

①

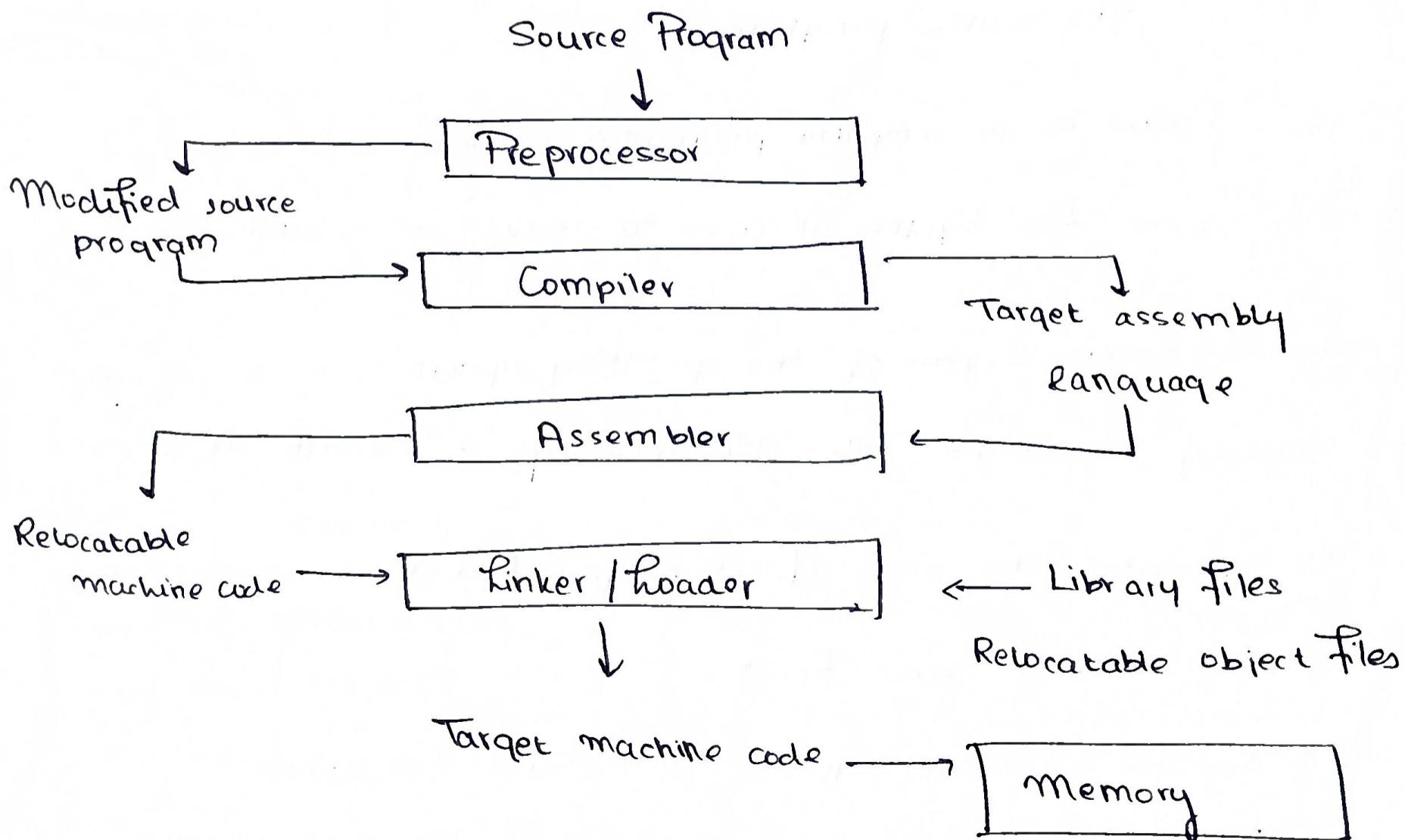
# UCS2702 COMPILER DESIGN

Unit 1  
mmmm

## Introduction to Compilers mmmmmm

Language processors - phases of compiler - role of lexical analyzer -  
 Input buffering - specification of tokens - recognition of tokens -  
 conversion of RE to DFA; Lexical analyzer generator; Structure of  
 lex program - lookahead operator and conflict resolution

### \* Language Processors



Translators : source language → Target language

Compiler : High level language → Target / Machine / Assembly  
 FORTRAN, PASCAL,  
 C/C++  
 lang  
 (Compiler also gives error msgs)

Language Processors : Assembly → Machine language  
 language

Preprocessor : High-level language → another high level language

Interpreter : Input → intermediate code

Instead of output, performs operations implied by  
 the source program

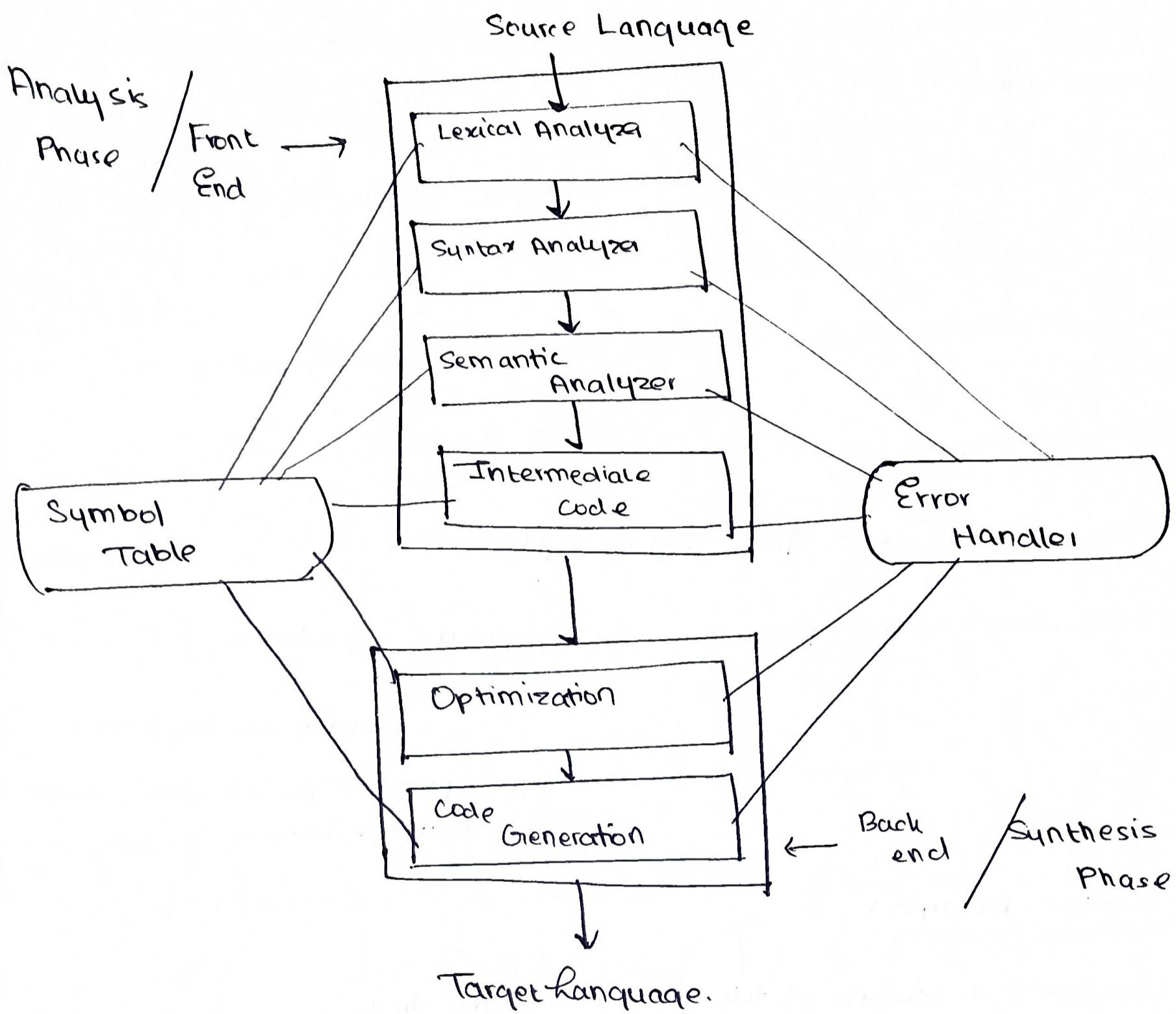
Linker: Linker is a computer program that links and merges  
 various object files together in order to make an executable file.

Loader : → Loader is part of the operating system & is responsible  
 for loading executable files into memory & execute them.

→ It calculates the size of a program (instructions and data)  
 & creates memory space for it

→ Initializes various registers to initiate execution

## \* Phases of a Compiler



There are 2 main phases of a compiler

### ① Analysis Phase

- an intermediate representation is created from the given source program
- Lexical analyzer, syntax analyzer, semantic analyzer & intermediate code generator are the parts of this phase

Q

## Synthesis Phase

- the equivalent target program is created from this intermediate representation
- Code generator and code optimizer are the parts of this phase.

