
- It should be easy to produce.
- It should be easy to translate into the target machine.

$$t1 = \text{int/float}(6.0)$$

$$t2 = id3 * t1$$

$$t3 = id2 + t2$$

$$- id1 = t3$$

Noting points:- ① 8

② Inst

③ The compiler
the value

④

⑤ Code Optim
M/C Indep
intermediate
result.

• Better mean
power.

(6) Code generation:

I/P & output
and maps

• If the target
memory loc used by the

• Crucial asp
variables.

LDR

MUL

LDP

AD

STF

• first operand
• "F" denotes
• "#" signified

- Noting points:-
- ① Each instruction has atmost one operator ^R m
 - ② Instructions fix the order.
 - ③ The compiler must generate a temporary name to hold the value computed by a three-address instruction.

⑤ Code Optimization:

M/C Independent Code-optimization, improve the intermediate code so that better target code will result.

- Better means faster, shorter than consumes less power.

$$f1 = id3 * 60.0$$

$$id1 = id2 + f1$$

⑥ Code generation :- The code generator takes an I/p ~~is~~ an intermediate representation of source program and maps it into the target language.

- If the target lang. is m/c code, registers or memory locations are selected for each of the variables used by the program.
- Crucial aspect is assignment of registers to hold variables.

LDF R2, id3

MULF R1, R2, #60.0

LDF R1, id2

ADDF R1, R1, R2

STF id1, R1

- first operand of each instruction specifies a destination.
- "F" denotes that it deals with floating-point numbers.
- "#" signifies that '60.0' is to be treated as a constant.

→ Compiler- Construction TSD/8 :

Most successful tools are those that hide the details of the generation algorithm and produce components that can be easily integrated into the remainder of the compiler.

Tools - ① Parser Generators

Automatically produce syntax Analyzers from a grammatical description of a programming language. Example:- EAM, PTC

② Scanner generators

produces Lexical analyzers from a regular expression description of the tokens of a language.

Ex:- LEX

③ Syntax-directed translation engines.

produce collections of routines for walking a parse tree and generating intermediate code.

④ Code-generator

Produces a Code generator from a collection of rules for translating each operation of the intermediate lang. into the m/c lang. for a target m/c.

⑤ Data Flow analysis engines

- Facilitate the gathering of information about how values are transmitted from one part of a program to other part.

• key po

⑥ Compil
provide
various

→ The R

- Main task
character
lexeme
- whenever
identifier
symbol

Source → LA
program

Intermediate

- It performs
- ① Skipping a tab

- ② Correlating with the LA may keep characters error message

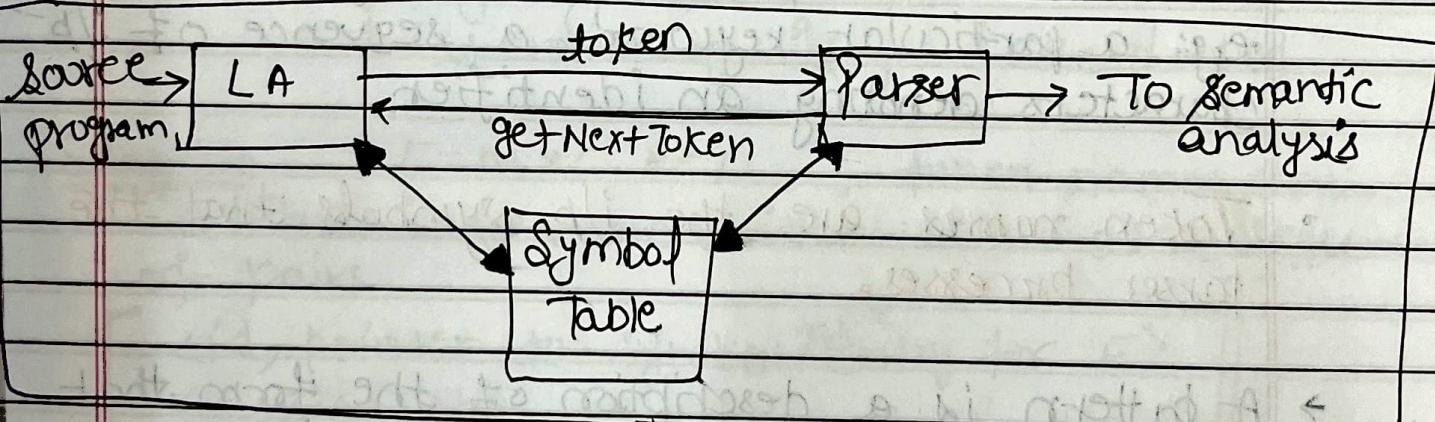
- Key part of code optimization,

⑥ Compiler-construction toolkits

provide an integrated set of routines for constructing various phases of a compiler.

