

**COMPILER FUNCTIONS:**

- \* check errors and does validation
- \* construct a model for conversion
- \* Does memory allocation for intermediate code.
- Programmers can concentrate on effective code to solve a problem.

"COMPILER" is a software that contains several lines of code and uses all components of a computer to generate the machine code.

**TRANSLATORS OF HELLO WORLD PROGRAM**

- \* Assemblers
- \* Interpreter
- \* Compiler

Interpreter

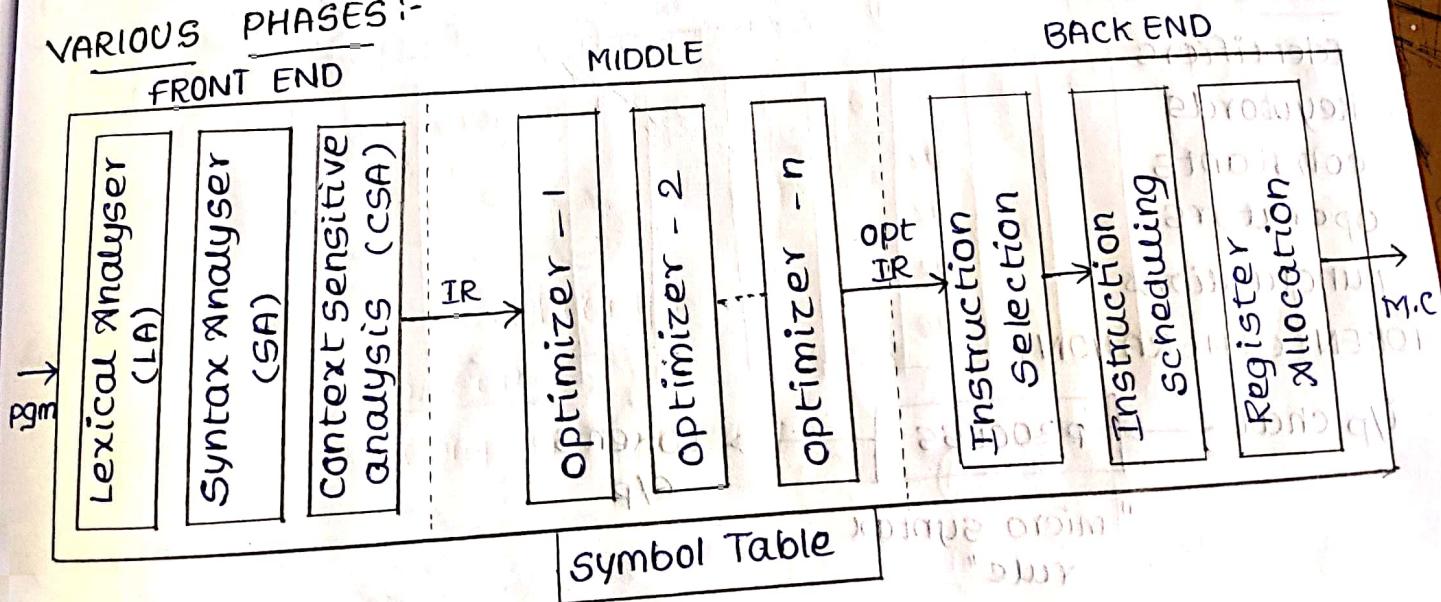
Line by Line  
at the end  
result of the  
program

compiler

full code at  
the end

Machine code  
(below)

Assembly Language → **Assembler** → Machine Code

**VARIOUS PHASES:-**

LA (scanner) → tokens : keywords, identifiers, constants, operators.

SA (parser) → check grammar (syntax) → generate IR  
IR → Intermediate Representation

opt → optimize

CSA → context sensitive analysis

Instruction Selection → selects machine instruction for code.

Instruction Scheduling → schedule machine instruction for code.

RA → Register Allocation for IR

Symbol Table → Repository for variables in code

PRINCIPLES OF A COMPILER :- (2 MA)

\* It shouldn't change the meaning of the code written by programmer

\* It should build an effective machine code.

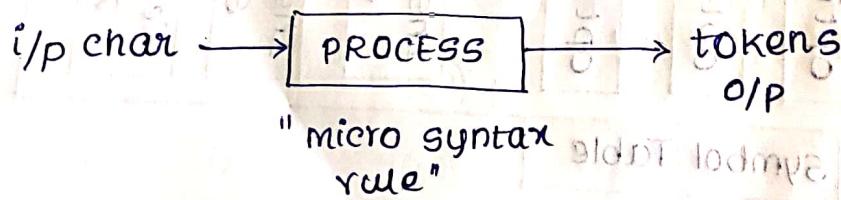
#### LEXICAL ANALYSIS [SCANNER]



#### TOKENS

- Identifiers
- keywords
- constants
- operators
- punctuations

#### TOKEN GENERATION:



"micro syntax rule"

ADVANTAGES

- \* parser generates syntax tree, parse tree, etc.
- \* LA is fast
- \* Every run time
- Scanner can

#### 1. STATE

(S<sub>0</sub>)

coding:

→ nextchar

New Scan

c ← n

GF(c)

begin

c ←

if

be

else

#### STATE

#### 1. NE

(S<sub>0</sub>)

2.0

variables, constants

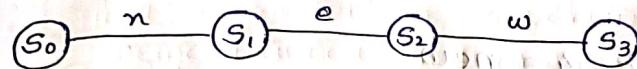
generate IR  
↓  
syntax tree,  
parse tree, TAC

construction  
functions  
for  
scanner  
construction  
for  
parser  
and  
lexer  
written

### ADVANTAGES OF LA OVER PARSER:

- \* parser gets words as i/p and not letters
- \* LA is fully automatic than parser
- \* Every rule added in LA shrinks parser construction.

#### 1. State transition



coding:

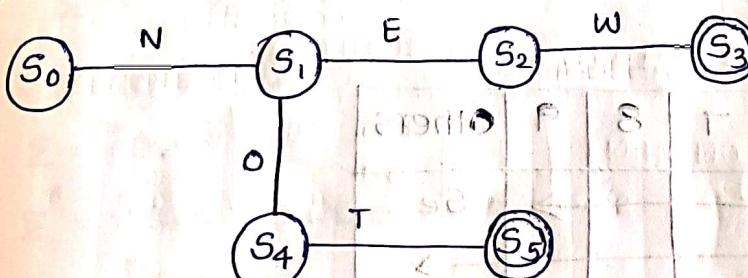
→ nextchar() function to get char i/p. e.g. below scanning

New scanning:

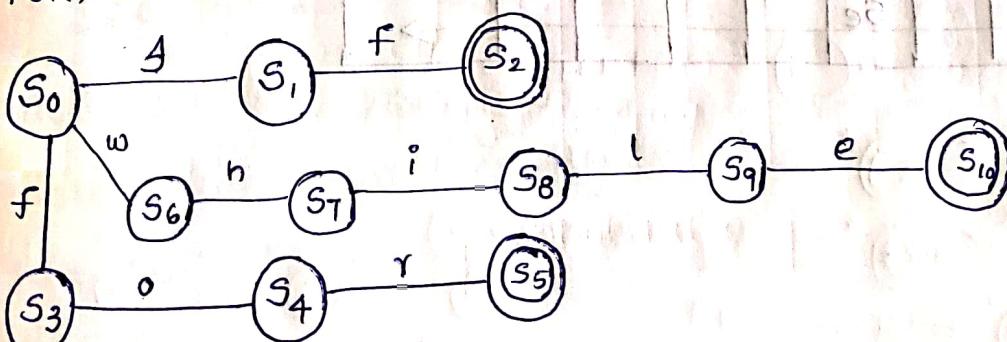
```
c ← nextchar();  
if (c = 'n') then  
begin  
    c ← nextchar();  
    if (c = 'e') then  
        begin  
            c ← nextchar();  
            if (c = 'w') then  
                report 'success'  
            else  
                do something else
```

#### STATE TRANSITION FOR:

##### 1. NEW, NOT



##### 2. FOR, WHILE, IF



## IMPLEMENTING SCANNER:

1. Table driver
2. Hand coded

FINITE AUTOMATA (FA):

\* FA is a collection of states, set of transition between states, alphabets, a start state and one or more final states.

\* 5 Tuples:  $\{S, \Sigma, S_0, S_f\}$  new

$S \rightarrow \{S_0, S_1, S_2, S_3\}$  states

$\Sigma \rightarrow \{n, e, w\}$  alphabets

$S_f \rightarrow \{(S_0 \xrightarrow{n} S_1), (S_1 \xrightarrow{e} S_2), (S_2 \xrightarrow{w} S_3)\}$  transition

$S_0 \rightarrow S_0$

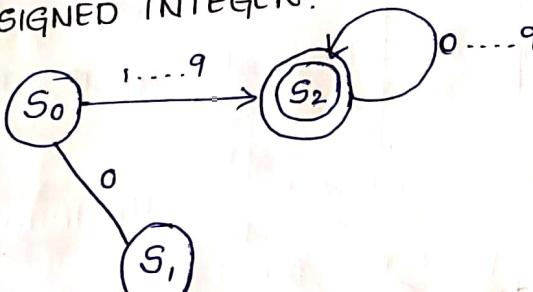
$S_f \rightarrow S_3$

Write FA for new, not

$S \rightarrow \{(S_0 \xrightarrow{n} S_1), (S_1 \xrightarrow{e} S_2), (S_2 \xrightarrow{w} S_3), (S_1 \xrightarrow{o} S_4), (S_4 \xrightarrow{t} S_5)\}$

$S_f \rightarrow \{S_5, S_3\}$

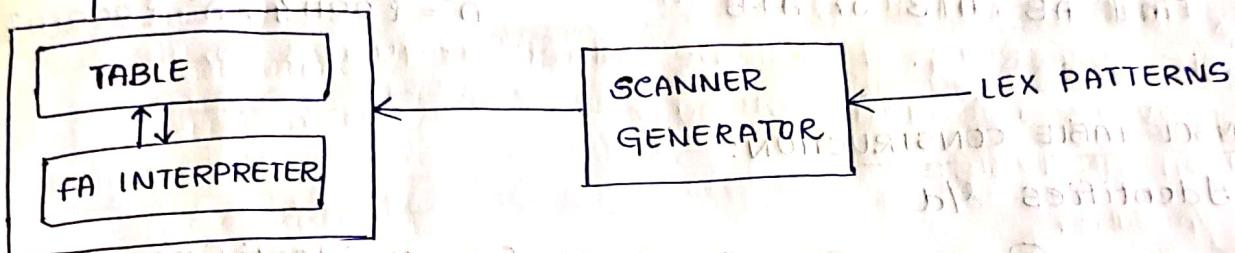
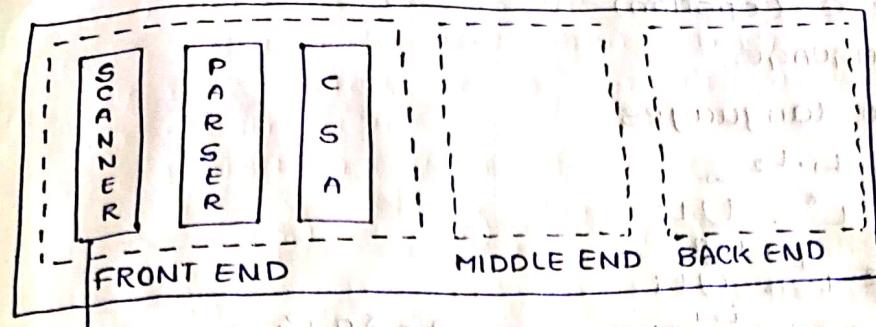
## UNSIGNED INTEGER.



## TABLE DRIVER::

	0	1	2	3	4	5	6	7	8	9	Others
S0	$S_1 \leftarrow$				$S_2$				$S_e \rightarrow$		
S1		$\leftarrow$			$S_e$				$\rightarrow$		
S2			$\leftarrow$		$S_2$			$\rightarrow$			
S3				$\leftarrow$	$S_e$			$\rightarrow$			

## AUTOMATIVE SCANNER GENERATION (Table)



**DIRECT CODED SCANNER:**

USING RE:- Regular Expression

- pattern matches a RE then token is recognized.