The individual stages in more detail are as follows:

### A. Lexical Analyzer

~~~~~

→ reads characters from left to right

→ converts source program into a stream of tokens  
↓  
identifiers, key words,  
operators, punctuation symbols,  
multi-character operators

### B. Syntax Analyzer

~~~~~

→ converts a stream of tokens into a parse tree

Performs 2 functions:

(i) Checks if the token if they occur in patterns specified by the source language

(ii) Construct a tree-like structure that can be used by the subsequent phases

(Parser also known as  
hierarchical analysis)

The hierarchical structure of programs is usually described by recursive rods as follows:

1. Any identifier is an expression
2. Any number is an expression
3. If expression<sub>1</sub> & expression<sub>2</sub> are expressions, then so are:
  - a. Expression<sub>1</sub> + Expression<sub>2</sub>
  - b. Expression<sub>1</sub> \* Expression<sub>2</sub>
  - c. (Expression<sub>1</sub>)
4. If identifier<sub>1</sub> is an identifier and expression<sub>2</sub> is an expression  
then:  
Identifier := Expression<sub>2</sub> is a statement
5. If Expression<sub>1</sub> is an expression & Statement<sub>2</sub> is a statement,  
then:
  - a. while (Expression<sub>1</sub>) do Statement<sub>2</sub>
  - b. if (Expression<sub>1</sub>) then Statement<sub>2</sub> are statements
  - c. Semantic Analyzer - checks for more static semantic errors  
parse tree → modified / annotated parse tree
- D. Intermediate Code Generator - variety of intermediate representations

## Properties of 3 Address Code

- at most one operator other than the assignment operator
- Compiler must generate a temporary name to hold the value computed by each instruction.
- Some 3 address instruction can have less than 3 operands

E) Code Optimizer

    mm mmm

    Tries to run code to

        → run faster

        → be smaller

        → consume less energy

F) Code Generator

    mm mmm

    → converts optimized code into target program

## Other components

① **Symbol Table** → keep track of names declared in the program

    → separate levels for each scope

    → Linear list → slow, but easy to implement

    Hash Table → complex to implement, but fast.

② **Error Handler**

    - detection of errors

    - reporting errors

    - recovery of errors

## Applications of Compilers

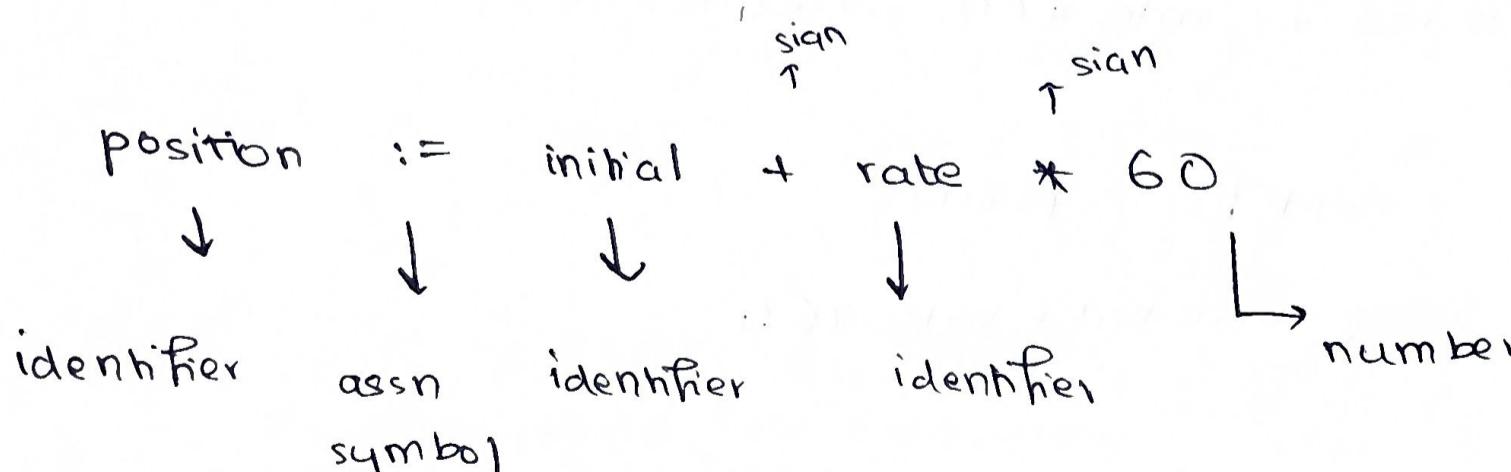
- Techniques used in a lexical analyzer can be used in text editors, information retrieval system & pattern recognition programs
- Parser - query processing system such as SQL
- Many software having a complex front-end may need techniques used in compiler design
- Most techniques used in compiler design can be used in NLP systems

(-X.)

Example: For the given expression, show its transformation through each stage

position := initial + rate \* 60

Step 1 : Lexical analyzer } L → R parse - get tokens



Step 2 : Syntax Analyzer → check patterns, construct parse tree

Referring the rules on pg. 5

Position := initial + rate \* 60

By Rule ① position  
initial      }  
rate              }  $\Rightarrow$  expression

By Rule ② 60 }  $\Rightarrow$  expression

By Rule ③ a exp1 = rate

exp2 = 60

rate \* 60  $\Rightarrow$  expression

By rule ③ b exp1 = rate \* 60

exp2 = initial

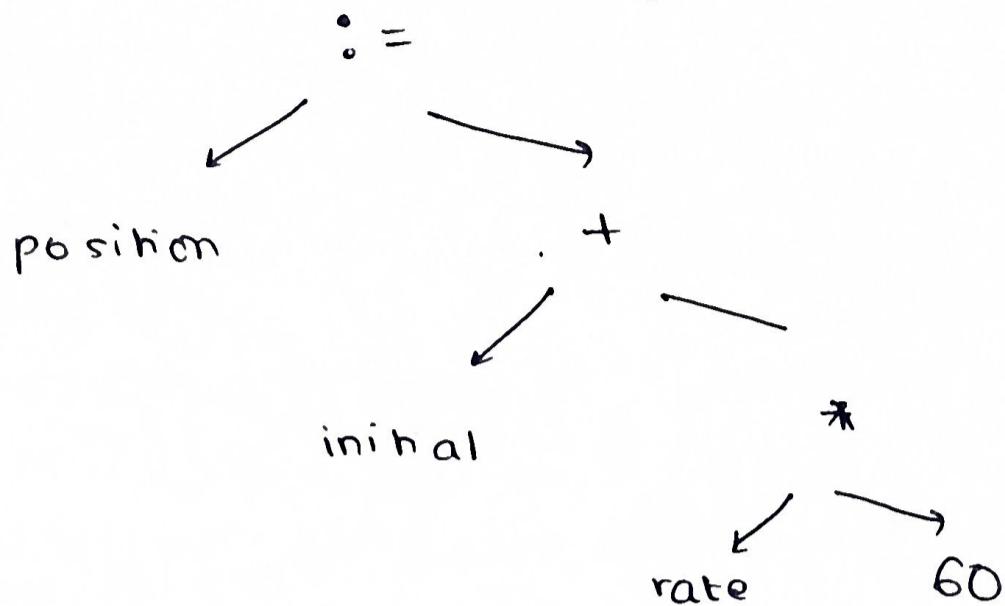
initial + rate \* 60  $\Rightarrow$  expression