⇒ The Role of Lexical Analyzer :

- Main task of lexical analyzer is to read the i/p characters of the source program, group them into lexemes & produce a sequence of tokens for each lexeme.
- Whenever LA discovers a lexeme constituting an identifier, it needs to enter that lexeme into the symbol table.



Interaction b/w the LA and the Parser.

- It performs certain other tasks :-
- ① Skipping out comments and whitespace (blank, newline, tab).
 - ② Correlating error messages generated by the compiler with the source program.
- LA may keep track of the number of newline characters seen, to relate the line no. with each error message.

Sometimes LAs are divided into cascade of two processes:-

- ① Scanning : simple processes that do not require tokenization, such as deletion of comments and compaction of consecutive whitespace characters into one.
- ② Lexical Analysis
 - More complex portion.
 - which produces tokens from the o/p of the scanning process.

① Lexical Analysis versus Parsing

2) ↳ Tokens, Patterns and Lexemes:

- A Token is a pair consisting of a token-name and an optional attribute value.
 - Token name represents a kind of lexical unit e.g. a particular keyword, a sequence of i/p characters denoting an identifier.
 - Token names are the i/p symbols that the parser processes.
- A pattern is a description of the form that the lexemes of a token may take.
 - In case of token keyword as a token, the pattern is just the sequence of the characters that form the keyword.
- A lexeme is a sequence of characters in the source program that matches the pattern for a token and is identified by the LA as an instance of an that token.

Example:- `printf ("Total = %d \n", score);`

Here, both printf and score are lexemes matching the pattern for token id and "Total = %d \n" is a lexeme matching literal.

Tokens	Informal description	Sample Lexemes
if	Characters i, f	if
else	characters e, l, s, e	else
comparison	< or > or <= or >= or == or !=	<=, !=
id	letter followed by letters and digits	hi, score, 72
number literal	any numeric constant anything but "surrounded by "" "	3.14159, 6, 6.02e23 "core dumped"

Example:- `E = m * c ** 2`, token names and attribute values are written below as a sequence of pairs

<id, pointer to symbol table entry for E>
 <assign_op>
 <id, pointer to symbol table entry for m>
 <molt_op>
 <id, pointer to symbol table entry for c>
 <exp_op>
 <number; integer value 2>

③ ↳ Attributes for Tokens

When more than one lexeme can match a pattern, the LA must provide the compiler phases additional information about the particular lexeme that matched.

The pattern for token number matches both 0 and 1.

LA returns \langle token name, attribute value \rangle that describes the lexeme represented by token. Token-name influences parsing decisions.

In certain pairs, especially operators, punctuation and keywords, there is no need for an attribute value.

④ Lexical Errors

- Suppose a case arises in which LA is unable to proceed because none of the patterns for tokens matched any prefix of the remaining i/p.

Simplest recovery strategy is "panic mode" recovery. LA delete successive characters from the remaining i/p, until the analyzer can find a well-formed token at the beginning of what i/p is left.

Other act

- ① Delete
- ② Insert a
- ③ Replace a
- ④ Transpo

⑤ Input Buff

Tale cannot
an identifier
is not a
part of th

- A single ch
also be the
operator
- for this w
handled

⑥ Buffer pa

Introduced
process a
alternately

Two pointe

- ① Pointer l
marks t
whose ext
- ② Pointer for
Scans a

Other actions are:-

- ① Delete one character from the remaining i/b
- ② Insert a missing " " into the remaining i/p.
- ③ Replace a character by another
- ④ Transpose two adjacent characters.

The lexical - Analyzer Generator - Lex

- Lex or FLEX, allows to specify a lexical analyzer by specifying regular expressions to describe pattern for tokens.
- The I/p notation for the Lex tool is referred to as the Lex language and the tool itself is the Lex Compiler.
- Lex Compiler ~~transforms~~ the i/p pattern into a transition diagram and generates code, in a file called Flex.yy.c that simulates this transition diagram.

Finite Automata

Lex turns its I/P program into a FA.

① Finite automata are recognizers, they simply "yes" or "no" about each possible string.