FEW NOTATION:

1. Alphabet - Set of symbols -  $\Sigma$

$$\Sigma = \{0, 1\}$$

$$\Sigma = \{0, 1, 2, \dots, 9\}$$

$$\Sigma = \{A - Z\}$$

$$\Sigma = \{a - z, A - Z, 0 - 9\}$$

2. String -  $\Sigma^*$  - Sequence of numbers from alphabets.

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

$$\Sigma^* = \{a, b, ab, aab, abb, \dots\}$$

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

$$\Sigma^* = \{ab, bc, ac, aa, bb, cc, \dots\}$$

3. Length of the string  $|w|$

→ Total number of strings

$$w = abc$$

$$|w| = 3$$

S.T

OKEN

4. Empty word - without any symbol  
 $w = \emptyset$ , or  $\epsilon$  (epsilon)

# OPERATION ON LANGUAGE: LANGUAGES

### OPERATION

#### 1. UNION OF two languages.

2. CONCATENATION

2. CONCATENATION  
3. KLEEN CLOSURE,  $L^* = \bigcup_{i=0}^{\infty} L^i$

4. positive closure  $L^+ = \bigcup_{i=1}^{\infty} L_i$

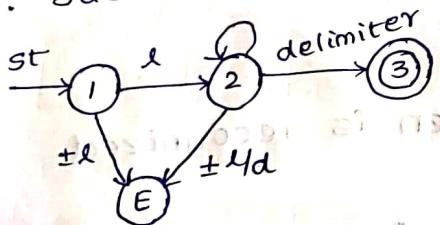
4. POSITIVE

EG:  $A = \{ab, bc\}$      $B = \{21, 43\}$   
 $A/B$  OR  $A+B$

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
1																										
2																										
3																										
4																										
5																										
6																										
7																										
8																										
9																										
10																										
11																										
12																										
13																										
14																										
15																										
16																										
17																										
18																										
19																										
20																										
21																										
22																										
23																										
24																										
25																										
26																										
27																										
28																										
29																										
30																										
31																										
32																										
33																										
34																										
35																										
36																										
37																										
38																										
39																										
40																										
41																										
42																										
43																										
44																										
45																										
46																										
47																										
48																										
49																										
50																										
51																										
52																										
53																										
54																										
55																										
56																										
57																										
58																										
59																										
60																										
61																										
62																										
63																										
64																										
65																										
66																										
67																										
68																										
69																										
70																										
71																										
72																										
73																										
74																										
75																										
76																										
77																										
78																										
79																										
80																										
81																										
82																										
83																										
84																										
85																										
86																										
87																										
88																										
89																										
90																										
91																										
92																										
93																										
94																										
95																										
96																										
97																										
98																										
99																										
100																										
101																										
102																										
103																										
104					</td																					

## SYMBOL TABLE CONSTRUCTION

1. Identifies id



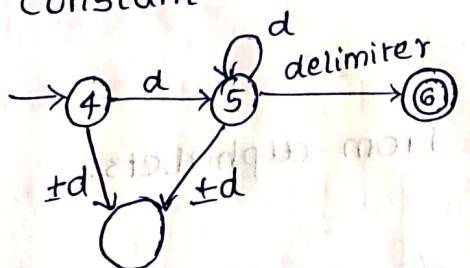
\* repeat for all identifiers.

```
return C(i, install(i))
```

↓  
in symbolTable

TOKEN	code	value
identifier	1	ptr. to st

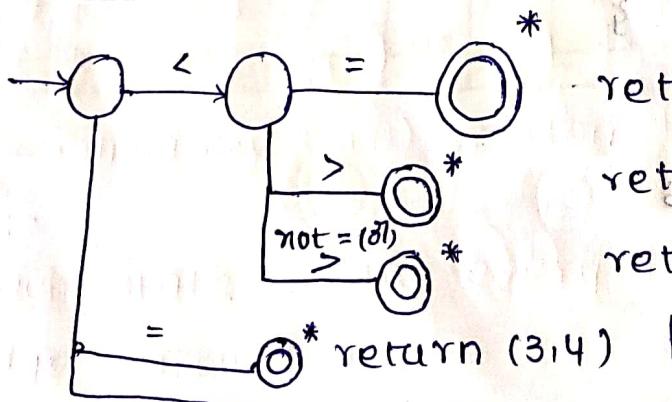
## 2. Constant



```
return(2,install())
```

CONSTANT	2	ptr to S.T
----------	---	---------------

## RELATIONAL OPERATIONS:



return (3,1)

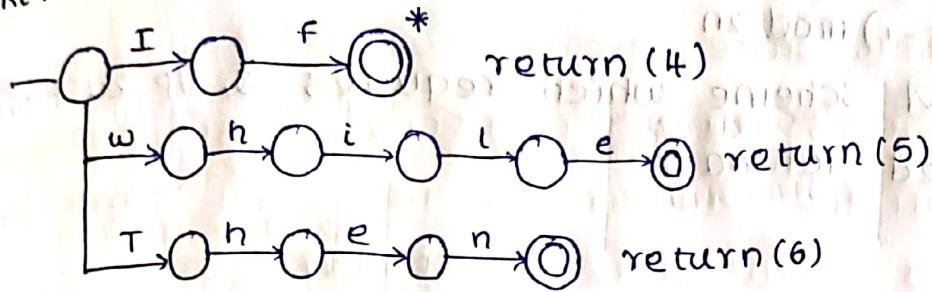
return (3,2);

return (3,3)

 \*  
return (3,5)  
return (3,6)

S.T	TOKEN	CODE	VALUE
	<=	3	1
	<>	3	2
	<	3	3
	=	3	4
	>	3	5
	>=	3	6

KEYWORDS:



TOKEN	CODE	V
IF	4	-
WHILE	5	-
THEN	6	-

WRITE! STRINGS:

1.  $a|ba^*$  : {a, b, aa, ab, baaa, ---}
2.  $(a|b)(a|b)(a|b)$  : {aaa, bbb, abb, bbb, bab, ---}
3.  $\epsilon|a|b$  : { $\epsilon$ , a, b}
4.  $(a|b)(a|b)(a|b)^*$  : {aa, bb, aaa, ---}
5.  $\bigcup_{i=0}^{\infty} (a,b)^i$  (Kleene closure) = { $\epsilon$ , aa, bb, ---}

HAND CODED SCANNER:  
 Instead of table driven and direct coded scanners, calling the procedures overhead is reduced by using buffered  $\gamma_0$ . i.e., Each read operation returns a string of characters or buffers and the scanner indexes through the buffer. i.e., the scanner maintains a pointer into the buffer and the next character function keeps the buffer filled and traps the current location in the buffer.

IMPLEMENTING NEXTCHAR & ROLLBACK:

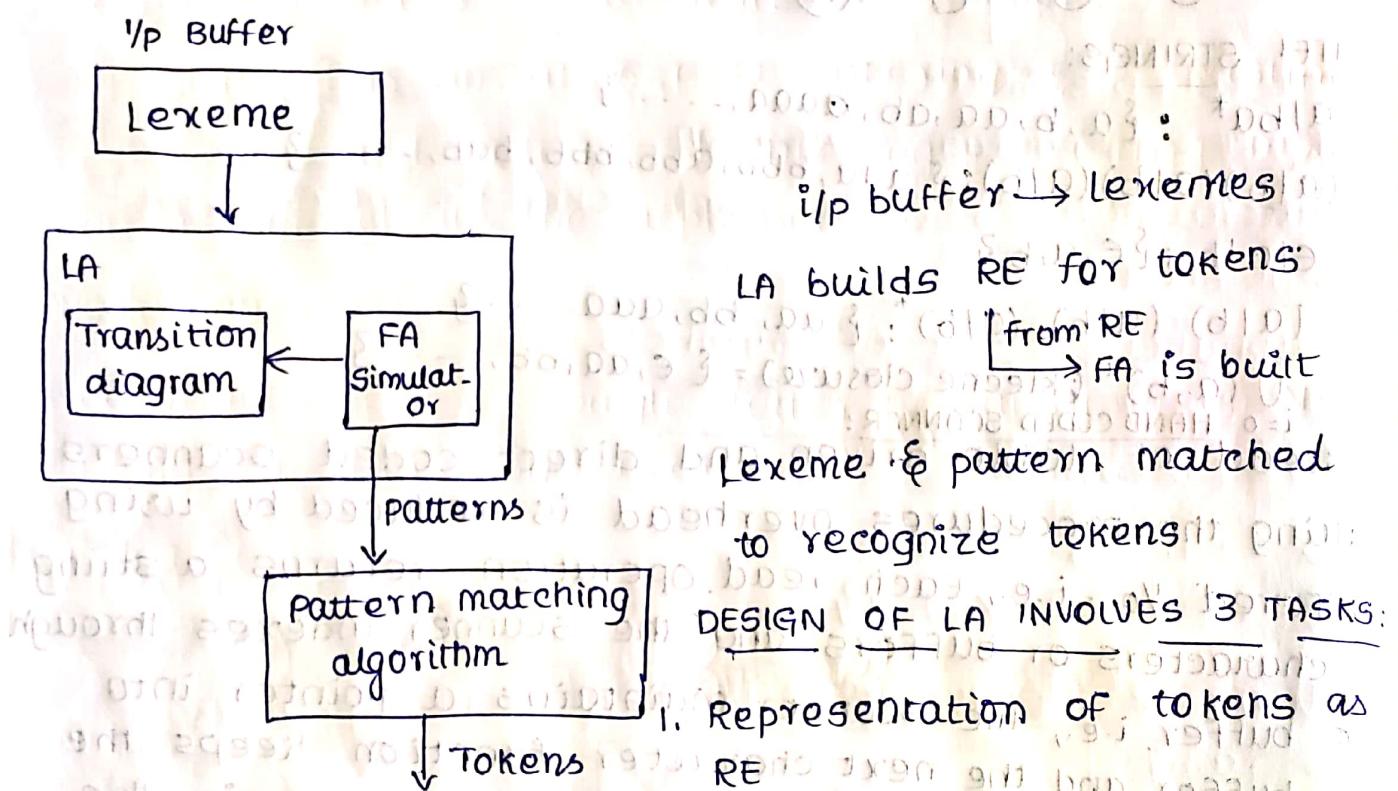
```

char ← Buffer[input];
input ← (input + 1) mod 2n
IF (input mod n == 0) then
begin
    fillbuffer[input : input + n - 1];
    Fence ← (input + n) mod 2n

```

```

    end
    return char;
INITIALIZATION:
    input ← 0
    fence ← 0 ;
    fillbuffer [0:n];
IMPLEMENT ROLL BACK:
    IF (input = fence) then
        signal rollback error,
        input ← (input - 1) mod 2n
imp DOUBLE BUFFERING: A scheme which requires 2 i/p buffers
    to provide bounded roll back.
LEXICAL ANALYSIS:-
```



- DESIGN OF LA INVOLVES 3 TASKS:
1. Representation of tokens as REs
  2. representation of RE by FA or FSM
  3. Simulation of FA

FA:-

- \* Deterministic (DFA) or Non Deterministic ( $\text{NFA}$ ).
- DFA:-
- \* Each state  $s$ , for every input  $a \in \Sigma$  at most one edge (labelled  $a$ )

### FEATURES:-

1. One start state
2. One or more acceptors state
3. No null ( $\epsilon$ ) Transition
4. One transition on a particular symbol.

Eg:  $a^*$



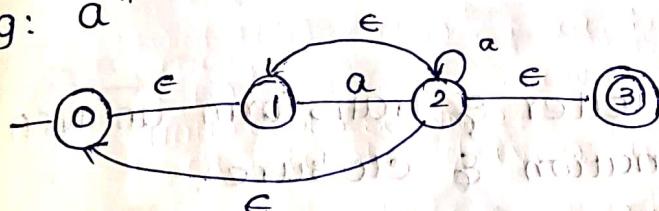
### NFA:-

- \* Start state, few or more Final state and for single input their can be multiple transition.