By rule ④ identifier = position

exp = initial + rate \* 60

$\Rightarrow$  position := initial + rate \* 60  $\Rightarrow$  statement

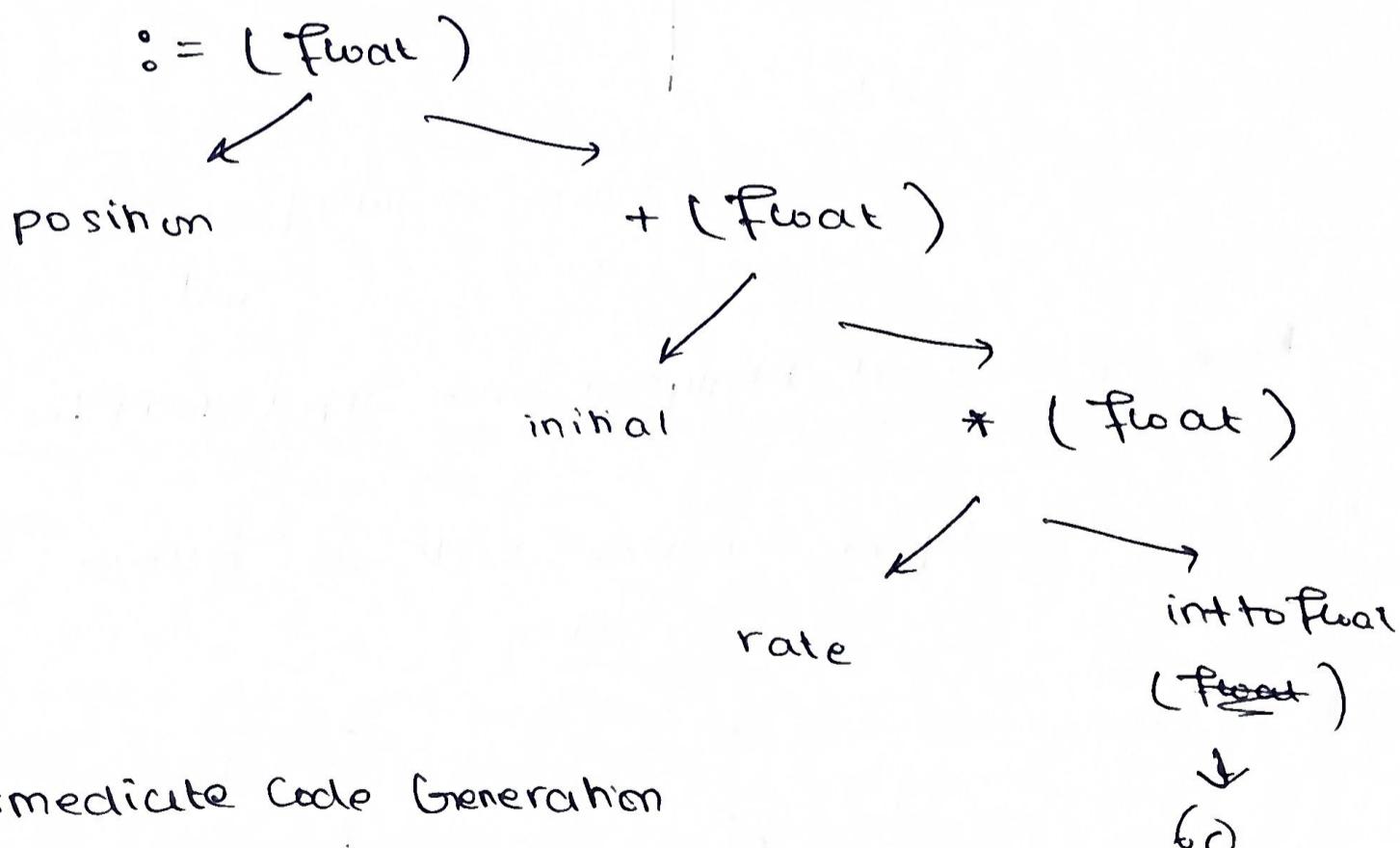
Now constructing a parse tree



Step 3 : Semantic Analyzer

→ convert to parse tree with annotations

(write type / type conversions)



Step 4: Intermediate Code Generation

(use the properties of 3 address code)

temp1 := int to float (60)

temp2 := rate \* temp1

temp3 := initial + temp2

position := temp3

## Step 5: Code Optimizer

temp1 : id3 \* 60.0  
 id1 : id2 + temp1  
 position → initial

Move vars to registers R1 & R2

OP stored in record → reg

## Step 6: Code Generator

MOVF id3, R2 # move value of ~~register~~ id3 to register R2

MULF ~~id2~~ #60.0, id3

(with rate value) into ~~id2~~ id3

Registers

MOVF id2, R1 # move value of ~~register~~ id2 (w/ initial value) into ~~id2~~ reg R1

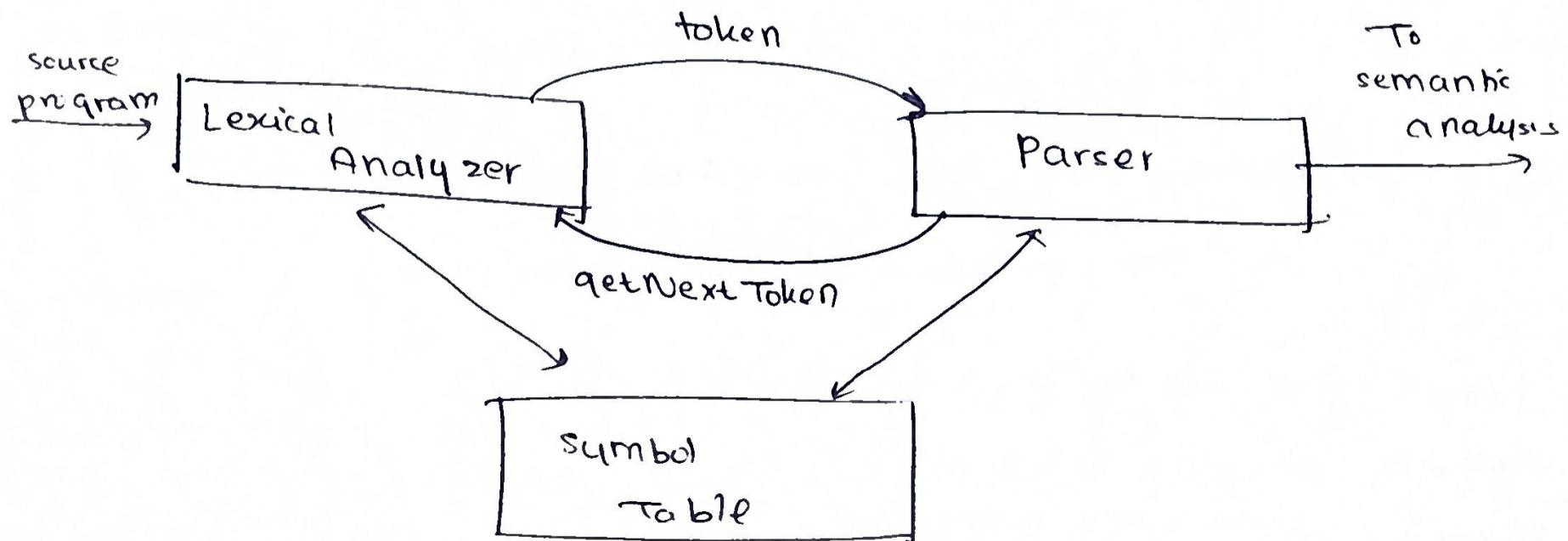
ADDF R2, R1

MOVF R1, id1

move result of addition stored in R1 to id1

## \* Role of Lexical Analyzer and Input Buffering

A lexical analyzer & parser work as follows:



## \* Why separate Lexical analysis & parsing

- ① Simplicity of Design → Removing white space by Lexical analyzer allows for easy implementation of parsing
- ② Improving Compiler Efficiency → Large amount of time is spent in reading the source program. Buffering techniques are used for reading input characters
- ③ Enhancing Compiler Portability → The representation of special symbols can be isolated in Lexical analyzer.

## \* Tokens, Patterns and Lexemes

- ① Token - A token is a pair of a token name and an optional token value.
- ② Pattern - A pattern is a description of the form that the lexemes of a token may take.
- ③ Lexeme - A Lexeme is a sequence of characters in the source program that matches the pattern for a token.

<u>e.g.</u>	<u>Token</u>	<u>Lexemes</u>
	if	if
	else	else
	comparison	=, <, >
	id	pi, score, variable
	number	1, 3.14, -6

## \* Attributes for Tokens

e.g.  $E = M * C \ast \ast 2$

- <id, pointer to symbol table entry for E>
- <assign-op>
- <id, pointer to symbol table entry for M>
- <mult-op>
- <id, pointer to symbol table entry for C>
- <exp-op>
- <number, integer value 2>

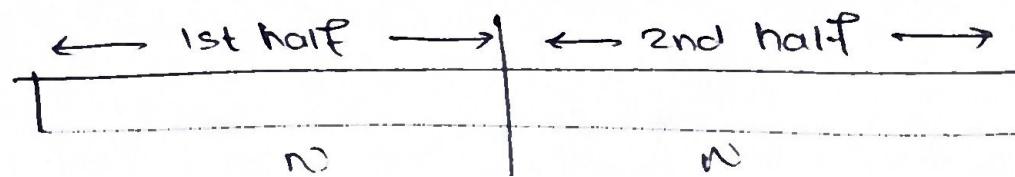
## \* Input Buffering

- Sometimes, the lexical analyzer needs to look ahead some symbols to decide about the token to return.
- A two-buffer schema is introduced to handle large lookahead safely

## Buffer Pairs

- The buffer is divided into two N-character halves.

N is the no. of characters 1024 or 4096



## Buffer Pair Algorithm

There are two pointers - lexeme-beginning & forward

lexeme-beginning is moved only once a token has been recognized. It then moves to the next character to parse.

forward moves one step forward for every character

### Algorithm

```
IF forward at end of first half then begin
    reload second half
    forward := forward + 1
end
```

```
else if forward at end of second half then begin
    reload first half
```

```
    move forward to beginning of first half
```

```
end
```

```
else forward := forward + 1
```

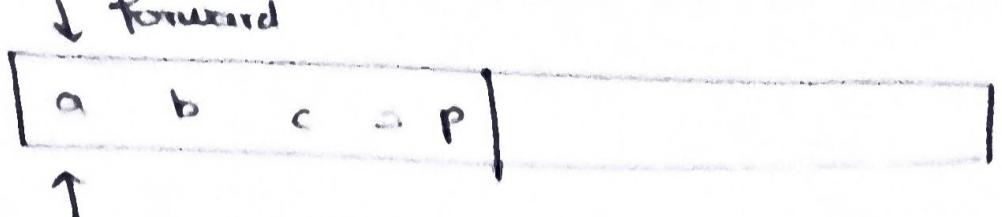
Example consider  $abc = pqr \times 24^2$

divide the characters into half

First half of buffer = abc = P

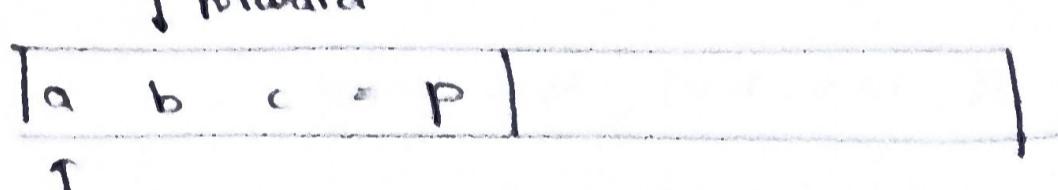
second half of buffer = qr  $\times 24^2$