② Two types of FA

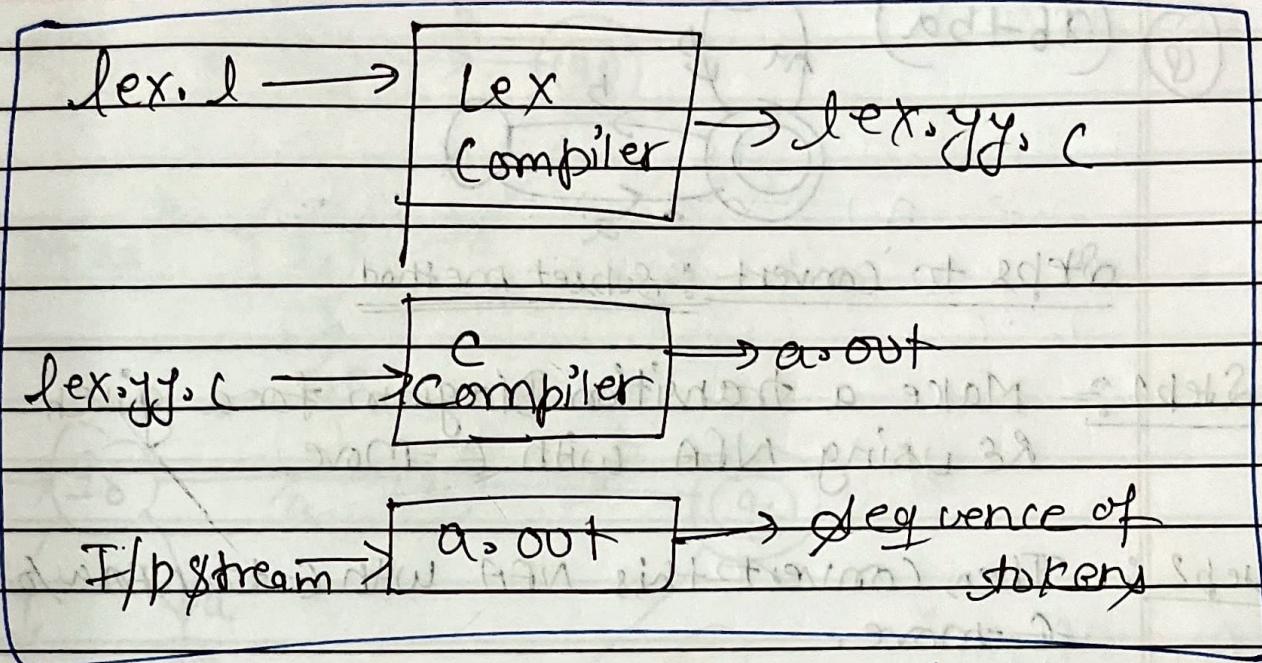
a) (NFA) Non-Deterministic Finite Automata
No restrictions on the labels of their edges.

A symbol can label several edges out of the same state, and ϵ the empty string is a possible label.

b) Deterministic finite Automata (DFA)

For each state and for each symbol of its I/P alphabet exactly one edge with that symbol leaving that state.

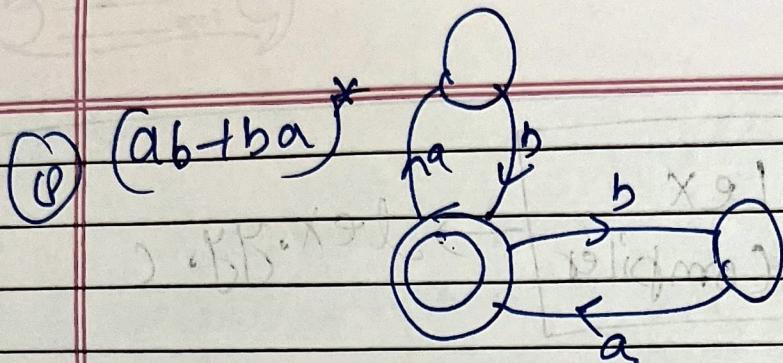
Both are capable of recognizing the same languages.



Creating a LA with lex tool

⇒ Regular Expression to FA

- ① \emptyset → 0
- ② a → 0 → a
- ③ $a+b$ → 0 → a, b
- ④ ab → 0 → a → b
- ⑤ a^* → 0 → a
- ⑥ ab^* → 0 → a → 0 → b
- ⑦ $(ab)^*$ → 0 → a → 0 → b