### FEATURES:-

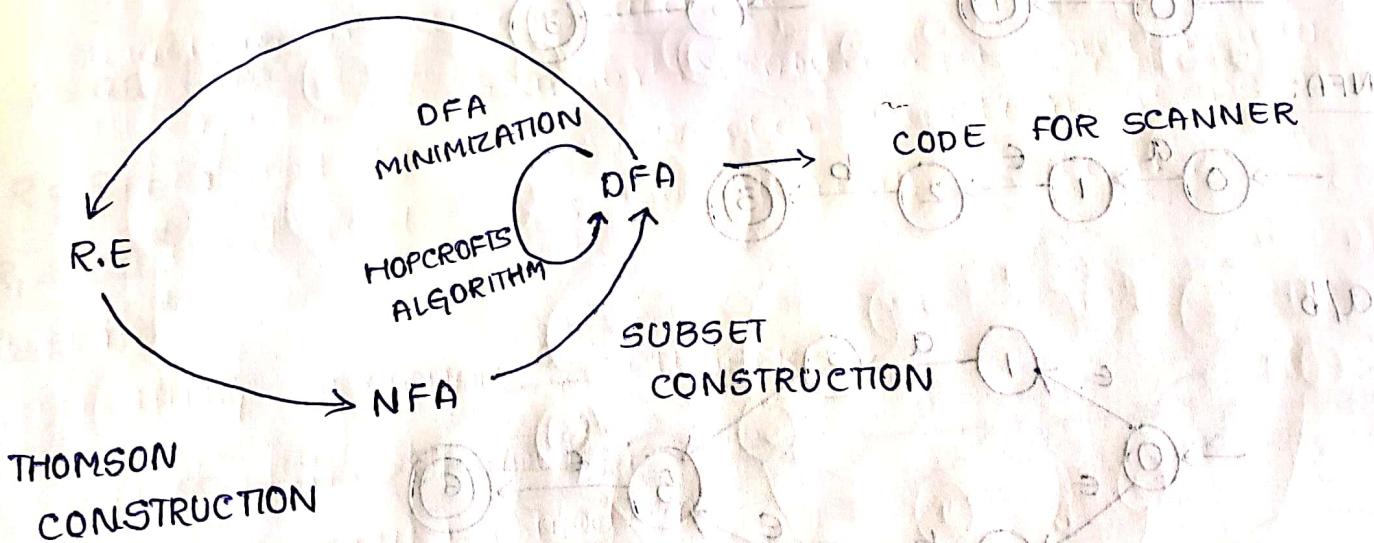
1. One start state
2. One or more accepting state
3. can have ' $\epsilon$ ' transition
4. one or more transition for a input symbol.

Eg:  $a^*$



### CYCLE OF CONSTRUCTION FINITE AUTOMATA :-

#### KLEENE'S CONSTRUCTION



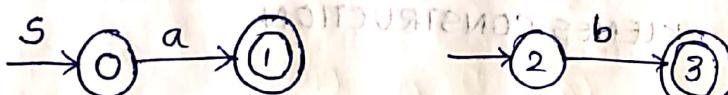
R.E to NFA [Thomson]:

- \* Thomson construction moves to connect NFA's for expressions.
- \* NFA derived from this construction which has several useful properties.
  1. Each NFA has a single start state and a single final state. The only transition that enters the start state is the initial transition. It is the one to leave the final state.
  2. An epsilon ( $\epsilon$ ) move always connects 2 states that were earlier in the process i.e., the start state of a process and a final state of another.
  3. A state has atmost 2 entering & 2 existing ' $\epsilon$ ' moves and atleast 1 entering & existing move on a symbol in the alphabet.

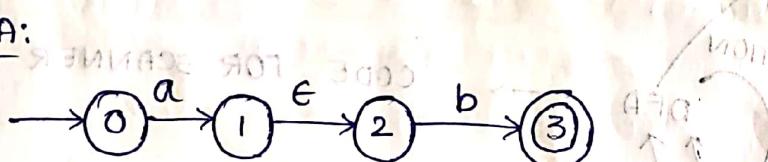
### CONSTRUCTION:-

- \* Construct NFA for Single letter & transform an NFA using concatenation, alternation & closure.

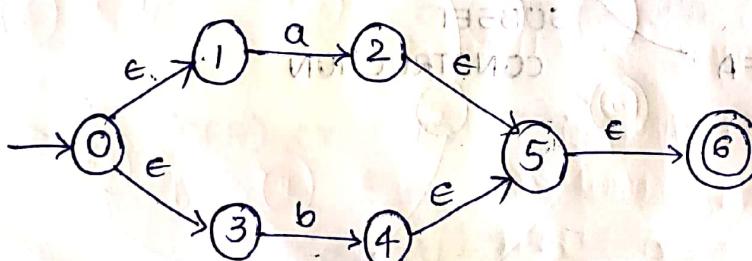
Eg:- 1. a.b



NFA:



2. a/b



1. construct NFA to  $a(b|c)^*$

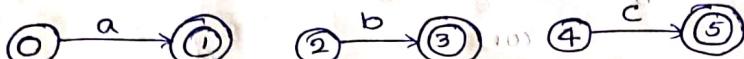
160

$$a \in R_2$$

$$R = \frac{R_1}{\theta}$$

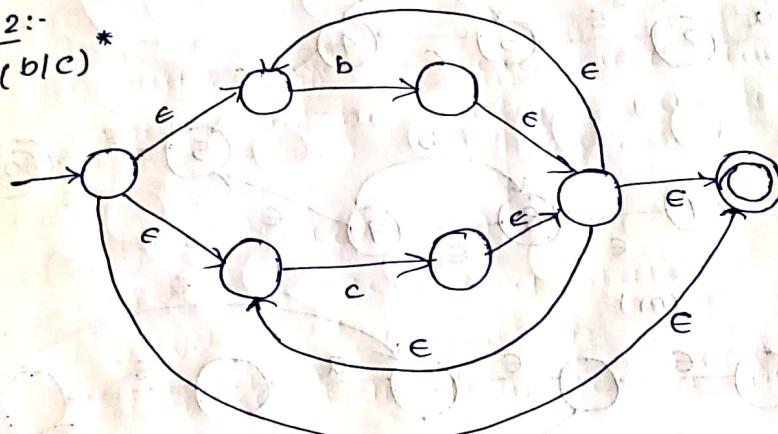
$$\frac{R_1}{R_2} = \frac{a}{(b/c)} *$$

ST:1:



ST:2:-

$$\frac{1}{(b/c)^*}$$



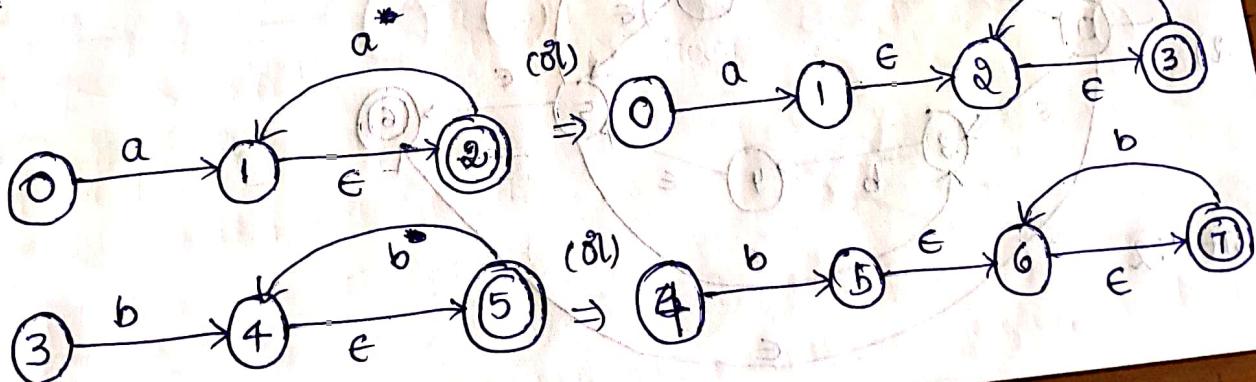
$$2. R = aa^* \mid bb^*$$

$$R = R_1 \mid R_2$$

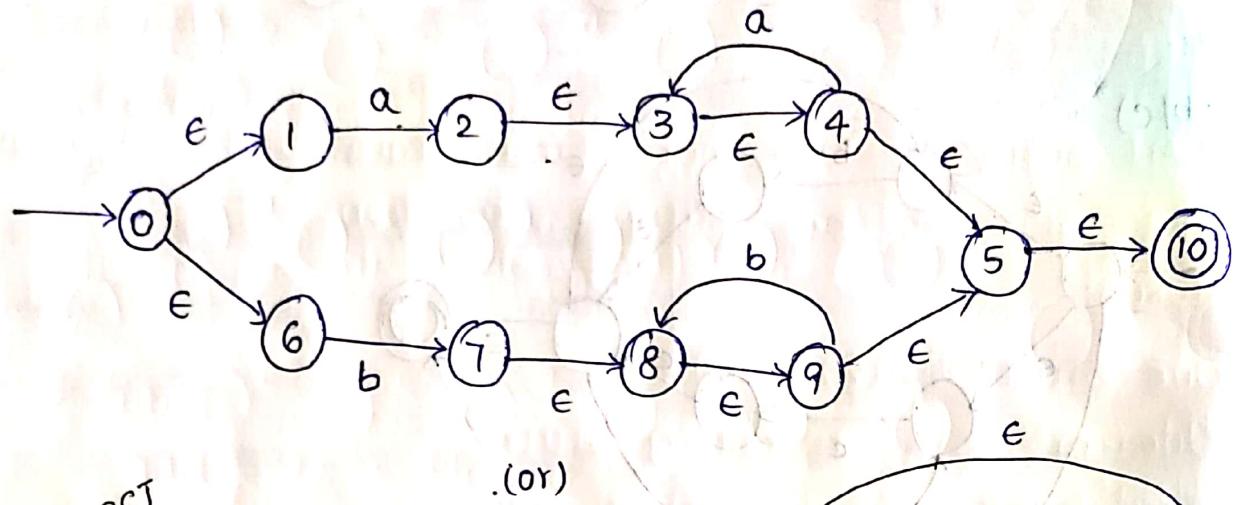
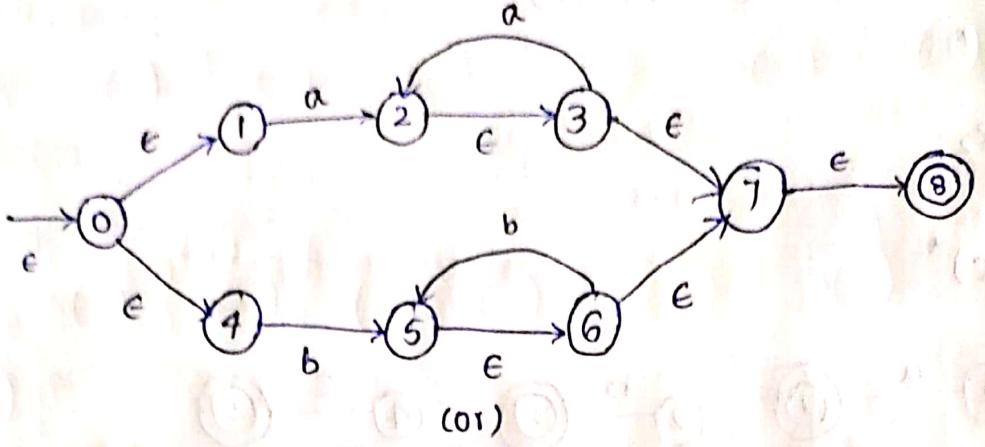
$$R_1 = \alpha a^*$$