(i)



lexeme - beginning

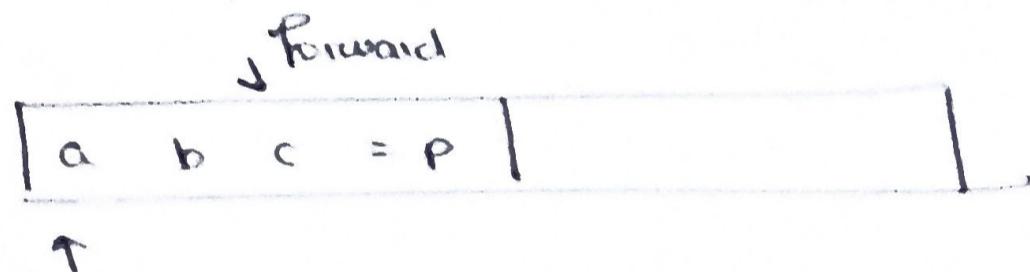
(ii)



lexeme -

beginning

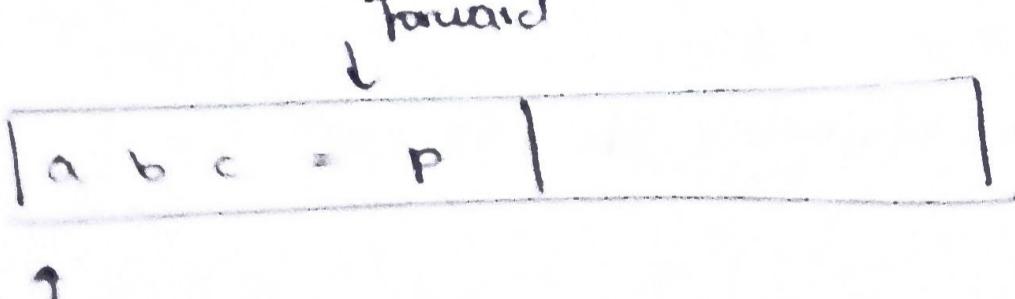
(iii)



lexeme

beginning

(iv)

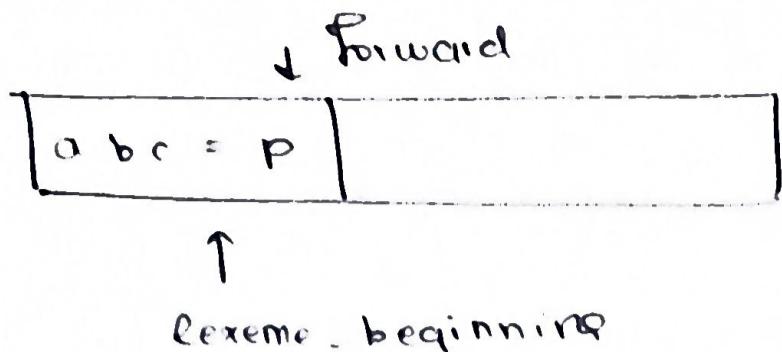


lexeme beginning

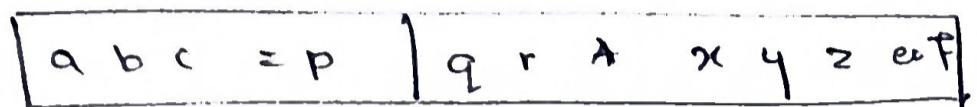
=  $\Rightarrow$  assignment operator

$\Rightarrow abc \rightarrow$  identifier

(v)



(vi) now load 2nd half



continue pointer movement

**\* Sentinels**

→ Sentinels are used to mark the end of input in the buffer without requiring explicit checks after each character.

Algorithm

move forward pointer & check  
character at location

Switch (\* forward ++ ) {

case eof :

if (forward is at the end of first buffer) {

reload second buffer

forward := beginning of second buffer;

{

else if (Forward is at the end of second buffer) {  
    reload first buffer  
    Forward := beginning of first buffer  
}

else {  
    # means eof within buffer  $\Rightarrow$  marks end of input  
    terminate lexical analysis  
    break;  
}

case for other characters

#### \* Lexical Errors

→ Some errors are out of power of the Lexical analyzer to recognize:

eg  $f_i(a == f(x))$

→ However, it may be possible to recognize errors like,

$d = \partial x$

Such errors are recognized when no pattern for tokens matches a character sequence.

## \* Error Recovery

- Panic mode : successive characters are ignored until we reach a well-formed token
- Delete one character from the remaining input
- Insert a missing character into the remaining input
- Replace a character by another character
- Transpose two adjacent characters