Steps to Convert :- Subset method

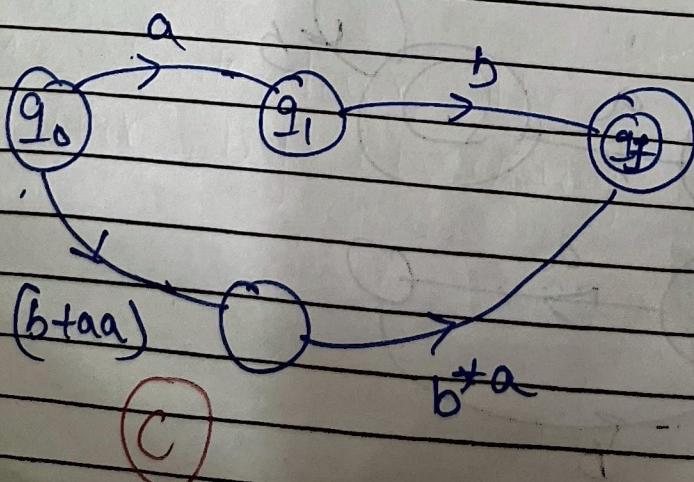
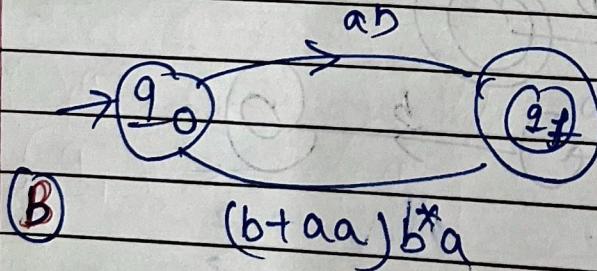
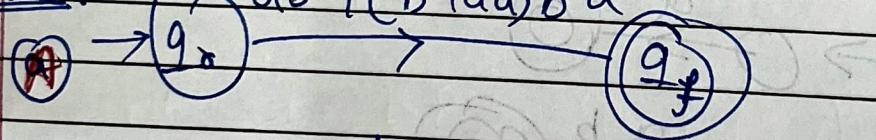
Step 1 :- Make a transition Diagram for a given RE using NFA with ϵ -move.

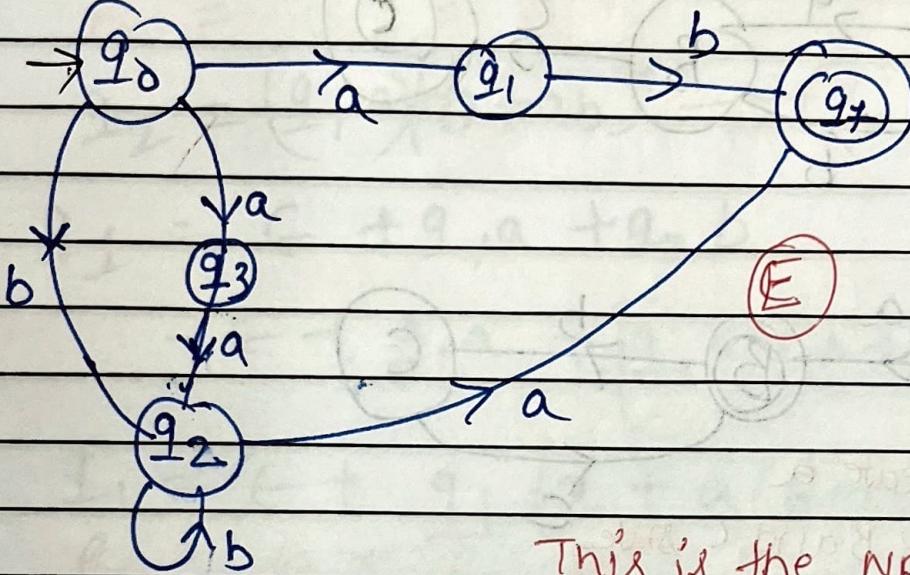
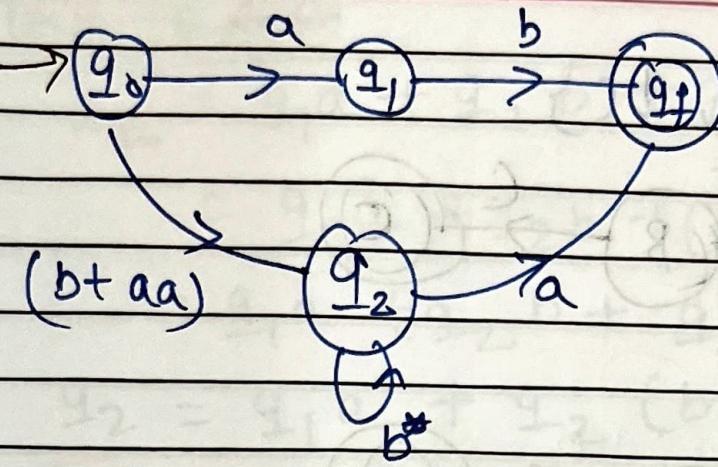
Step 2 :- Then convert this NFA with ϵ to NFA w/o ϵ -move.