$$R_2 = bb^*$$

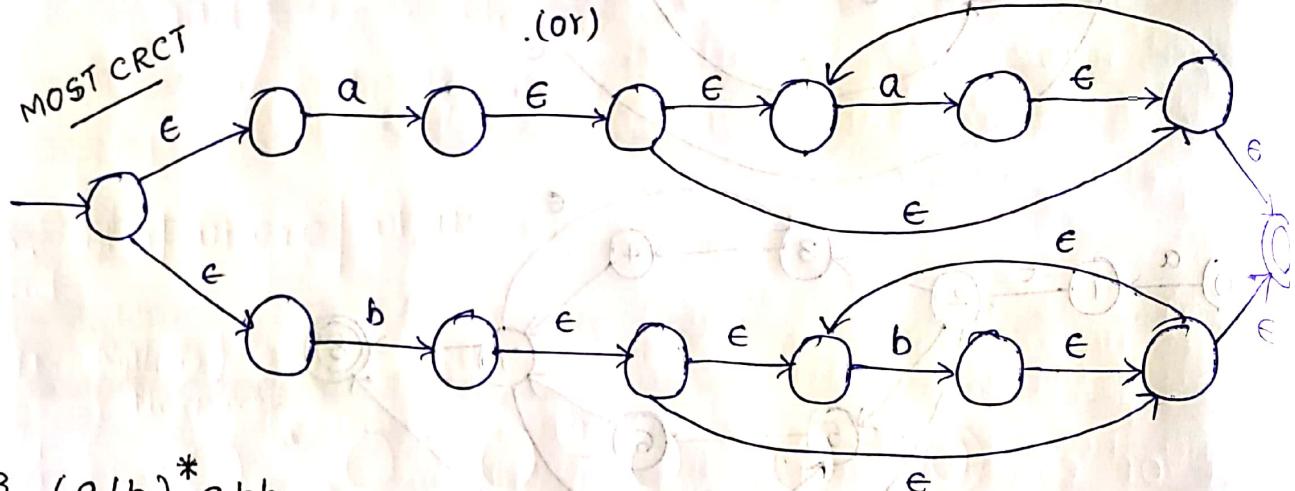
ST:1



ST:2



MOST CRCT



$$3. (a/b)^*abb$$

$$R = R_1 \cdot R_2 \cdot R_3 \cdot R_4$$

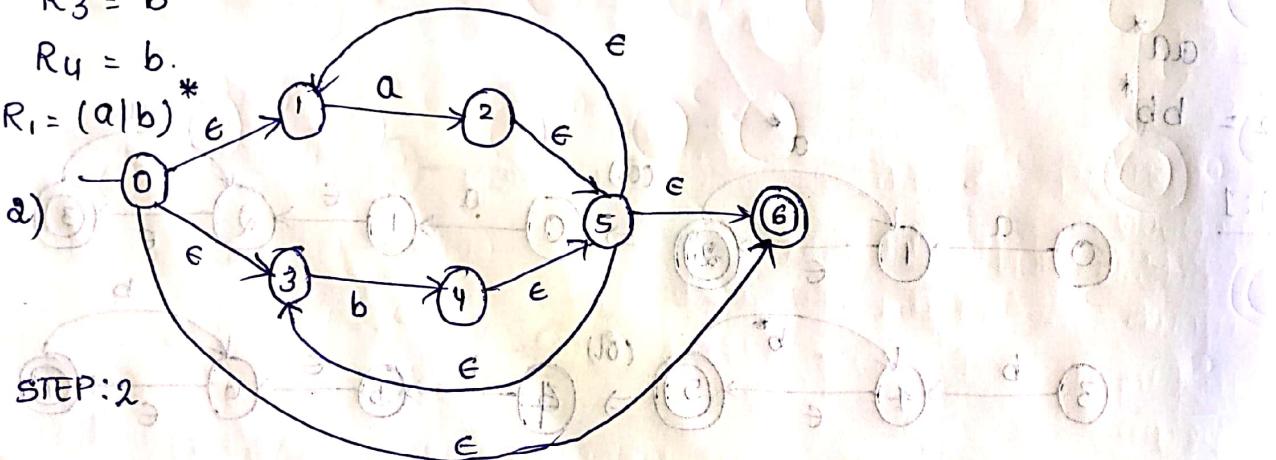
$$R_1 = (a/b)$$

$R_2 = \alpha$

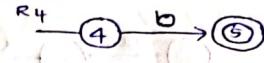
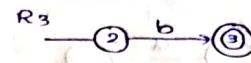
$$R_3 = b$$

$$Ry = b$$

$$R_1 = (a|b)^*$$



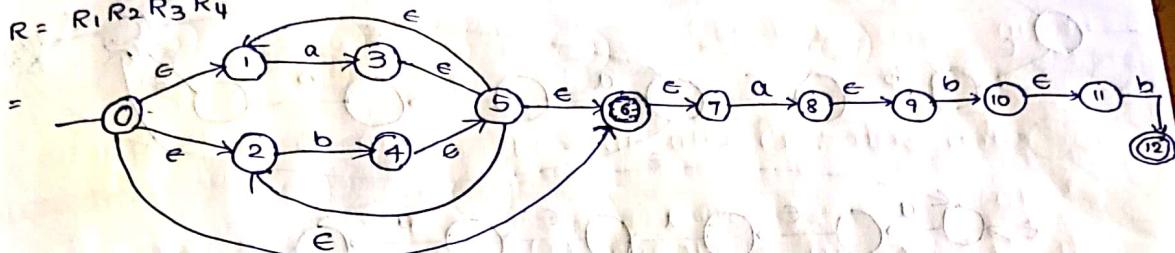
STEP: 1



$$R = R_1 R_2 R_3 R_4$$

$$R_2 R_3 R_4 =$$

$$R = R_1 R_2 R_3 R_4$$



$$4. \quad 1(0+1)^* 0$$

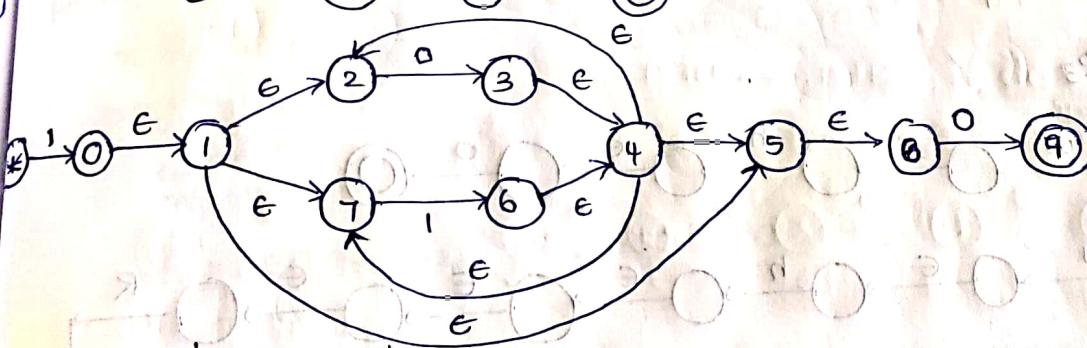
$$= 1(0|1)^* 0$$

$$R = R_1 R_2 R_3$$

$$R_1 =$$

$$R_3 =$$

$$R_2 =$$



$$5. \quad (aaa)^+ \cup b(ab)^*$$

$$R = R_1 R_2 R_3$$

$$R_1 =$$

$$\quad \quad \quad$$

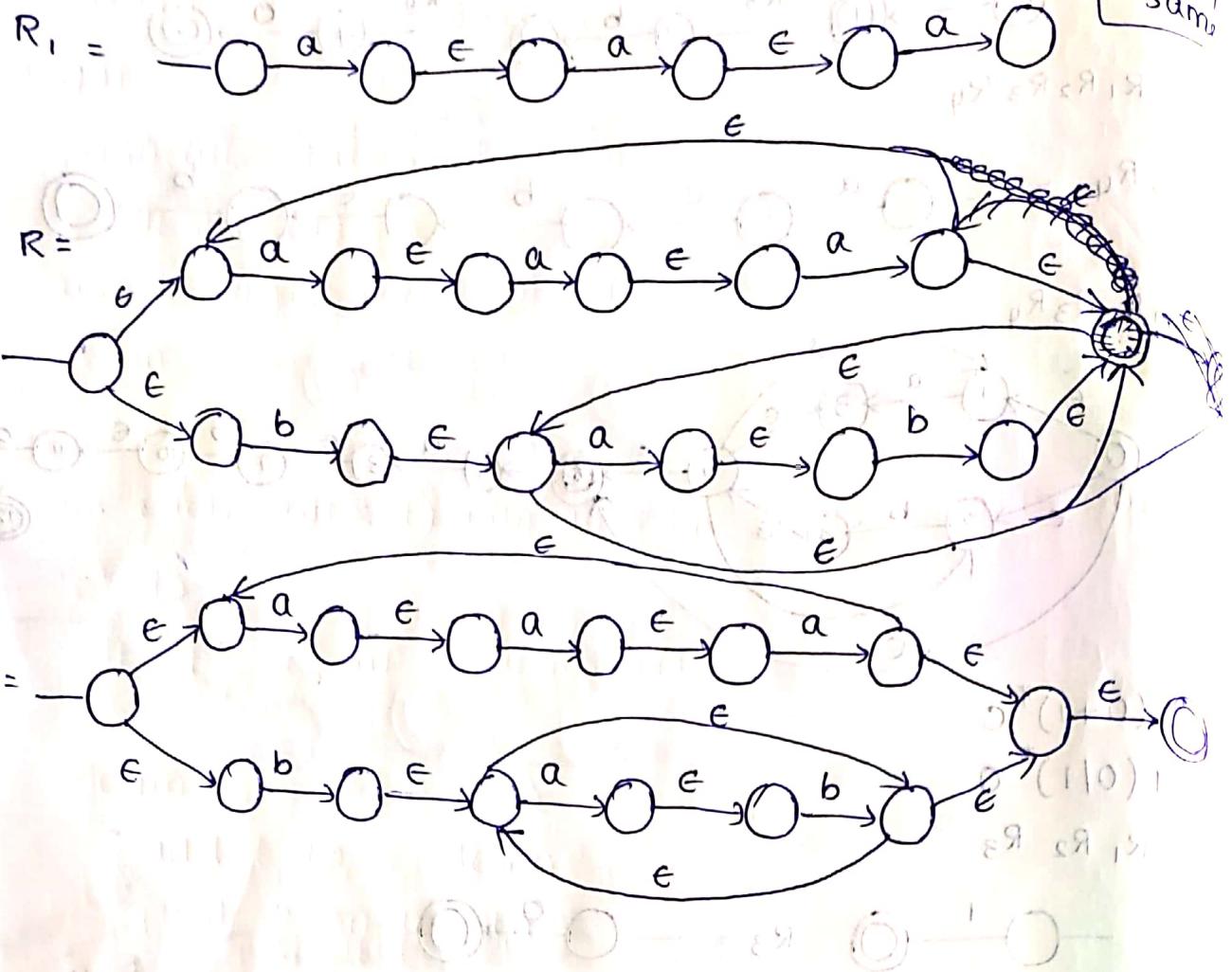
$$R_2 =$$

$$\quad \quad \quad$$

$$R_3 = (ab)^*$$

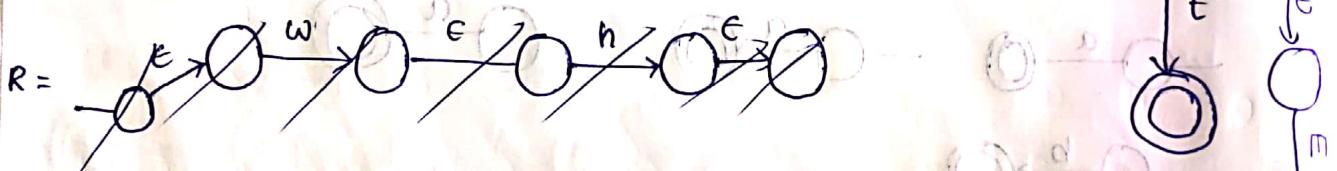
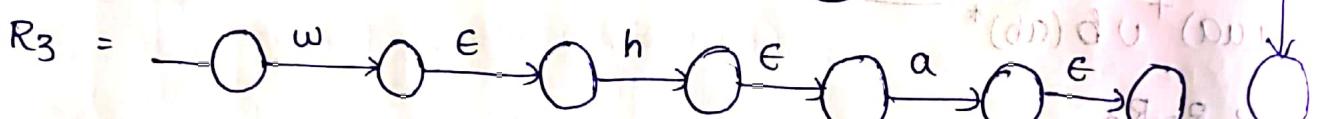
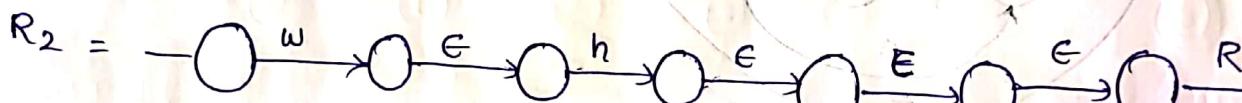
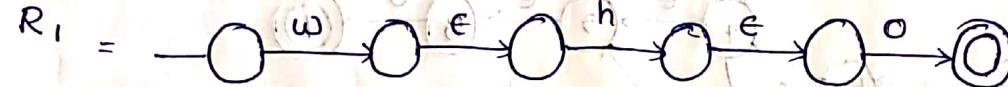
$$=$$