Example : Tokenize the following expression

a + b / 2.6

a : Identifier

+ : Operator

b : Identifier

/ : Operator

2.6 : Numeric literal

## \* Specification of Tokens

### ① Alphabet or Character Class

~~~~~ m ~~~~~ m

→  $\Sigma$  is a finite set of symbols

→  $\{0,1\}$  is a binary alphabet

### ② String

→ A finite sequence of symbols from  $\Sigma$

→  $|s|$  = length of string

→  $\epsilon$  = empty string  $| \epsilon | = 0$

### ③ Language

→ A language is a specific set of strings over some fixed alphabet  $\Sigma$ .

### More Specifications

① Prefix of  $s$  - A string obtained by removing 0 or more trailing symbols of  $s$

② Suffix of  $s$  - A string formed by deleting 0 or more leading symbols of  $s$

③ Substring of  $s$  - A string obtained by removing the suffix and prefix from  $s$

④ Proper prefix and proper suffix - Any prefix or suffix other than the string itself

⑤ Subsequence of  $s$  - Any string formed by deleting zero or more not necessarily contiguous symbols from  $s$

banana  $\rightarrow$  baas

### \* Language Operations

① Union :  $L \cup M = \{s \mid s \in L \text{ or } s \in M\}$

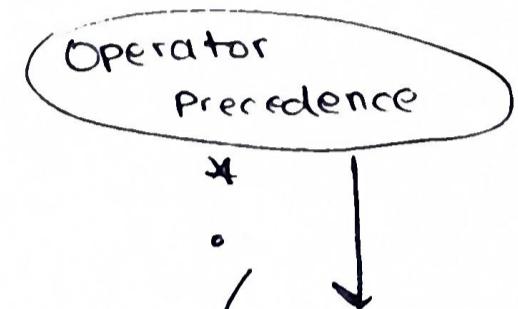
② Concatenation:  $L \cdot M = \{xy \mid x \in L \text{ & } y \in M\}$

③ Kleene closure:  $L^* = \bigcup_{i=0 \dots \infty} L^i$

④ Positive closure:  $L^+ = \bigcup_{i=1 \dots \infty} L^i$

## \* Regular Expressions

1.  $\epsilon$  is a regular expression,  $L(\epsilon) = \{\epsilon\}$
2. If  $a$  is a symbol in  $\Sigma$ , then  $a$  is a regular expression  
 $L(a) = \{a\}$
3.  $(r) / (s)$  is a RE denoting the language  $L(r) \cup L(s)$
4.  $(r)(s)$  is a RE denoting the language  $L(r)L(s)$
5.  $(r)^*$  is an RE denoting  $(L(r))^*$
6.  $r$  is an RE denoting  $L(r)$



## \* Recognition of Tokens

Step 1: Find the starting point, to understand the grammar the tokens denote

e.g.  $\text{stmt} \rightarrow \text{if expr then stmt} \mid \text{if expr then stmt else stmt}$

expr

$\text{expr} \rightarrow \text{term relop term} \mid \text{term}$

$\text{term} \rightarrow \text{id} \mid \text{number}$

Step 2: Formalize the patterns

9. digit  $\rightarrow [0-9]$

digits  $\rightarrow \text{digit}^*$

number  $\rightarrow \text{digit} (\cdot \text{digit})^* ( \in [0-7] ? \text{digit} )^*$

letter  $\rightarrow [A-Z a-z]$

id  $\rightarrow \text{letter} (\text{letter} | \text{digit})^*$

if  $\rightarrow \text{if}$

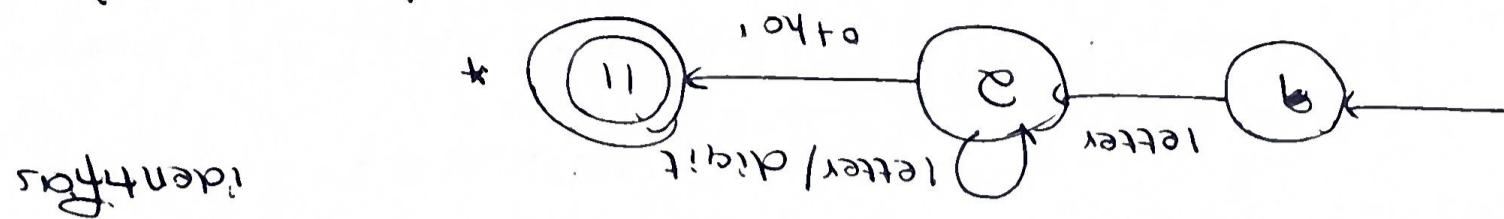
then  $\rightarrow \text{then}$

else  $\rightarrow \text{else}$

relop  $\rightarrow < | \geq | = | \leq | \geq | < >$

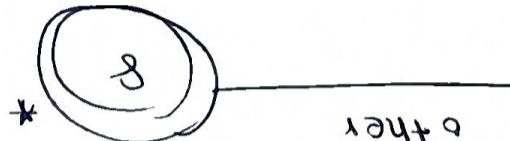
ms+allZB()

return (defToNone(),



Example : Transition diagram for received words and

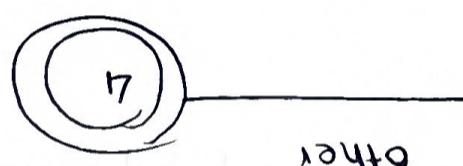
return (relop, GT)



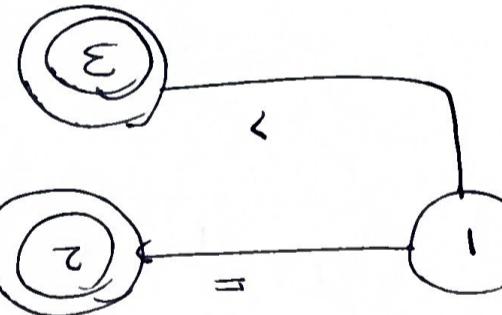
return (relop, GE)



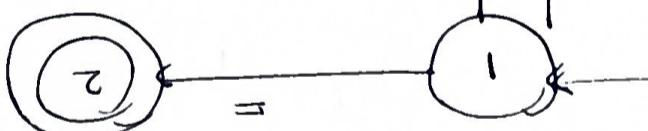
return (relop, LT)



return (relop, NE)



return (relop, LE)



Example : Draw a transition diagram for relational operators

Step 1 Draw a transition diagram

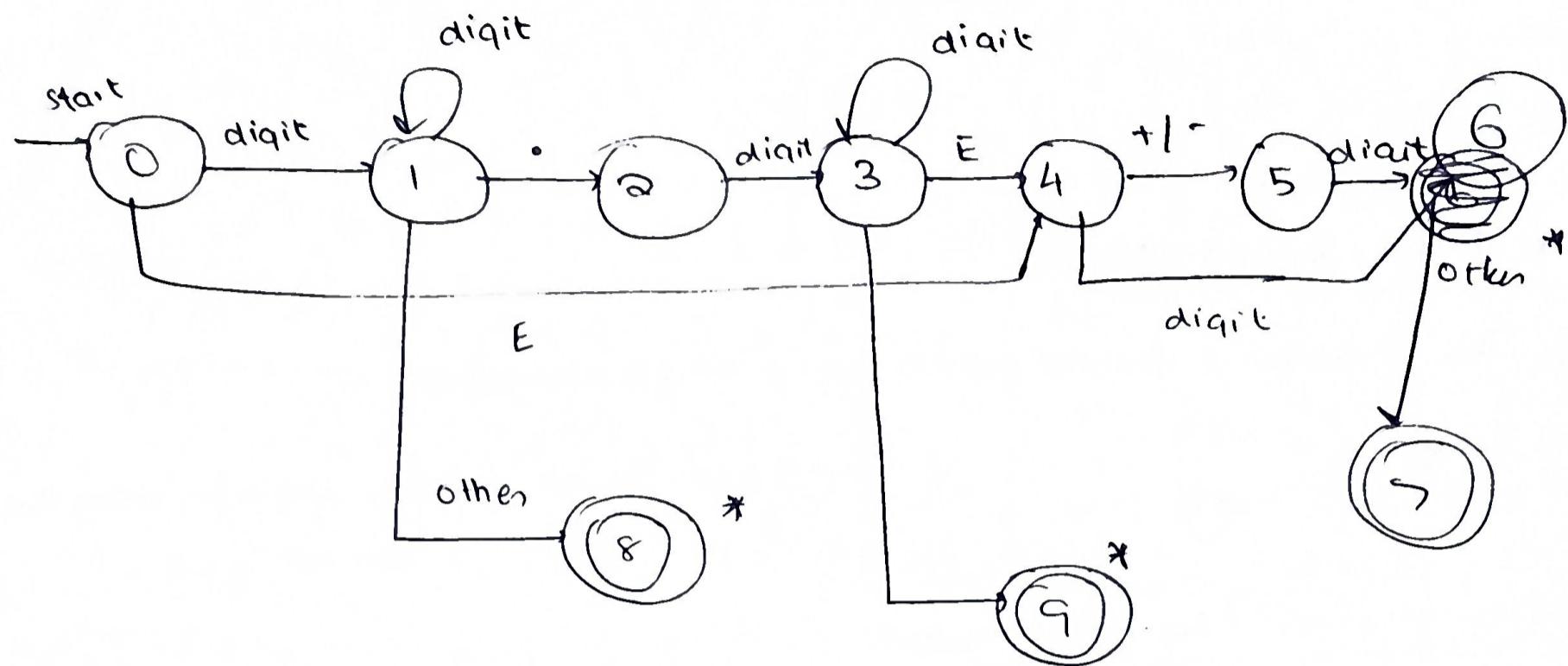
ws  $\hookrightarrow$  (binary lab) newline ) +

Handle white spaces

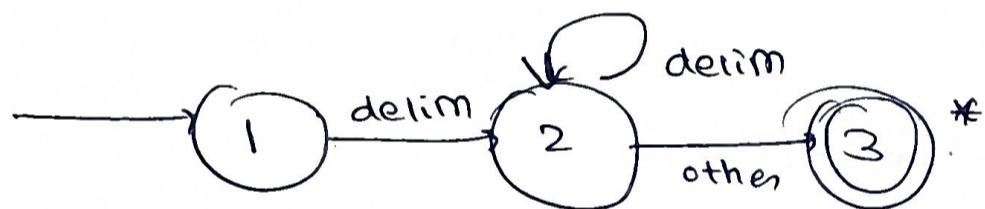
Step 3

(21)

**Example 3** Transition diagram for unsigned numbers

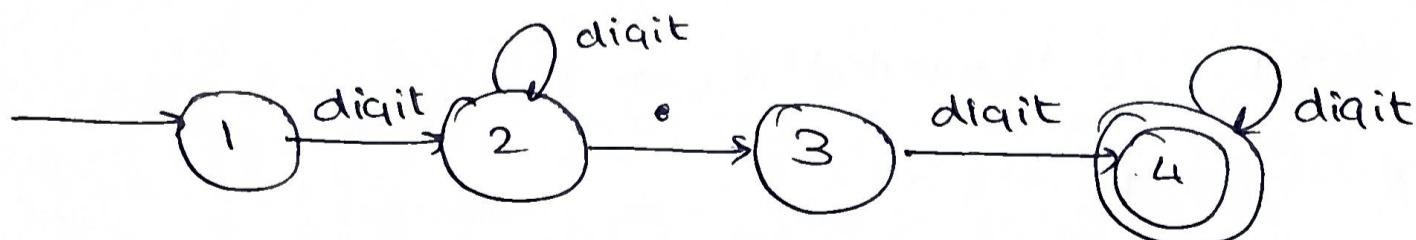