Step 3 :- Finally, convert the obtained NFA into DFA.

Ex:- RE = $(ab + (b+aa)b^*a)$ into FA.

Step 1 :- $ab + (b+aa)b^*a$



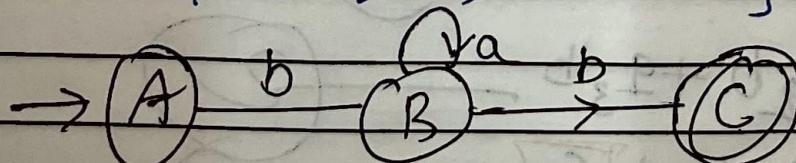


This is the NFA w/o ϵ move,
now we will convert it into required DFA if ask
to do so.

NFA		DFA			
State	a	b	State	a	b
$\rightarrow q_0$	$\{q_1, q_3\}$	q_2	$\rightarrow q_0$	$\{q_0, q_3\}$	q_2
q_1	\emptyset	q_f	q_1	$\{q_1, q_3\}$	q_2
q_2	q_f	q_2	q_2	q_f	q_2
q_f	\emptyset	\emptyset	q_f	\emptyset	\emptyset
q_3	q_2	\emptyset	\emptyset	\emptyset	\emptyset

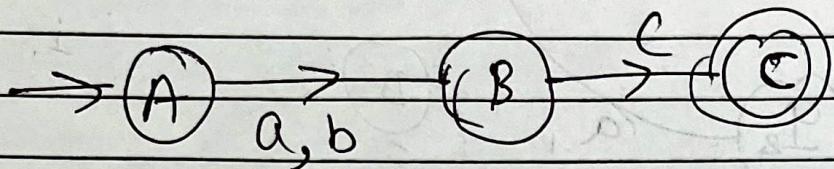
RE: $b \ a^* b$

$L = \{baab, bb, bab -\}$

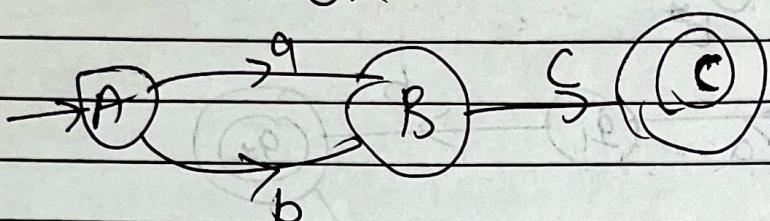


$$\textcircled{Q} - (a+b)c$$

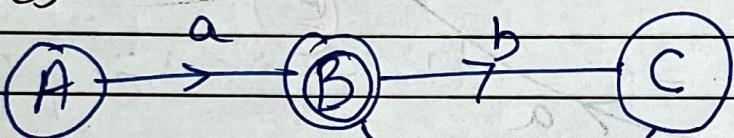
$$L = \{ac, bc\}$$



OR



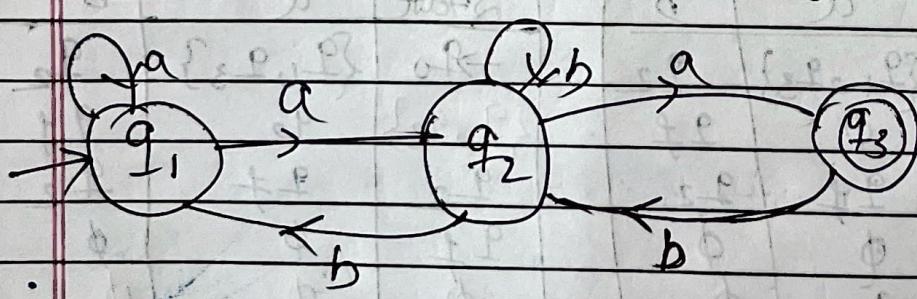
$$\textcircled{Q} - a(bc)^*$$



bc^* will create a self loop btw B and C states.