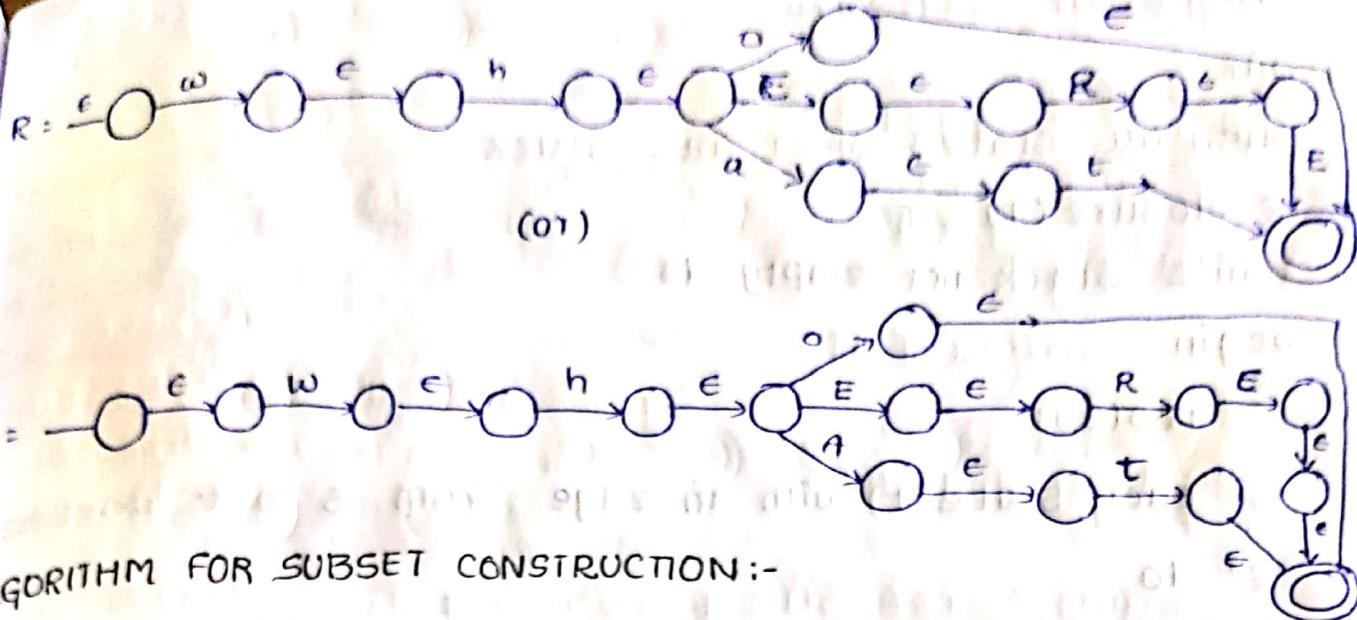
$$\quad \quad \quad$$



6. who | where | what

$$R_0 = R_1 \mid R_2 \mid R_3$$





## ALGORITHM FOR SUBSET CONSTRUCTION:-

Input : NFA ( $N$ )

Input :  
output : DFA (D)

Accepting the same language with the initial state of '0' is the set of  $\epsilon$ -closure ( $S_0$ )

The other states are constructed as follows  
we assume each state of 'N' as unwanted and  
perform the following.

$x = \{s_1, s_2, \dots, s_n\}$  OF D do

begin

1507

mark x

for each i/p symbol 'a' do

begin

Let ' $t$ ' be the set of symbol to which there is a transition 'a'. For some state  $S_1$ , in  $x$ .

$y = \text{e\_closure}(T);$

If 'y' is not added to the set of states of b  
then make "unmarked state of b"  
add a transition for x to y.

end

end

Algorithm for  $\epsilon$ -closure:

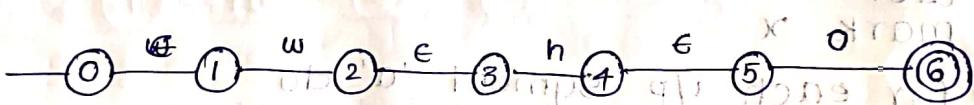
```
Begin
    push all states of T into stack
     $\epsilon$ -closure(T) = T
    while stack not empty do
        begin
            pop(s);
            for each state(t) with an edge from s to t labelled
            'e' do
                if 't' not in  $\epsilon$ -closure(T)
                    do
                        begin
                            add t to  $\epsilon$ -closure
                            push t on to stack
                        end
                    end
                end
            end
        end
    end
```

Eg:-  $(a/b)^*abb$

NFA  $\rightarrow$  DFA

Next page

Eg:



STEP:1:  $\epsilon$  -

$\{0,1\} = A$  (initial states)

STEP:2:

$A \rightarrow w \rightarrow \epsilon$

$\{2,3\} = B$

STEP:3:

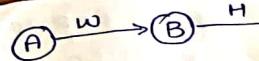
$B \rightarrow h \rightarrow \epsilon$

$\{4,5\} = C$

STEP:4:

$C \rightarrow x \rightarrow \epsilon$

$\{6\} = D$  (final state)



Write a program  
regular expres

# include < std

# include < str

int q[30][3]

int i,j,k,x,op

void print\_NFA();

void main()

{ char reg

int a,b,i

for(i=0

{

for(j=

{

q

y

printf(

scanf(

len =

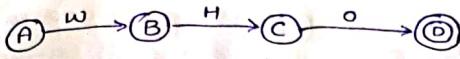
i=j=0

while

{

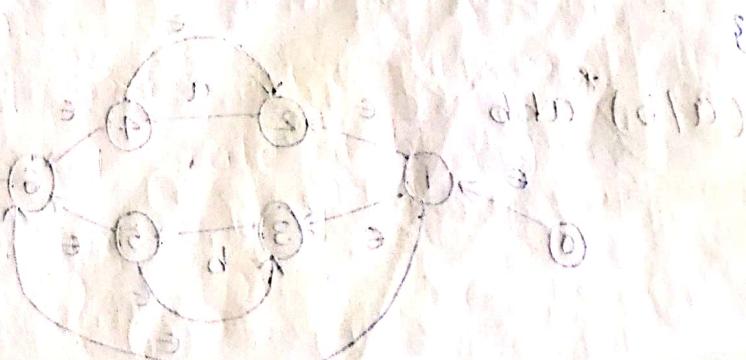
if

{



Write a program for thomson construction for converting a regular expression  $(a/b)^*$  to NFA.

```
#include <stdio.h>
#include <string.h>
int q[30][3];
int i, j, k, x, open;
void print_NFA();
void main()
{
    char reg[20]; // reading input string
    int a, b, len;
    for(i=0; i<30; i++)
    {
        for(j=0; j<3; j++)
            q[i][j] = -1; // initialising
    }
    printf("ENTER REGULAR EXPRESSION:");
    scanf("%s", reg);
    len = strlen(reg);
    i = j = 0;
    while(i < len)
    {
        if (reg[i] == 'a' && reg[i+1] == 'a' && reg[i+2] == 'b')
        {
            q[j][2] = ((j+1)*100)+(j+3);
            j++;
            q[j][0] = j+1;
            q[j][2] = j+3;
            j++;
        }
    }
    print_NFA();
}
```



$q[j][1] = j+1$

$j++;$

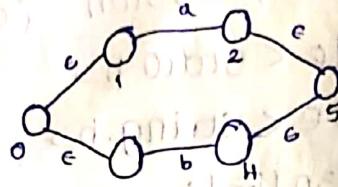
$q[j][2] = j+1;$

$j++;$

$i = i + 2;$

3  
print\_NFA();

o/p:



ENTER THE REGULAR EXPRESSION:

TRANSITION FUNCTION

$q[0, e] \rightarrow 1 \& 3$

$q[1, a] \rightarrow 2 \quad q[2, e] \rightarrow 5$

$q[3, b] \rightarrow 4$

$q[4, e] \rightarrow 5$

void print\_NFA()

{

printf(" TRANSITION FUNCTION");

for(i=0; i<=j; i++)

{

if(q[i][0] != -1)

printf("\n In q[%d, a] → %d", i, q[i][0]);

if(q[i][1] != -1)

printf("\n In q[%d, b] → %d", i, q[i][1]);

if(q[i][2] != -1)

{

if(q[i][2]<100)

printf("\n In q[%d, e] → %d", i, q[i][2]);

else

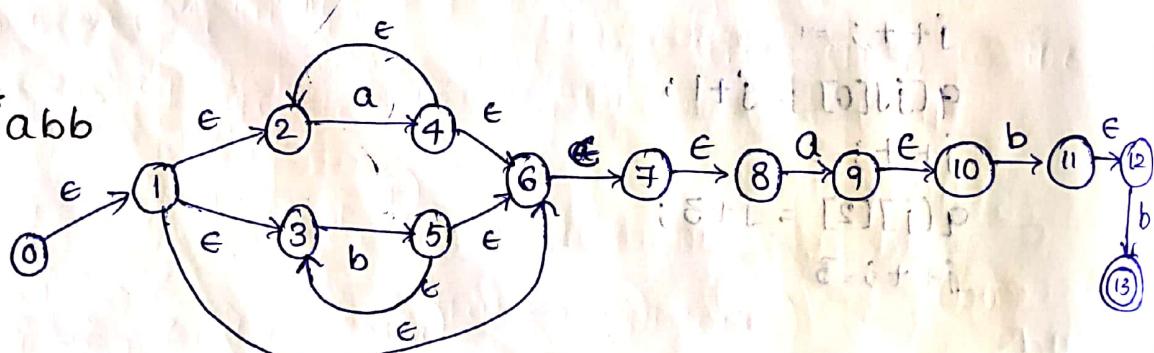
printf("\n In q[%d, e] → %d", i, q[i][2]/100,

          & %100);

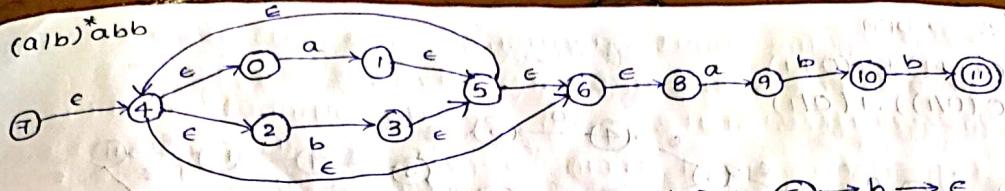
3

3

$(a/b)^*abb$



$$\begin{array}{l} 0 \rightarrow a \\ 1 \rightarrow b \\ 2 \rightarrow \epsilon \end{array}$$



Step: 1:  $\text{A} \rightarrow \text{Start state}$

$$1 \rightarrow \epsilon \quad \{4, 0, 2, 6, 8\} \quad (A)$$

Step: 2:  $(A) \rightarrow a \rightarrow \epsilon$

$$\{0, 2, 5, 6, 1, 4, 8, 9\} \quad (B)$$

Step: 3:-

$$(A) \rightarrow b \rightarrow \epsilon$$

$$\{3, 5, 6, 8, 4, 0, 2\} \quad (C)$$

Step: 4:-

$$(B) \rightarrow a \rightarrow \epsilon$$

$$\{1, 5, 6, 8, 4, 0, 2, 9\} \quad (B)$$

Step: 5:-

$$(B) \rightarrow b \rightarrow \epsilon$$

$$\{3, 5, 6, 8, 4, 0, 2, 10\} \quad (D)$$

Step: 6:-

$$(C) \rightarrow a \rightarrow \epsilon$$

$$\{5, 0, 1, 2, 6, 8, 9, 4\} \quad (B)$$

Step: 7:-

$$(C) \rightarrow b \rightarrow \epsilon$$

$$\{5, 0, 2, 3, 6, 8, 4\} \quad (C)$$

Step: 8:-

$$(D) \rightarrow a \rightarrow \epsilon$$

$$\{5, 6, 8, 9, 4, 0, 2, 1\} \quad (B)$$

Step: 9:-

$$(D) \rightarrow b \rightarrow \epsilon$$

$$\{5, 6, 8, 4, 0, 2, 3, 11\} \quad (E)$$

Step: 10:-

$$(E) \rightarrow a \rightarrow \epsilon$$

$$\{5, 4, 0, 2, 1, 6, 8, 9\} \quad (B)$$

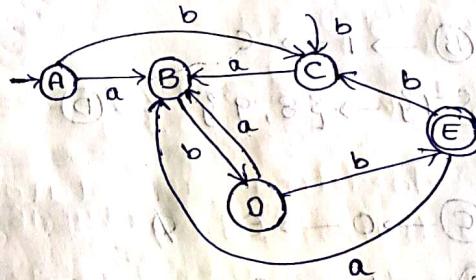
Step: 11:  $(E) \rightarrow b \rightarrow \epsilon$  (C)