**Example 4** Transition diagram for whitespaces

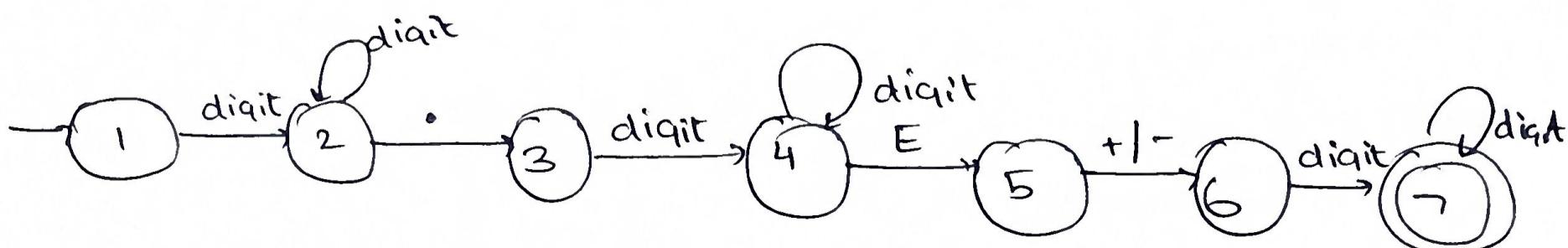


**Example 5** Draw the transition diagram for the following

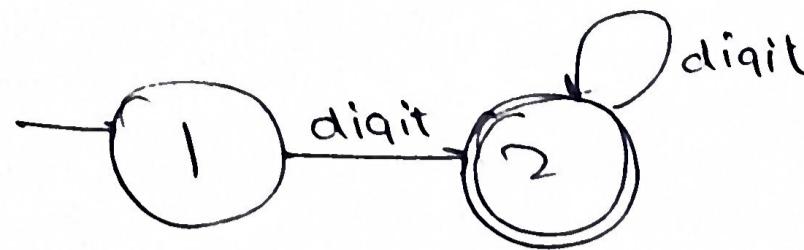
(i) ~~3.45E.3.4.56~~



(ii) 3.45E-5



(ii) 3456



### \* Finite Automata

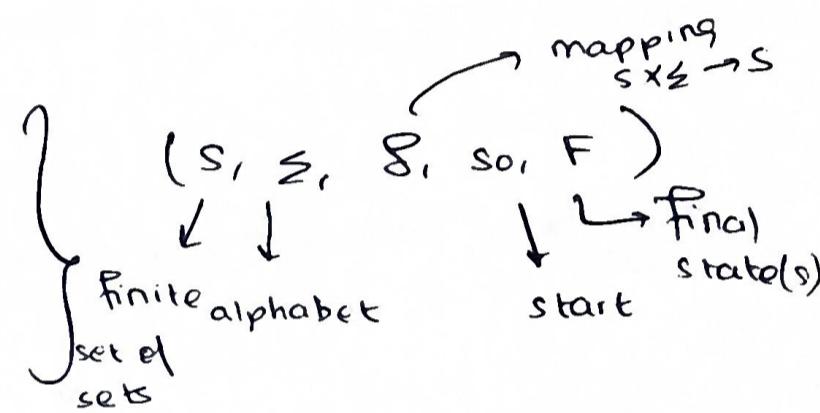
→ RE = outline, Finite automata = implementation

→ 2 types → NFA  
→ DFA

↳ have finite memory, encode only the current state

#### Deterministic Finite Automata

- one transition per input per state
- No  $\epsilon$ -moves



#### Non-deterministic finite automata

- can have multiple transitions for one input in a given state
- can have  $\epsilon$ -moves

### \* Simulating a DFA

→ Input string  $x$  is terminated by an EOF, for a DFA  $D$

→ Output yes if accepts, else no

$S := s_0$

$a := \text{nextchar}()$

while  $a \neq \text{eof}$  do

$s := \text{move}(s, a)$

$a := \text{nextchar}()$

end do

If  $S$  is in  $F$ , then

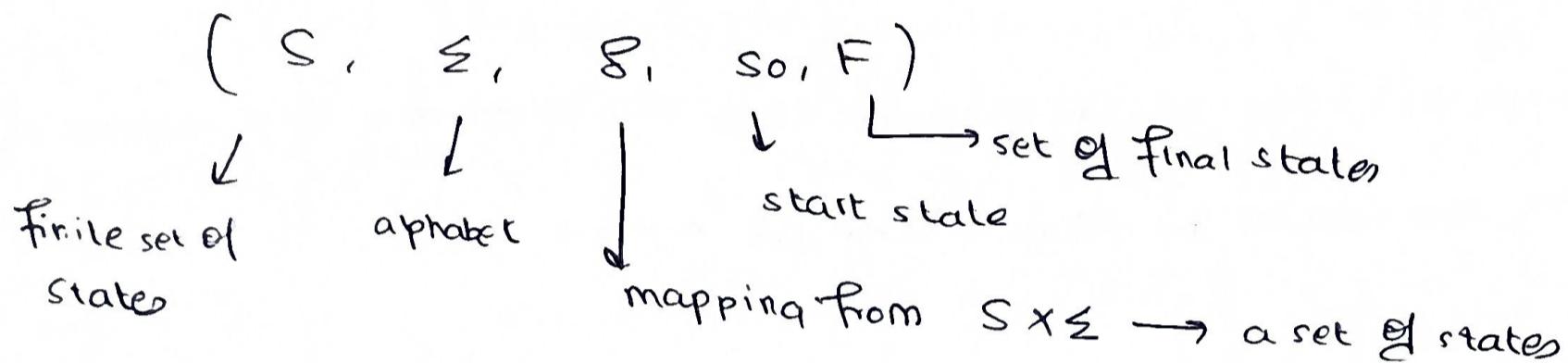
return 'yes'

else

return 'no'

### \* Nondeterministic Finite Automata

An NFA is a 5-tuple



### \* Direct Conversion from RE to DFA

Step 1 : Augment the regular expression  $r$  with a special end symbol  $\#$

Step 2 : Construct a syntax tree for  $(r)\#$

alphabets = leaf node

operator = inner node

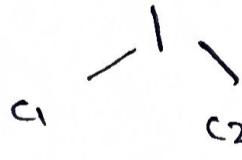
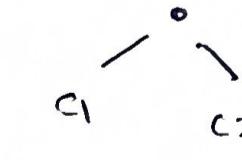
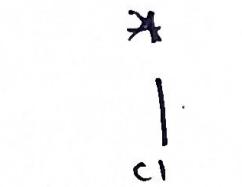
Step 3 : Number each alphabet including  $\#$

Step 4 : Traverse the tree to construct the functions  
 nullable, firstpos, lastpos & followpos

## Terms

- (i)  $\text{nullable}(n)$  : the subtree at node  $n$  generates languages including the empty string
- (ii)  $\text{firstpos}(n)$  : set of positions that can match the first symbol of a string generated by a subtree at node  $n$
- (iii)  $\text{lastpos}(n)$  : the set of positions that can match the last symbol of a string generated by the subtree at node  $n$
- (iv)  $\text{followpos}(i)$  : the set of positions that can follow position  $i$  in the tree

### \* Annotating the tree

| Node $n$                                                                            | $\text{nullable}(n)$                                    | $\text{firstpos}(n)$                                                                                                              | $\text{lastpos}(n)$                                                                                                             |
|-------------------------------------------------------------------------------------|---------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------|
| leaf $\epsilon$                                                                     | true                                                    | $\emptyset$                                                                                                                       | $\emptyset$                                                                                                                     |
| leaf $i$                                                                            | false                                                   | $\{i\}$                                                                                                                           | $\{i\}$                                                                                                                         |