⇒ ~~logics~~ Designing RE

NFA to RE using Arden's Theorem



Step 1 - Write down the equations for each state, by checking incoming edges.

$$q_3 = q_2 \cdot a \quad \text{--- } \textcircled{1}$$

$$q_2 = q_1 \cdot a + q_2 \cdot b + q_3 \cdot b \quad \text{--- } \textcircled{2}$$

$$q_1 = q_1 \cdot a + q_2 \cdot b + \boxed{E} \quad \text{--- } \textcircled{3}$$

DFA to RE

$$q_2 = q_1 a + q_2 b + q_3 b$$

$$= q_1 a + q_2 b + (q_2 \cdot a) b \quad [q_3 \text{ from eq.1}]$$

$$= q_1 a + q_2 b + q_2 ab$$

$$q_2 = q_1 a + q_2 (b+ab)$$

$$R = q + RP \quad [R = q + RP]$$

$$q_2 = (q_1 a)(b+ab)^* \quad \text{--- (4)}$$

$$q_1 = \epsilon + q_1 a + q_2 b$$

$$= \epsilon + q_1 a + [q_1 a (b+ab)^*] b \quad (\text{from eq.4})$$

$$q_1 = \epsilon + q_1 [a + a(b+ab)^*] b$$

$$R = q + RP \quad [R = q + RP]$$

$$q_1 = \epsilon + a + a$$

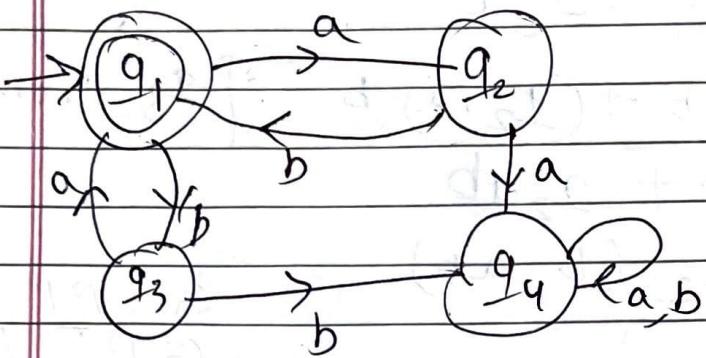
$$q_1 = (a + a(b+ab)^* b)^* \quad \text{--- (5)}$$

$$q_3 = q_2 a$$

$$q_3 = [(q_1 a)(b+ab)^*]^a$$

$$q_3 = [(a + a(b+ab)^* b)]^a \cdot (b+ab)^*$$

RE

DFA TO RE

R.E = q_1 / final state

$$q_1 = \epsilon + q_2 b + q_3 a \quad \text{--- (1)}$$

$$q_2 = q_1 a \quad \text{--- (2)}$$

$$q_3 = q_1 b \quad \text{--- (3)}$$

$$q_4 = q_2 a + q_3 b + q_4 a + q_4 b \quad \text{--- (4)}$$

→ Analyze and solve it,

$$q_1 = \epsilon + q_2 b + q_3 a$$

$$q_1 = \epsilon + (q_1 a)b + (q_1 b)a$$

~~$$q_1 = \epsilon + q_1 ab + q_1 ba$$~~

$$q_1 = \epsilon + q_1(ab+ba)$$

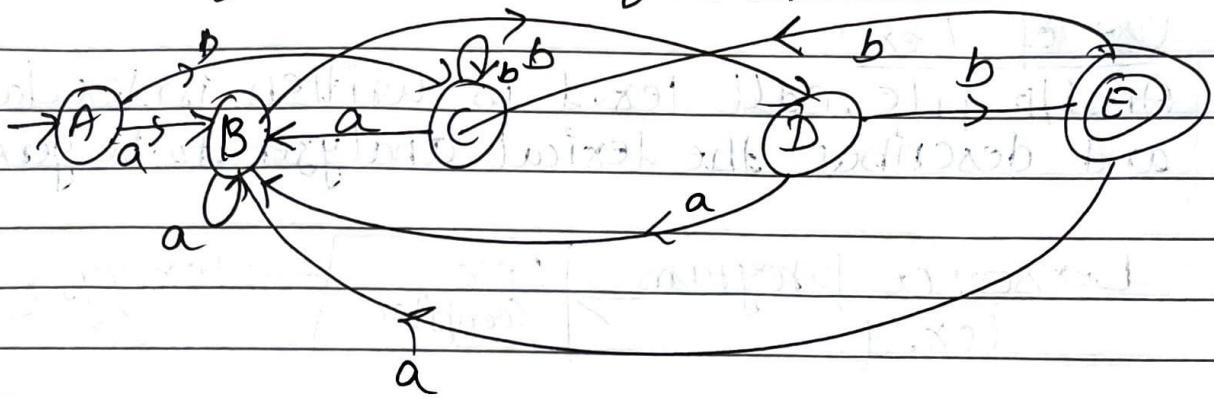
$$R = Q + R P$$

$$q_1 = \epsilon + (ab+ba)^*$$

$$q_1 = (ab+ba)^*$$

$$R.E = q_1 = (ab+ba)^*$$

Q Convert DFA to minimized DFA



~~patterns~~ → {E} ∪ {AB, C}, {B} ∪ {D}

Patterns for tokens

digit → [0-9] number

digits → digit

number → digit (digits)? (E [+ -] digits)?