$$\{5, 4, 0, 2, 3, 6, 8\}$$

E' is the final state

TRANSITION TABLE

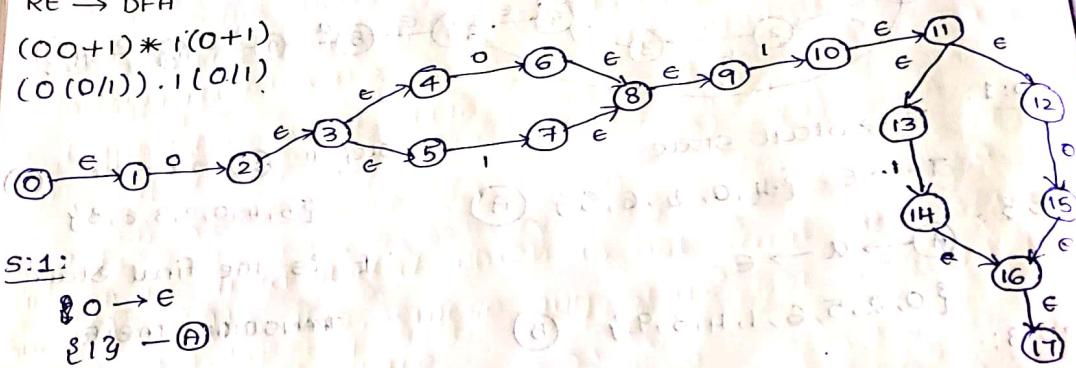
S	a	b
A	B	C
B	B	D
C	B	C
D	B	E
E	B	C



Next problem  
in last page

RE  $\rightarrow$  DFA

$$(00+1)*1(0+1) \\ (0(0/1)).1(0/1)$$



S:1:-

$$\{0 \rightarrow \epsilon \\ \{1\} \rightarrow A\}$$

S:2:-

$$A \rightarrow 0 \rightarrow \epsilon \\ \{2, 3\} \rightarrow \{3, 4, 5\} \rightarrow B$$

$$A \rightarrow 1 \rightarrow \epsilon \quad \emptyset$$

S:3:-

$$B \rightarrow 0 \rightarrow \epsilon \\ \{6\} \rightarrow \{8, 9\} \rightarrow C$$

$$B \rightarrow 1 \rightarrow \epsilon$$

$$\{7\} \rightarrow \{8, 9\} \rightarrow D$$

S:4:-

$$C \rightarrow 0 \rightarrow \epsilon \quad \emptyset$$

$$C \rightarrow 1 \rightarrow \epsilon$$

$$\{10\} \rightarrow \{11, 12, 13\} \rightarrow E$$

S:5:-

$$D \rightarrow 0 \rightarrow \epsilon \quad \emptyset$$

$$D \rightarrow 1 \rightarrow \epsilon$$

$$\{10\} \rightarrow \{11, 12, 13\} \rightarrow E$$

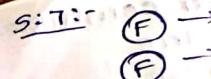
S:6:-

$$E \rightarrow 0 \rightarrow \epsilon$$

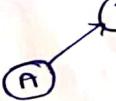
$$\{15\} \rightarrow \{16, 17\} \rightarrow F$$

$$E \rightarrow 1 \rightarrow \epsilon$$

$$\{14\} \rightarrow \{16, 17\} \rightarrow G$$



DFA:-



MINIMI

FINDI

DISTI

STA

INT

ALG

A

SI

D

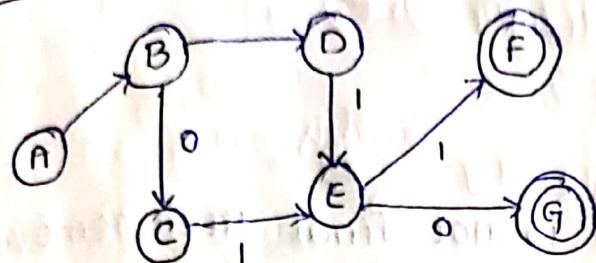
P

S

$$\text{S:7: } \begin{array}{l} F \xrightarrow{o} e \emptyset \\ F \xrightarrow{i} e \emptyset \end{array}$$

$$\text{S:8: } \begin{array}{l} G \xrightarrow{o} e \emptyset \\ G \xrightarrow{i} e \emptyset \end{array}$$

DFA:-



S	0	1
A	B	-
B	C	D
C	-	E
D	F	E
E	F	G
F	-	-
G	-	-

MINIMIZATION OF DFA:-

minimizing the no:of states of DFA works by finding the all the no:of states that can be distinguished by some input string. those states that can not be distinguish are then merge into a single state for the entire group.

ALGORITHM:

A DFA 'M' with a set of states 'S' and initial states 'S<sub>0</sub>', the set <sup>of</sup> final states S<sub>f</sub> the o/p is a DFA minimized accepting the same language but having few states as possible.

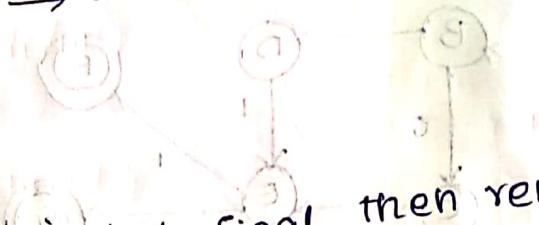
S:1: We construct the partition 'Π' of set of states with two groups i.e., Final state 'F' and non-final state 'S<sub>f</sub>'. Now, we construct a new partition Π<sub>new</sub> by using the following procedure.

for each group G(Π) do

Begin

partition G into subgroups such that 2 states 'S' & 'T' are in the same group iff for all input symbols S & T have same transition place the so formed subgroup as Π<sub>new</sub>.

STEP:2: In the final partition,  $\pi$  that has been constructed in Step:1 pick one representative. Let  $s_i$  be the representative state & for the i/p symbol if there is a transition on  $m$  from  $s \rightarrow t$  then it is replaced by  $s_i$ .



STEP:3  
IF 'm' is dead state which is not final then remove such a state. Also, remove any state not reachable from the initial state.

Eg:-  $(ab)^*abb$

a	b
A	B
C	
B	D
C	B
D	B
E	

{A, B, C, D, E}

S:1 → Identify all states

{A, B, C, D, E}

S:2 → Separate final & non-final states

{A, B, C, D} / {E}

S:3 → Identify any state that has transition from or to final state, then separate.

{ABC} / {DE}

S:4 → Identify the state that has transition to the final states.

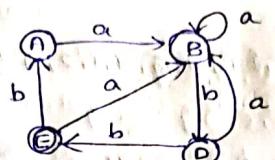
{AC} / {BDE}

S:5 →  $C \Rightarrow A$  (Should be replaced)

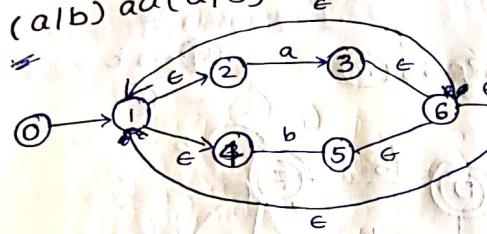
been constructed  
be the  
if there is  
replaced

en remove  
achable

a b  
A B A  
B B D  
D B E  
E B A



$$(a/b)^*aa(a/b)^*$$



$$0 \rightarrow \epsilon$$

$$\{0, 1, 6, 7, 2, 4\} - A$$

$$A \rightarrow a \rightarrow \epsilon$$

$$\{3, 8\} \{6, 7, 1, 2, 4, 9\} - B$$

$$A \rightarrow b \rightarrow \epsilon$$

$$\{5\} \{6, 7, 1, 2, 4\} - C$$

$$B \rightarrow a \rightarrow \epsilon$$

$$\{8, 3, 10\} \{9, 6, 7, 1, 2, 4, 11, 12, 14, 17\} - D$$

$$B \rightarrow b \rightarrow \epsilon$$

$$\{5\} \{6, 7, 1, 2, 4\} \rightarrow C$$

$$C \rightarrow a \rightarrow \epsilon$$

$$\{3, 8\} \{6, 7, 1, 2, 4, 17\} \rightarrow B$$

$$C \rightarrow b \rightarrow \epsilon$$

$$\{5\} \{6, 7, 1, 2, 4\} \rightarrow C$$

$$D \rightarrow a \rightarrow \epsilon$$

$$\{10, 8, 3, 13\} \{11, 12, 14, 17, 9, 6, 7, 1, 2, 4, 16, 17\} \rightarrow E$$

$$D \rightarrow b \rightarrow \epsilon$$

$$\{5, 15\} \{6, 7, 1, 2, 4, 16, 11, 12, 14, 17\} \rightarrow F$$

2020.01.20 13:26

$\oplus \rightarrow \alpha \rightarrow e$

$$\{13, 10, 8, 3\} \quad \{16, 11, 12, 14, 17, 9, 6, 7, 1, 2, 4\} \rightarrow \textcircled{E}$$

$$\textcircled{E} \rightarrow b \rightarrow e$$

$$\{15, 5\} \{6, 7, 1, 2, 4, 16, 11, 12, 14, 17\} \rightarrow \textcircled{F}$$

F → a → e

{8,3,13} {9,6,7,1,2,4,16,11,12,14,17} → ①

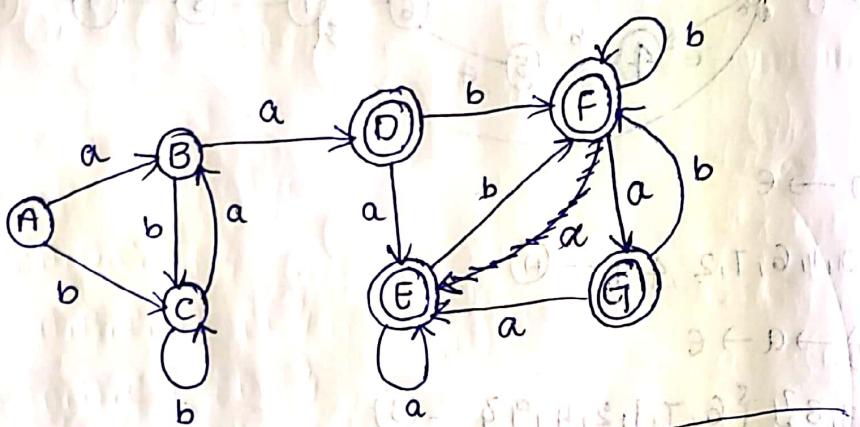
④  $\rightarrow a \rightarrow e$

(F) → b → e

$$\{5, 15\} \{6, 7, 1, 2, 4, 16, 11, 12, 14, 17\} \rightarrow F$$

G → b → e

	a	b
A	B	C
B	D	C
C	B	C
D	E	F
E	E	F
F	G	F
G	E	F

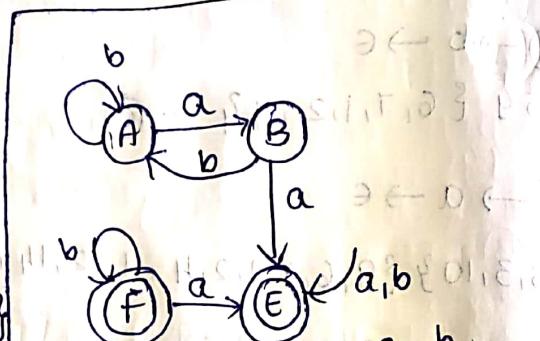


## MINIMIZATION:

{A, B, C, D, E, F, G}

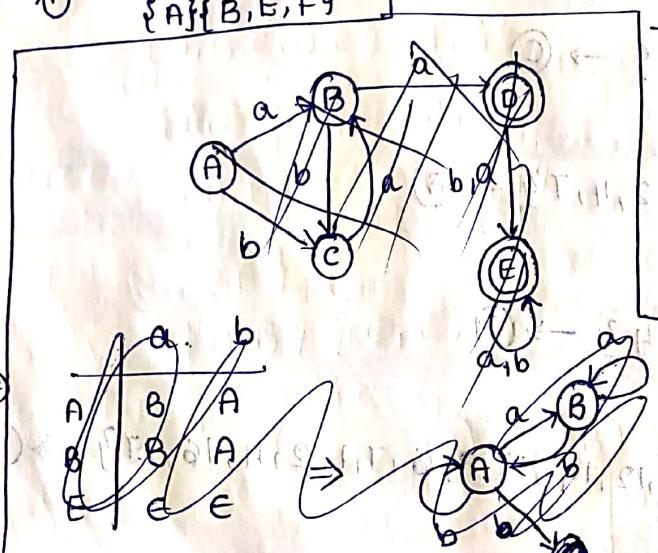
{A,B,C} {D,E,F,G}

$$\{A, C\} \quad \{B, D, E, F\} \Rightarrow \begin{cases} \{A\} \cup \{B, D, E, F\} \\ \{A\} \cap \{B, D, E, F\} \end{cases}$$