|  | $\text{nullable}(c_1)$<br>or<br>$\text{nullable}(c_2)$  | $\text{firstpos}(c_1) \cup \text{firstpos}(c_2)$                                                                                  | $\text{lastpos}(c_1) \cup \text{lastpos}(c_2)$                                                                                  |
|  | $\text{nullable}(c_1)$<br>and<br>$\text{nullable}(c_2)$ | $\text{if } \text{nullable}(c_1) \text{ then } \text{firstpos}(c_1) \cup \text{firstpos}(c_2) \text{ else } \text{firstpos}(c_1)$ | $\text{if } \text{nullable}(c_2) \text{ then } \text{lastpos}(c_1) \cup \text{lastpos}(c_2) \text{ else } \text{last pos}(c_2)$ |
|  | true                                                    | $\text{firstpos}(c_1)$                                                                                                            | $\text{lastpos}(c_1)$                                                                                                           |

## followpos(n)

for each node n in the tree do

if n is a cat-node with the left child c<sub>1</sub> and right child c<sub>2</sub> then

for each i in Lastpos(c<sub>1</sub>) do

$$\text{followpos}(i) := \text{followpos}(i) \cup \text{Firstpos}(c_2)$$

end do

else if n is a star-node

for each i in Lastpos(n) do

$$\text{followpos}(i) := \text{followpos}(i) \cup \text{Firstpos}(n)$$

end do

end if

end do

→ Calculate D<sup>Trans</sup> for each unique state starting from start  
 $(A1q0 \ pq \ Q1)$

A drawing showing two separate circles, each containing a black 'X' mark, representing female gametes.

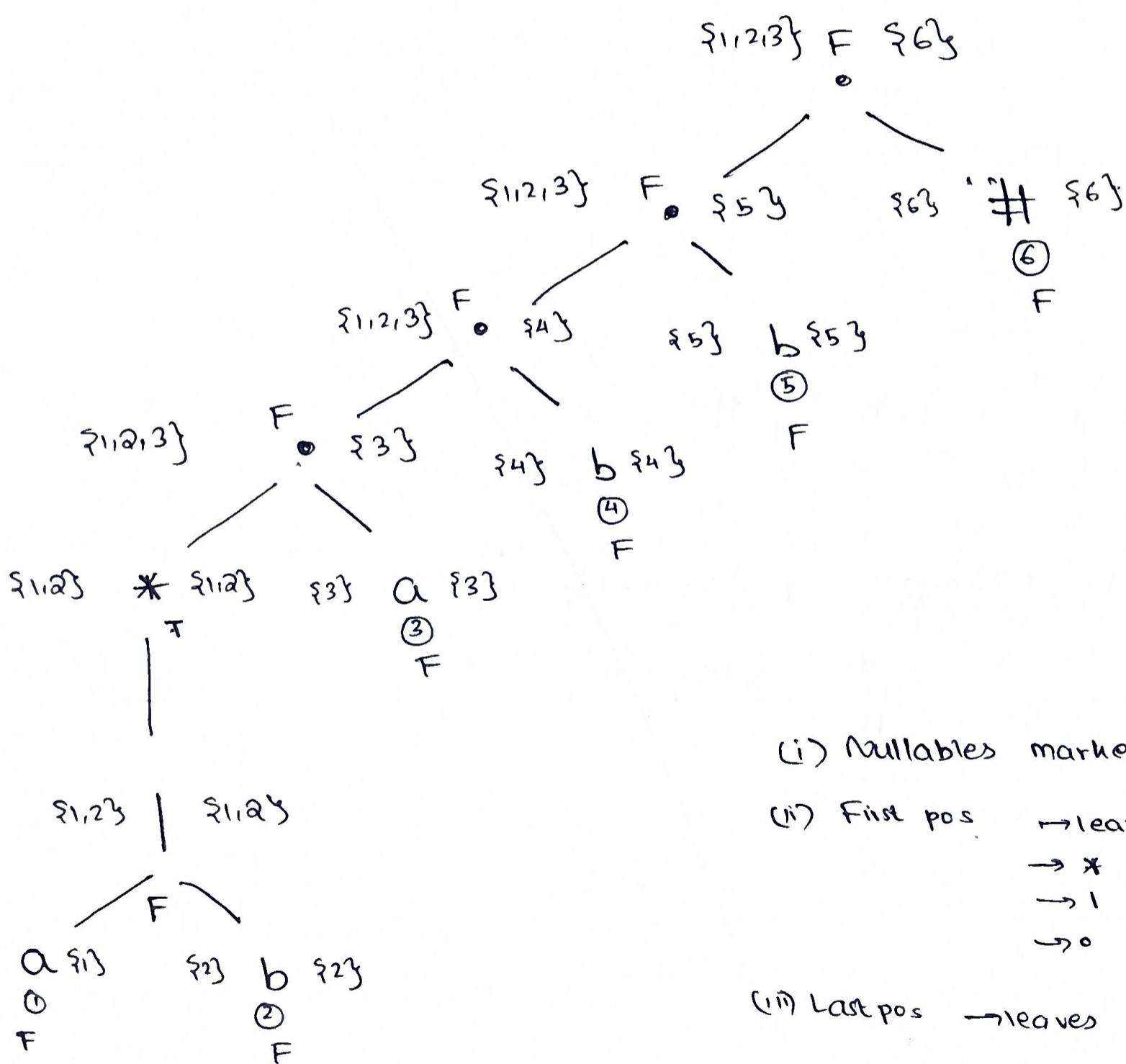
Convert the following Regular expressions into a DFA

23

$$\textcircled{1} \quad (a1b)^* abb$$

$$\Rightarrow (a|b)^* \cdot a \cdot b \cdot b \cdot \#$$

Ans



(i) Nullables marked

(ii) First pos  $\Rightarrow$  leaves

→ \*

(ii) Last pos  $\rightarrow$  always

→ \*

## Calculating Followpos

(i)

|   | Node  | Followpos                     |
|---|-------|-------------------------------|
| a | 1 (•) | <del>1, 2, 3</del><br>1, 2, 3 |
| b | 2 (•) | <del>1, 2, 3</del><br>1, 2, 3 |
| a | 3 (•) | 4                             |
| b | 4 (•) | 5                             |
| b | 5 (•) | 6                             |
| # | 6     | -                             |

$$\text{cat-node: } \text{followpos}(i) := \text{followpos}(i) \cup \text{firstpos}(c_2)$$

$$\text{star node: } \text{followpos}(c) :=$$

$$\text{followpos}(i) \cup \text{firstpos}(n)$$

$$(a/b)^* \cdot a \cdot b \cdot b \cdot \#$$

## Calculating DTrans

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

$$DTrans(A, a) = B = \{1, 2, 3, 4\}$$

$$DTrans(A, b) = A$$

$$DTrans(B, a) = B$$

$$DTrans(B, b) = C = \{1, 2, 3, 5\}$$

$$DTrans(C, a) = B$$

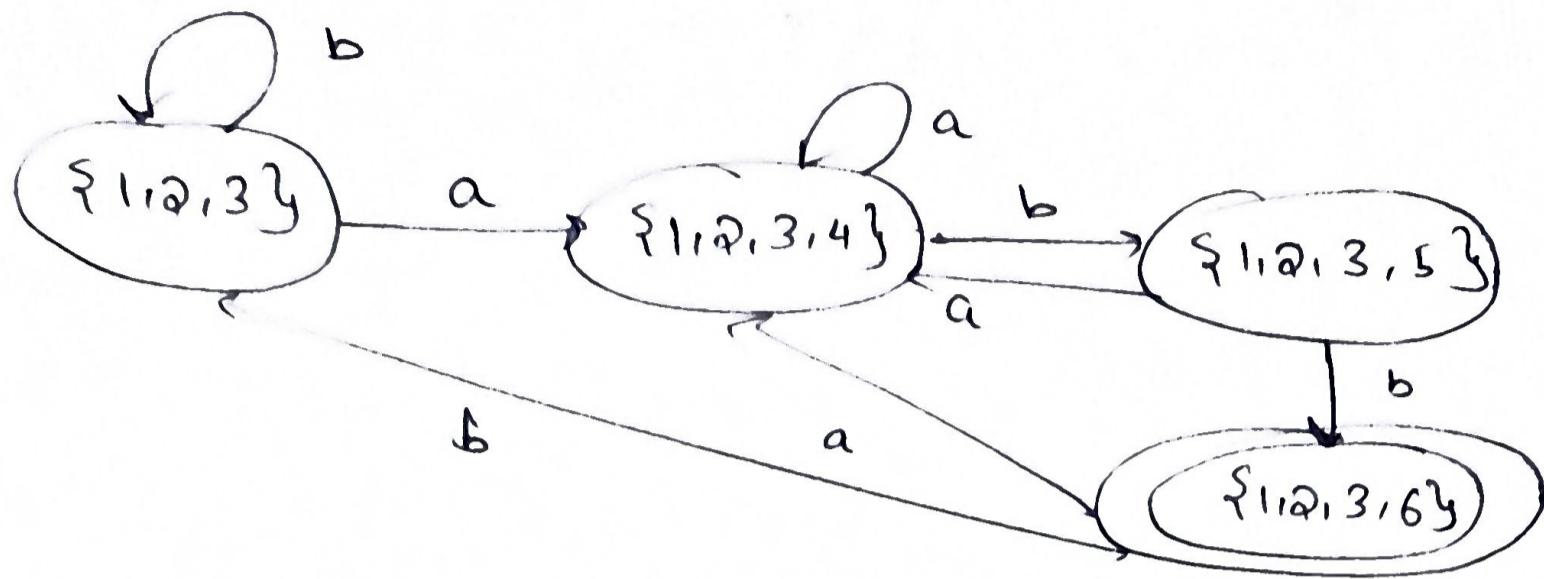
$$DTrans(C, b) = D = \{1, 2, 3, 6\}$$

$$DTrans(D, a) = B$$

$$DTrans(D, b) = D$$

# Drawing a DFA

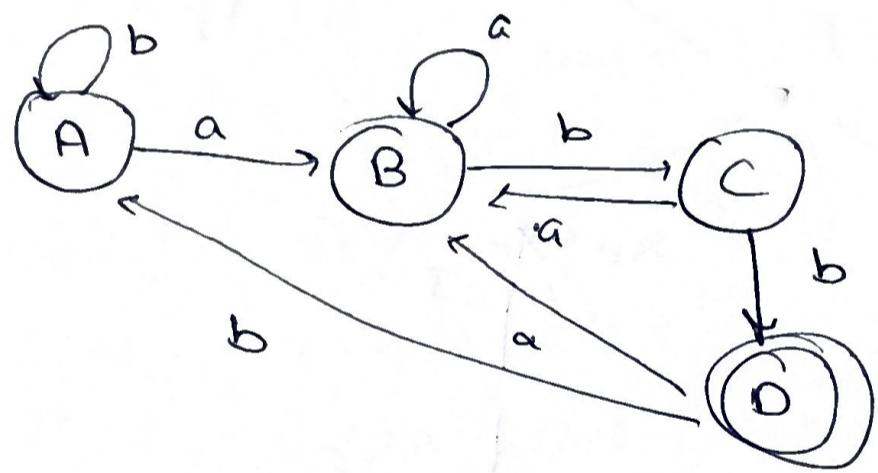
(29)



|   |   |     |
|---|---|-----|
| a | 1 | 123 |
| b | 2 | 123 |
| a | 3 | 4   |

Final DFA

# 6 -

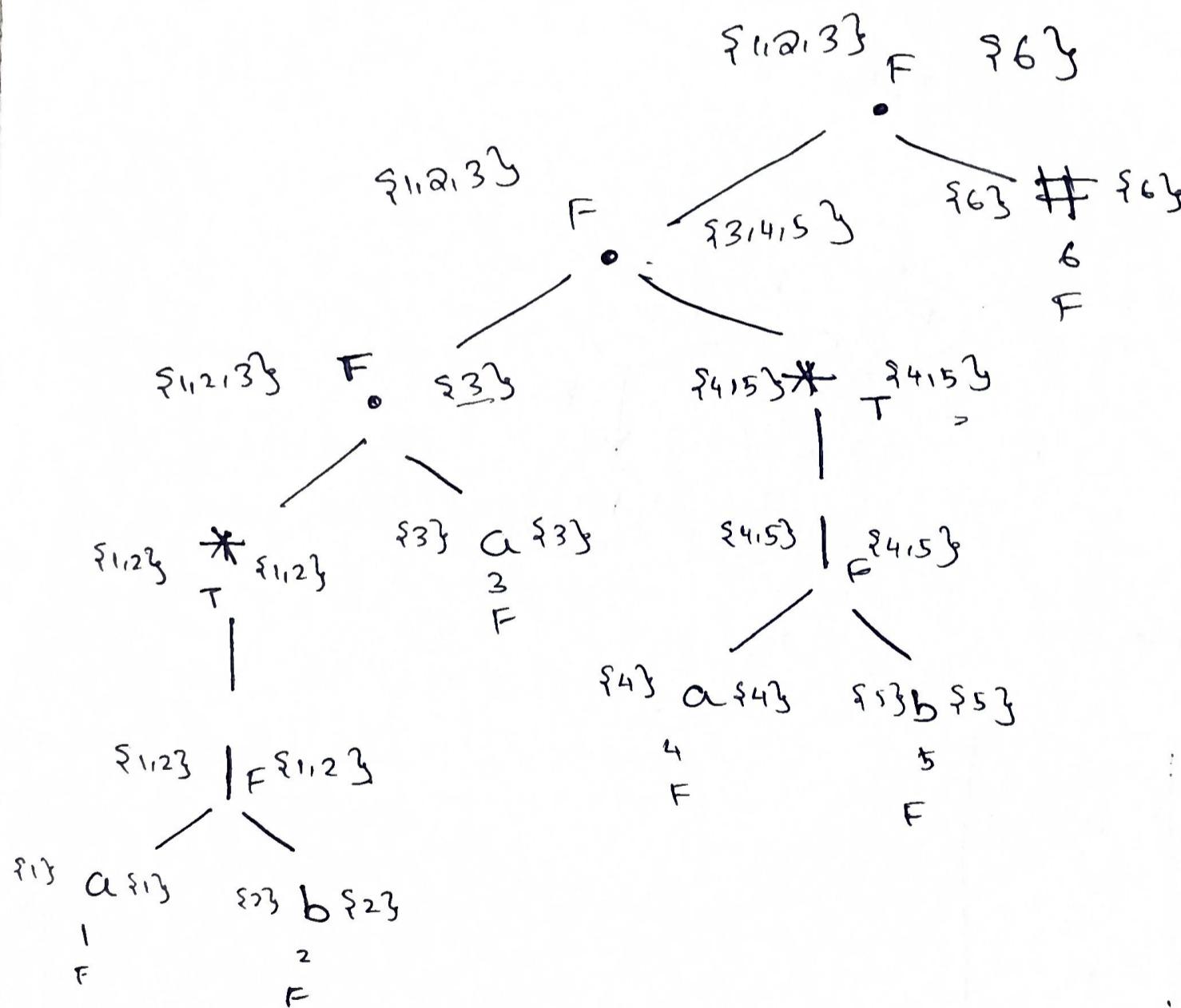


Transition Table

| S | a | b |
|---|---|---|
| A | B | A |
| B | A | C |
| C | B | D |
| D | B | A |

②  $(a/b)^* a (a/b)^*$