letter → [A-Z a-z]

id → letter (letter | digit)*

then → then

else → else

relOp → < | > | <= | >= | != | <= | >= | = | < | >

ws → (blank | tab | newline)*

Lexeme, TokenName → attribute value

Any ws

if

then

else

Any id

Any number

<

>

=

if

then

else

id

number

relOp

relOp

relOp

blank or tab or newline

blank or tab or newline

pointer to table entry

pointer to table entry

LJ

LE

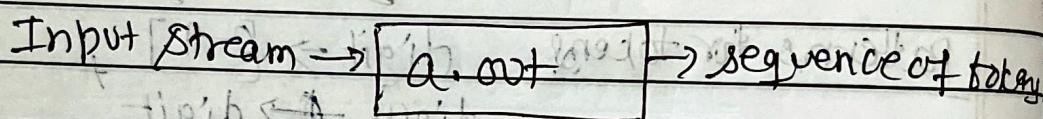
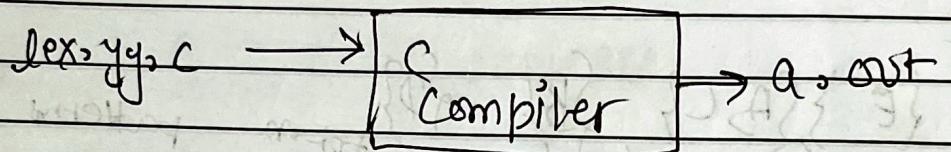
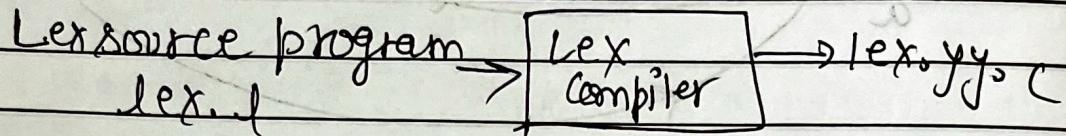
GE

EA

The Lexical-Analyzer Generator Lex

① Use of Lex

An i/p file, call lex.l is written in lex language and describes the lexical analyzer to be generated.



Creating a lexical analyzer with Lex

- Lex Compiler transforms lex.l to a C program, in a file that is always named lex.yy.c
- Later this file is compiled by C Compiler into a file called a.out as always.
- The C-compiler output is a working lexical analyzer that can take a stream of i/p char. and produce a stream of tokens.
- a.out is as a subroutine of the parser, it is a C function that returns an integer, which is a code for one of the possible token names.

- The attribute value, a numeric code, a pointer to the symbol table, or nothing is placed in a global variable `yyval`, shared btw the LA and parser

② Structure of Lex program

declarations

%%

translation rules

%%

auxiliary functions

// declarations of variables, constants

Translation rules each have the form:-

Pattern { Action }

Each Pattern is a regular expression, may be use the regular definitions of the declaration section.

The actions are fragments of code, typically written in C,

The 3rd section holds additional functions are used in actions, these functions can be compiled separately and loaded with the LA.

LA created by Lex behaves with the parser as follows:
When called by parser, the LA begins reading its remaining input, 1 character at a time.

Until it finds the longest prefix of the i/p that matches one of the patterns P_i ,

It then executes the associated action. ~~if~~

~~If~~ will return to the parser.

- But if it does not (because P_i describes comment or whitespace), then the LA parser proceeds to find additional lexemes until one of the corresponding actions causes a return to the parser.
- Until LA returns a single-value, the token-name but gives the shared integer variable yyval to pass additional information about the lexeme found, if needed.

Directly copied to the file lex.y, C and not treated as regular definition

/* { /* definition of manifest constants */

LT, LE, EQ, NE, GT, GE, IS, THEN,

ELSE, ID, NUMBER, RELOP */

}%

/* regular definitions */

delim

[\t\n]

ws

{ delim } +

letter

[A-Z a-z]

digit

[0-9]

id

{ letter } { letter } { digit } *

number

{ digit } + [. { digit }] + ?

(E { + - } ? { digit } *) ?

pattern

Action

return
token

IF

{WS}

{/* no action and no return */ }

if

{return (IF); }

then

{return (THEN); }

else

{id}

{yyval = (int) installID(); return (ID); }

{number}

{yyval = (int) installNum(); return (Number); }

<

{yyval = LT; return (RELOP); }

%%

int installID() /* function to install the lexeme
 whose first character is pointed to by
 yytext, and whose length is yylen, into
 the symbol table and return a pointer
 there to */

}

int installNum() /* similar to installID, but puts
 numerical constants into a
 separate table */

}

Anything in the auxiliary section is copied directly to file lex.yy.c, but may be used in actions.