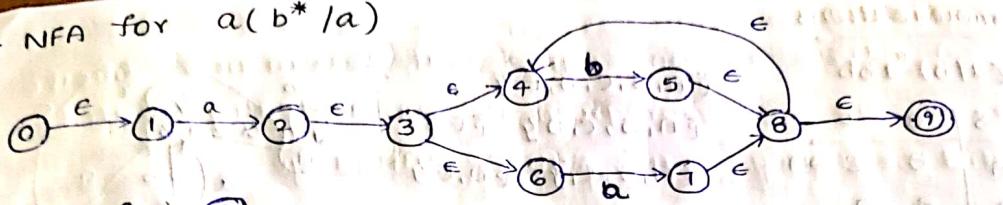
$$A \equiv C$$

	a	b
A	B	A
B	D E	A
D	C	E
E	E	E
F	E	F



	a	b	c
A	B	A	B
B	E	A	D
E	C	E	E
F	E	F	E

\* NFA for  $a(b^*/a)$



$0 \rightarrow \{1\} \rightarrow A$

$A \rightarrow a \rightarrow \epsilon$

$\{2\} \cup \{3, 4, 6\} \rightarrow B$

$A \rightarrow b \rightarrow \epsilon \quad \emptyset$

$B \rightarrow a \rightarrow \epsilon$   
 $\{7\} \cup \{8, 4, 9\} \rightarrow C$

$B \rightarrow b \rightarrow \epsilon$

$\{4, 5\} - \{8, 4, 9\} \rightarrow D$

$C \rightarrow a \rightarrow \epsilon, \emptyset$

$C \rightarrow b \rightarrow \epsilon$

$\{5\} \cup \{8, 4, 9\} \rightarrow D$

$D \rightarrow a \rightarrow \epsilon \quad \emptyset$

$D \rightarrow b \rightarrow \epsilon$

$\{5\} \cup \{8, 4, 9\} \rightarrow D$

$\begin{array}{c|cc} s & a, b \\ \hline A & B \\ B & C, D \\ C & - \\ D & - \end{array}$

$\begin{array}{c|cc} s & a, b \\ \hline A & B \\ B & C, D \\ C & - \\ D & - \end{array}$

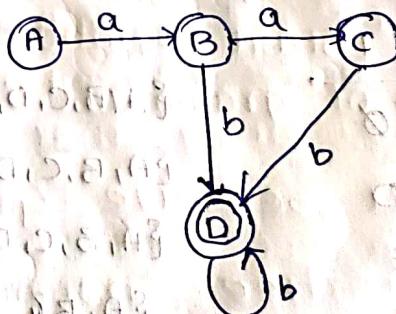
$\begin{array}{c|cc} s & a, b \\ \hline A & B \\ B & C, D \\ C & - \\ D & - \end{array}$

$\begin{array}{c|cc} s & a, b \\ \hline A & B \\ B & C, D \\ C & - \\ D & - \end{array}$

$\begin{array}{c|cc} s & a, b \\ \hline A & B \\ B & C, D \\ C & - \\ D & - \end{array}$

$\begin{array}{c|cc} s & a, b \\ \hline A & B \\ B & C, D \\ C & - \\ D & - \end{array}$

$\begin{array}{c|cc} s & a, b \\ \hline A & B \\ B & C, D \\ C & - \\ D & - \end{array}$



minimizations:

$(a/b)^*abb$

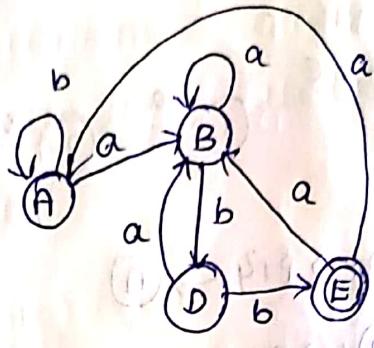
S	a	b
A	B	C → A
B	B	D
C	B	C
D	B	E
E	B	C → A

$\{A, B, C, D, E\}$

$\{A, B, C, D\}$   $\{E\}$

$\{A, B, C\}$   $\{D, E\}$

$\{A, B\}$   $\{D, E\}$



$a(a/b)^*ab \rightarrow$  see Back of the notes

$(a/b)^*abb$

Repeat

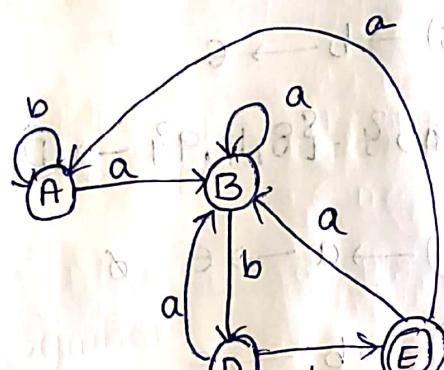
S	a	b
A	B	C
B	B	D
C	B	C
D	B	E
E	B	C

$\{A, B, C, D, E\}$

$\{A, B, C, D\}$   $\{E\}$

$\{A, B, C\}$   $\{D, E\}$

$\{A, C\}$   $\{B, D, E\}$



$(00+1)*1(0+1)$

S	0	1
A	B	$\emptyset$
B	C	D
C	$\emptyset$	E
D	$\emptyset$	E
E	G	F
F	$\emptyset$	$\emptyset$
G	$\emptyset$	$\emptyset$

$\{A, B, C, D, E, F, G\}$

$\{A, B, C, D, E\}$   $\{F, G\}$

$\{A, B, C\}$   $\{E, F, G\}$

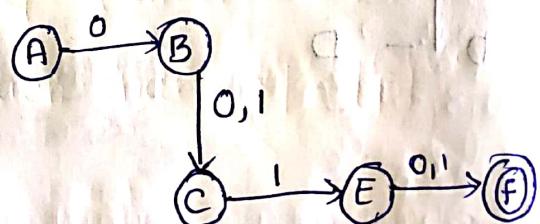
$\{A, B\}$   $\{C, D, E, F, G\}$

~~$\{A, B\}$~~

$\{A, B\}$

$\{C, E, F\}$

S	0	1
A	B	$\emptyset$
B	C	C
C	$\emptyset$	E
E	$\emptyset$	F
F	$\emptyset$	$\emptyset$



### DFA to RE

- \* Kleen's closure construction builds a RE to describe a set of strings generated by the language.
- \* corresponds to path problems

### KLEEN'S

for  $i=0$  to  $|D| - 1$

for  $j=0$  to  $|D| - 1$

$$R_{ij} = \{a \mid s(d_i, a) = d_j\}$$

IF ( $i=j$ ) then

$$R_{ii} = R_{ij} \cup \{\epsilon\}$$

for  $k=1$  to  $|D| - 1$

for  $i=0$  to  $|D| - 1$

for  $j=0$  to  $|D| - 1$

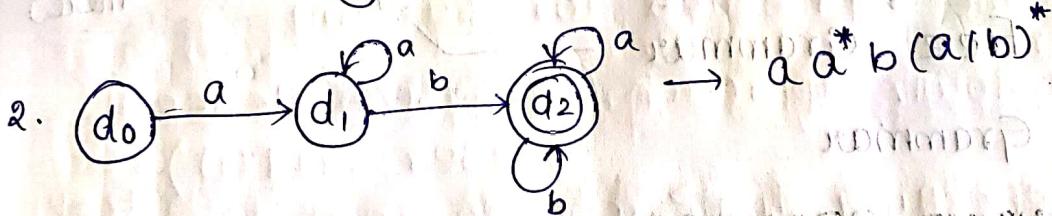
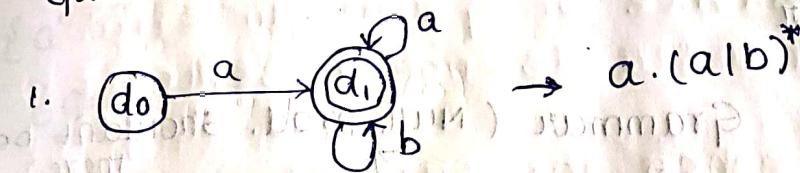
$$R_{ij}^k = R_{ik} (R_{kk}^{k-1})^* R_{kj}$$

$R_{ij}^k = R_{ik} (R_{kk}^{k-1})^* R_{kj}$

- \* Algorithm uses states numbered from 0 to  $|D| - 1$

- \* Systematically, constructs a path expression.

Given DFA denoted by



## PARSERS:

A parser is an engine that determines whether or not the input program is syntactically (valid) or not a valid sentence in the source language.

\* For one programming language, there can be many compilers.

\* To write efficient parsers, we need a formal mechanism for specifying the syntax of the input language.

\* The systematic method of the membership in this formally specified language.

## GRAMMAR:-

A grammar is a four tupled one.  $G$ , where

$$G = \{V_n, V_T, P, S\}$$

$V_n$  → set of non terminal symbols

$V_T$  → set of Terminal symbols

$P$  → set of productions

$S$  → start symbol.

## \* 3 types of grammar

1. Context Sensitive Grammar (Null prod. shouldn't be there)

2. Context free Grammar

3. Regular Grammar

## CONTEXT SENSITIVE GRAMMARS:-

A CSG given by  $G$ , where  $G = \{V_n, V_T, P, S\}$ , if for every production,  $\alpha \rightarrow \beta$  such that  $|\alpha| \leq |\beta|$

### CONTEXT FREE GRAMMARS:-

In CFG, where  $G = \{V_n, V_T, P, S\}$  for every production rule  $\alpha \rightarrow \beta$ , it should follow the below instructions.

1. ' $\alpha$ ' is a single non terminal.

2.  $\beta$  shouldn't have ( $\epsilon$ ) transitions. ( $\beta \neq \epsilon$ )

### REGULAR GRAMMARS:-

A Regular Grammar  $G = \{V_n, V_T, P, S\}$ , if every rule  $p$  is of the form  $A \xrightarrow{\epsilon} aA$  or  $A \xrightarrow{\epsilon} a$  where 'A' is a set of terminal symbols.

\* If a language has a set of strings given by  $\{a, aa, \dots\}$  then grammar is given by  $G = \{V_n, V_T, P, S\}$

$$P \rightarrow$$

$$S \rightarrow A$$

$$\begin{cases} A \rightarrow aA \\ A \rightarrow a \end{cases}$$

$$V_n \rightarrow \{A\}$$

$$V_T \rightarrow \{a\}$$

\* Generate a grammar for a language given by  $\{a^n b^n | n \geq 1\}$

$\{a^n b^n | n \geq 1\}$  we can do substitution in ab which gives  $\{ab, aabb, \dots\}$

$G = \{V_n, V_T, P, S\}$

$$P \Rightarrow S \xrightarrow{\alpha, A \rightarrow} aAb / ab$$

$$V_n \Rightarrow \{S, A\}$$

$$V_T \Rightarrow \{a, b\}$$

\* Generate a grammar for a language given by  $\{a^n b^m | n \geq 1, m \geq 1\}$

$$V_n = \{S, A, B\}$$

$$V_T = \{a, b\}$$

$$S \rightarrow A$$

$$A \rightarrow aAb / B$$

$$B \rightarrow B / \emptyset$$

\* The pro  
a specific

NOTE:

Most

ed with  
TOP DOWN

LL(1):

large  
backtr

IP

conv  
gra

\* H

PAR

\*

\* Generate a grammar for PALINDROME ( $wcw^R$ )

$wcw^R$

$G = \{V_n, V_T, P, S\}$

$S \rightarrow aSa / bSb / cSc / a / b / c$

$V_n = \{S\}$

$V_T = \{a, b, c\}$

\* Generate a grammar for  $(((())))$

$G = \{V_n, V_T, P, S\}$

$S \rightarrow \underline{\underline{S}} / (A) / ()$

$V_n = \{S\}$

$A \rightarrow A$

$A \rightarrow (A) / ()$

$V_T = \{(, )\}$

\* Generate a grammar for arithmetic operations.

$S \rightarrow S + S / S - S / S * S / S / S$

$V_n = \{S\}$

$V_T = \{+, -, *, /\}$

DEF

\* A parser for Grammar  $G$ , is a program that accepts an ip string  $w$ , & produces either the parse tree for  $w$ , if ' $w$ ' is a sentence of  $G$ , or an error message indicating that ' $w$ ' is not a sentence of  $G$ .

\* There are two approaches for parsing

1. Top down approach (Recursive descent parser)  $\xrightarrow{\text{left-lookup}}$

$\hookrightarrow$  checking from root to leaf.

2. Bottom up approach

\* Shift-reduce parser, Left recursive parser (LR parser)

\* The process of constructing a derivation from a specific input sentence is called PARSING.

NOTE:

Most programming languages can be constructed with LR(1) parser (or) LL(1) parser.

TOP DOWN APPROACH:

LL(1): The key to top down approach is from a large set of CFG, it can be passed without backtracking. Hence,

In top down approaches, it transforms (or) convert the grammar into the backtrack free grammar.

\* Hence, it is otherwise known as PREDICTIVE PARSE. (or) BACKTRACKFREE PARSER.