$(a/b)^* \cdot a \cdot (a/b)^* \cdot \#$



Calculating Followpos

$$(a|b)^* \underset{1}{a} (a|b)^* \underset{5}{\#} \underset{6}{}$$

| Node name | Node no | Followpos |
|-----------|---------|-----------|
| a         | 1       | 1, 2, 3   |
| b         | 2       | 1, 2, 3   |
| a         | 3       | 4, 5, 6   |
| a         | 4       | 4, 5, 6   |
| b         | 5       | 4, 5, 6   |
| #         | 6       | -         |

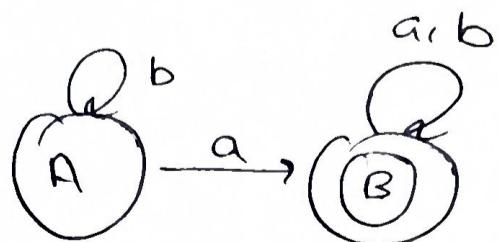
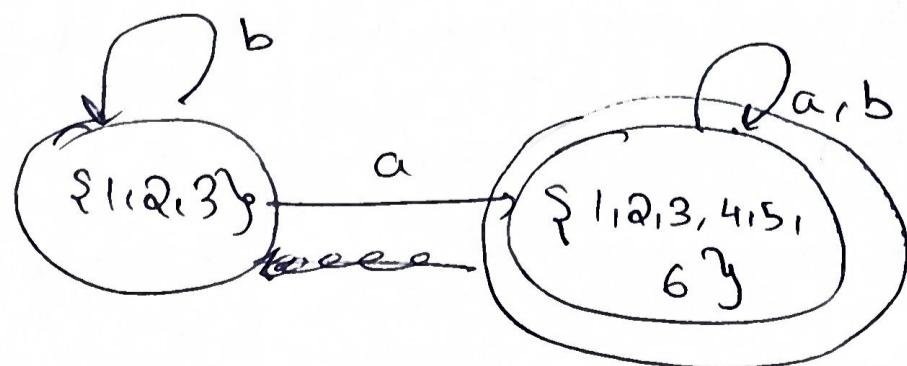
Calculating DTrans

$$DTrans(A, a) = \{1, 2, 3, 4, 5, 6\} = B$$

$$DTrans(A, b) = \{1, 2, 3\} = A$$

$$DTrans(B, a) = \emptyset$$

$$DTrans(B, b) = B$$

Constructing DFA

## \* Lex - Lookahead Operator and Conflict Resolution

### Conflict Resolution

→ Conflict arises when several prefixes of input matches one or more patterns. This can be resolved by the following:

- (i) Always prefer a longer prefix than a shorter prefix
- (ii) If two or more patterns are matched for the longest prefix, then the first pattern listed in the lex program is preferred.

e.g. For patterns:       $a \rightarrow P_1$   
                         $abb \rightarrow P_2$   
                         $a^*b^+ \rightarrow P_3$

For string abb,  $P_2$ ,  $P_3$  match

$P_2$  is chosen as it was listed first in the lex program

### Lookahead Operator

- An additional operator that is read by Lex in order to distinguish additional patterns for a token
- Lexical analyzer is used to read one character ahead of valid lexeme and then retracts to produce a token.
- At times, it is needed to have certain characters at the end of input to match with a pattern.
- In such cases / is used to indicate end of part of pattern that matches the lexeme.

e.g. 1 If language keywords are not reserved :

IF (I,J) = 5    8

IF (condition) THEN

results in a conflict whether to produce IF as an arrayname or a keyword.

This can be resolved as :

IF / \(.\*)\} {keyword}  
IF \(.\*)\} {arr name}

e.g 2 {IDENT} /ab

I/p = bracadabra

o/p = bracad for 44text

e.g 3 {IDENT}

I/p = bracadabra

o/p = bracadabra for 44text

e.g 4 Rule: a+b/c

I/p : aaabcc

o/p = aaab for 44text

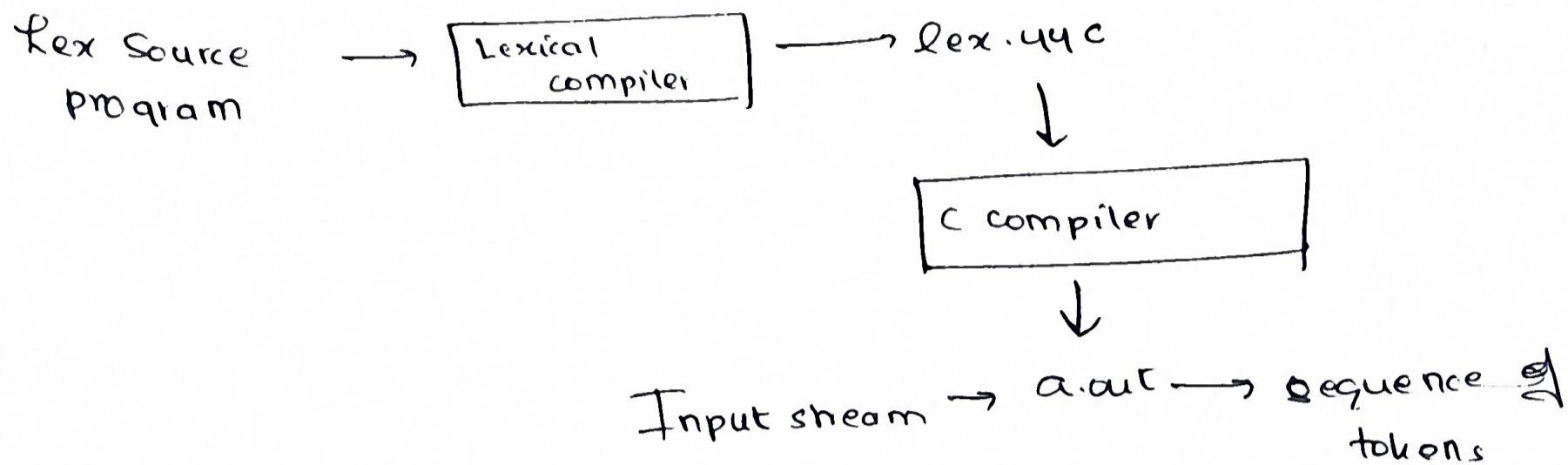
e.g 5 Rule: x\* | xy

I/p = xxxy

o/p = xxxy for 44text

## \* Introduction to LEX

Lexical Analyzer Generator



## \* Lex Specification

Has 3 parts :

- (i) regular definitions, C definitions
- (ii) Translation rules
- (iii) User-defined auxiliary procedures

## \* Sample Lex Specification

```
{0-9} +      { printf ("%s\n", yytext ); }  
· | \n      { }  
main ()  
{ 4ylex(); }  
}
```

compile it by:

```
lex spec.l
gcc lex.yy.c -ll
./a.out spec.l
```

Ex1 Lex program to count no. of lines, words & characters

1. {

```
#include <stdio.h>
int cn = 0, wd = 0, nl = 0;
```

1. }

delim [ \t ] +

1n

```
{ cn++; wd++; nl++ }
```

```
{delim} { cn+=44lenq; wd++; }
```

• { ch++; }

1. 1.

main() {

u4lex();

```
printf("%d %d %d %d\n", nl, wd, ch),
```

}

Example 2 RE LEX Program to recognize digits, letters & identifiers

1. {

```
#include <stdio.h>
```

2. }

digit [0-9]

letter [A-zA-Z]

id {letter} ({letter} | {digit})\*

3. .

```
{digit}* {printf("number: %s\n", yytext); }  
{id} {printf("ident: %s\n", yytext); }  
. {printf("other: %s\n", yytext); }
```

4. .

```
main() {  
    yylex();  
}
```