Action taken when id is matched is threefold:

① Function `installID()` is called to place the lexeme found in the symbol table.

② `f` returns a pointer to the symbol table, which is placed in global variable `yyval`, `f` is available two variables of LA :-

③ `yytext` : pointer to the beginning of the lexeme similar to `beginptr`.

④ `yylen` : length of the lexeme found.

⑤ Token name is returned to the parser.

Conflict Resolution in Lex

Ruled that Lex used to decide a proper lexeme to select, when several prefixes of ip match one or more patterns :-

① Always prefer a longer prefix.

② If the longest prefix matches two or more patterns, prefer the pattern listed first in the Lex program.

Design
Gen

① Structure

Lex →
Program

A Lex pro
table an
finite-ai

The rest
that are
by Lex it

① A transiti

② Those f

Lex to t

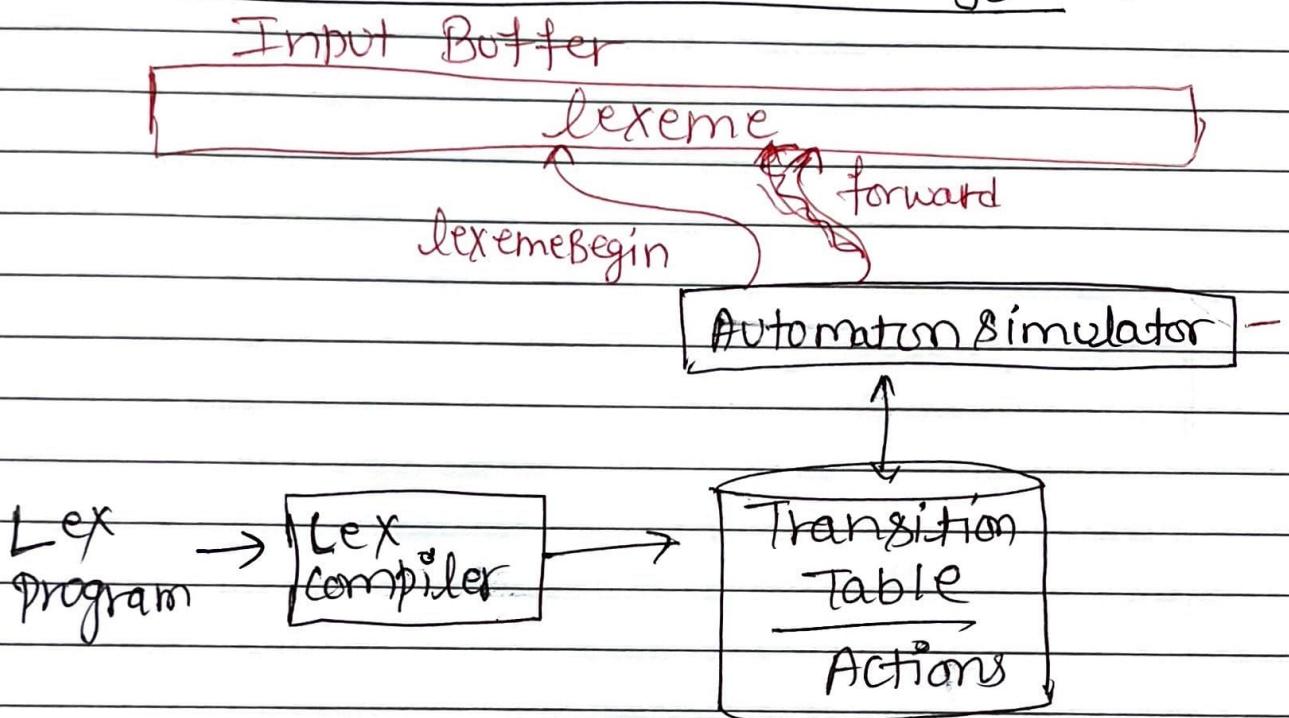
③ Actions fr

of code to

by to aut

Design of a lexical Analyzer Generator

① Structure of the Generated Analyzer :



- A Lex program is turned into a transition table and actions, which are used by a finite-automaton simulator.

The rest of the LA consists of components that are created from the Lex program by Lex itself. Components are:-

- ① A transition table for the automaton.
- ② These fs that are passed directly through Lex to the output.
- ③ Actions from the i/p program, fragments of code to be invoked at appropriate time by the automaton simulator.