\* There are 2 common topdown process

1. LL(1) - LEFT LOOKUP

2. Hand Coded Recursive Descent parser.

LL(1):

STEPS FOR CONSTRUCTION:

1. Eliminate ambiguity in the grammar.
2. Eliminate Left recursion in the grammar.
3. Eliminate Left factorization in the grammar.
4. calculate 'FIRST' for every production and follow for every non-terminal in the grammar.
5. compute parsing table
6. Finally, check the grammar is LL(1).

Given

### ELIMINATING AMBIGUITY:-

\* Consider the following Grammar G.

$$E \rightarrow E * E$$

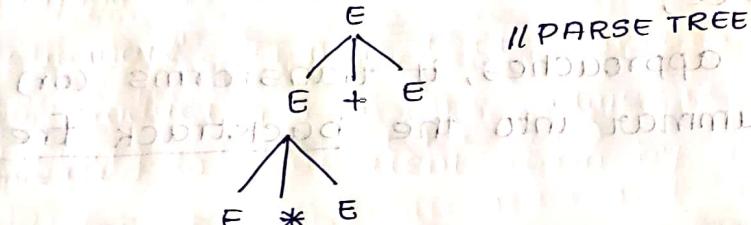
$$| E + E$$

| id

A sentential form  $E * E + E$  has 2 derivations

$$E \rightarrow E * E \rightarrow E * E + E \rightarrow \begin{array}{c} E \\ | \\ E * E \\ | \\ E + E \end{array}$$

$$E \rightarrow E + E \rightarrow E * E + E \rightarrow \begin{array}{c} E \\ | \\ E + E \\ | \\ E * E \end{array}$$



### REASONS FOR AMBIGUITY:

No precedence followed.

No associativity

### TO REMOVE AMBIGUITY:

Ambiguity can be removed by introducing several different variable each representing a binding string from the grammar.

for eg: For the above grammar, 2 variables can be introduced, i.e., 'F' and 'T', if F is a factor, is an identifier and an parenthesized function and cannot be broken by an operator.

't' is a term by ~~is~~  $\Rightarrow$  one or more factor that cannot be broken down by the parser.

The above grammar can be rewritten as

$$E \rightarrow \emptyset$$

$$| E + T$$

$$| T$$

$$T \rightarrow \emptyset$$

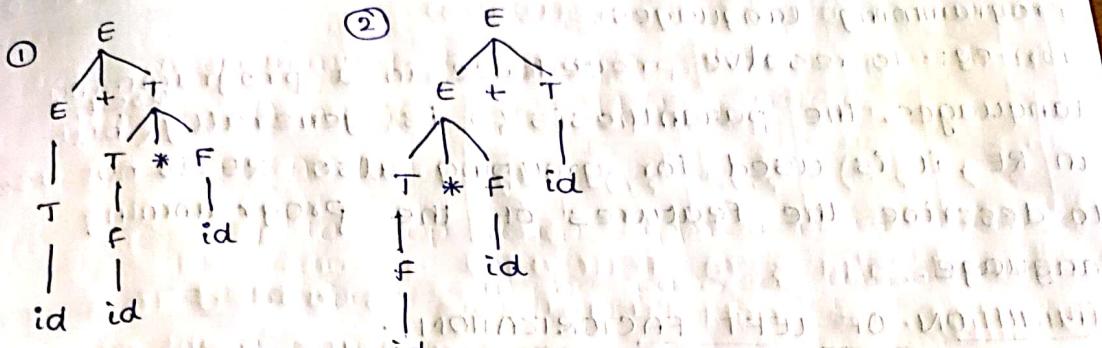
$$| T * F$$

$$| F$$

$$F \rightarrow id$$

$$| (E)$$

Given,  $id + id * id$  ①       $id * id + id$  ②



TO REMOVE LEFT RECURSION: The product of 'T' is left recursive if it has the form  $A \Rightarrow A\alpha/\beta$ , where ' $\alpha$ ' is a grammar symbol and 'A' is the non-terminal to eliminate left recursive. The following rules are followed.

1.  $A \rightarrow \beta A'$       2.  $A' \rightarrow \alpha A' + E$

Now, we have to rewrite the expression in terms of above rules,

$$E \rightarrow E + T$$

$$A \quad A \alpha \beta$$

$$T \rightarrow FT'$$

$$E \rightarrow TE' \quad T \rightarrow *T|e$$

$E' \rightarrow +E'/\epsilon$  (WITERSKI & JACQUES) FOUNDATION FOR  
SOL. SUCH FORMS

\* If there are several such forms  
     $\text{B}_n$  and

\* If there are several add them along the same page in bold

replace by  $B_A$

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

$A' \rightarrow \alpha_1 A' | \alpha_2 A' \dots$  (dmA') ∈

e.g.: -  $A \rightarrow AC \mid Aad \mid bd \mid \epsilon$

NOTE: Why not RE in parsing?

\* RE lacks the power to derive the full syntax of programming languages.

for eg:- To resolve precedence in a programming language, the parenthesis () is generally used, in RE, it is used for grouping. also RE is not used to describe the features of the programming language.

### ELIMINATION OF LEFT FACTORIZATION:-

To eliminate, common left factors that appears in two productions of the same non-terminal. it is converted to equivalent form by removing the common left factors.

For instance ,  $A \rightarrow qb \mid qc$ : The common left factors are removed as the factor is one of the above two productions and if a wrong production is chosen, it leads to backtracking.

Hence, the grammar is converted by introducing a variable that takes out the common production prefixes in the production rule.

### FIRST:

For each grammar symbol,  $\alpha$  FIRST( $\alpha$ ) is the set of terminal symbols. that can appear as the first word in some string derived from  $\alpha$ . The domain of FIRST is the set of grammar symbols including Terminals  $\cup$  Non-Terminal  $\cup \{ \epsilon, \text{eof} \}$

\* If ' $\alpha$ ' is either a terminal,  $\epsilon$  or eof then FIRST( $\alpha$ ) has exactly one member  $\alpha$ ,

\* For a non-terminal A, FIRST(A) contains the complete set of terminal symbols that can appear as a leading symbol in a sentential form derived from A.

(ii) (unambiguous)  $\rightarrow E \rightarrow E + T$

Given,

$$E \rightarrow E + E$$

$$\quad | \quad E^* E$$

$$\quad | \quad id$$

$$T \rightarrow T * E$$

$$\quad | \quad F$$

$$F \rightarrow id$$

$$\quad | \quad (E)$$

remove left recursion:

$$E \rightarrow TE'$$

$$E' \rightarrow +TE'| \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT'| \epsilon$$

$$F \rightarrow id$$

$$F \rightarrow (E)$$

FIRST (T & NT)

TERMINALS:	+	*	$\epsilon$	id	( )
	+	*	$\epsilon$	id	( )

NON-TERMINALS:	E	$E'$	T	$T'$	F
	id, (	+ , $\epsilon$	id, (	* , $\epsilon$	id, (

$$E \rightarrow TE'$$

$$T \rightarrow FT'$$

$$F \rightarrow id$$

FOLLOW FOR NON-TERMINALS: for a Non-Terminal  $\alpha$ , FOLLOW( $\alpha$ ) contains a set of words that occur immediately after  $\alpha$  in a sentence.

Steps:

1. Insert follow(start)  $\leftarrow \$$  (eof)
2. If there is a production  $A \rightarrow \alpha B \beta$  &  $\beta \neq \epsilon$  then everything in FIRST( $\beta$ ) but  $\epsilon$ , in FOLLOW(B).

- Q. If there is a production rule,  $A \rightarrow \alpha B \beta$  or  $A \rightarrow \alpha B^*$ , then everything in FOLLOW(A) will be in FOLLOW(B).

### FOLLOW

E	$E'$	T	$T'$	F
$\$, )$	$\$, )$	$\$, ), +$	$\$, ), +$	<del><math>\\$, +, *</math></del> $\$, +, *$

- \* By rule 1, add  $\$$  to FOLLOW(E)
- \* By rule 2, The right parenthesis is added to E,
- \* By rule 3, on  $E \rightarrow TE'$ , the  $\$$  and the right parenthesis are in FOLLOW( $E'$ )
- \* FOLLOW(E) = FOLLOW( $E'$ ) =  $\$, )$
- \* Similarly, FOLLOW(T), FOLLOW( $T'$ ) and FOLLOW(F) are computed where FOLLOW(T) = FOLLOW( $T'$ ) =  $\$, ), +$

### PREDICTIVE PARSING TABLE:

Input: Grammar  $g$ ; without ambiguity, left recursion and left factorization

Output: parsing table for LL(1).

METHOD 1: For each production,  $A \rightarrow a$  for the grammar do the following steps.

cont'd: Find FIRST & FOLLOW for  $S \rightarrow Bb | Cd$

$S \rightarrow Bb | Cd$        $B \rightarrow aB | \epsilon$        $C \rightarrow CC | \epsilon$

$B \rightarrow aB | \epsilon$        $a, \epsilon$        $b$

$C \rightarrow CC | \epsilon$        $c, \epsilon$        $d$

$(S) \rightarrow Bb$        $a, c, b, d$        $\$, b, d$

$(S) \rightarrow Cd$        $a, \epsilon$        $b$

or  $A \rightarrow \alpha B \beta$

thing in to

point to

suspension

span

3 + 3

3 + 3

3 + 3

3 + 3

to E,

t parenthesis

) are

\$), +

) 78

left

sum

ui

mp.

grammar,

OJ.

rot

be

99

	FIRST	FOLLOW
$S \rightarrow ABCDE$	$a, b, c, \epsilon$	$\$, b, c$
$A \rightarrow aE$	$a, \epsilon$	$b, c$
$B \rightarrow bE$	$b, \epsilon$	$c$
$C \rightarrow C$	$c$	$a, \epsilon, \$$
$D \rightarrow dE$	$d, \epsilon$	$e, \$$
$E \rightarrow eE$	$e, \epsilon$	$\$$

CONT'D:

STEP:2: For each terminal,  $\alpha$  in the FIRST( $\alpha$ ), add  $A \rightarrow \alpha$  to  $M[A, \alpha]$

STEP:3: If ' $\epsilon$ ' is in the FIRST( $\alpha$ ), add  $A \rightarrow \epsilon$  to  $M[A, \epsilon]$  for each terminal in the FOLLOW( $\alpha$ )

STEP:4: If ' $\epsilon$ ' in the FIRST( $\alpha$ ) and the '\$' in the FOLLOW( $\alpha$ ), add  $A \rightarrow \epsilon$  to  $M[A, \$]$

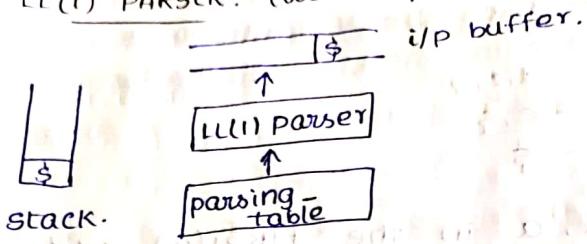
STEP:5: each undefined entity of M is errors

PREDICTIVE TABLE:

	id	+	*	(	)	\$
E	$E \rightarrow TE'$				$E \rightarrow TE'$	
E'		$E' \rightarrow T'E'$			$E' \rightarrow E$	$E \rightarrow E$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T \rightarrow E$	$T \rightarrow *FT$		$T \rightarrow E$	$T \rightarrow E$
F	$f \rightarrow id$				$F \rightarrow (E)$	

	FIRST	FOLLOW
$E \rightarrow TE'$	$id, l$	$\$, )$
$E' \rightarrow +TE'  \epsilon$	$+, \epsilon$	$\$, )$
$T \rightarrow FT'$	$id, l$	$+,\$, )$
$T' \rightarrow *FT'  \epsilon$	$*, \epsilon$	$+,\$, )$
$f \rightarrow id   (E)$	$id, ($	$*  +, \$ )$

## LL(1) PARSER: (WORKING PRINCIPLE)



\* The parser takes 4 possible actions

1. shift
2. reduce
3. accept
4. error

SHIFT: In this action, the next i/p symbol is shifted to the top of the stack.

REDUCE: In this action, the parser knows the right end of the handle is at the top of the stack and must then locate the left end of the handle and decide what non-terminal to replace with handle.

ACCEPT: parser successfully complete parsing.

ERROR: The parser announces error.

check  $id + id * id$

<u>STACK</u>	<u>I/P BUFFER</u>	<u>ACTION</u>
\$ E	$id + id * id \$$	$E \rightarrow TE'$
\$ E' T	$id + id * id \$$	$T \rightarrow FT'$
\$ E' T' F	$id + id * id \$$	$f \rightarrow id$
\$ E' T' id	$id + id * id \$$	reduce
\$ E' T' +	$+ id * id \$$	$T' \rightarrow E$
\$ E' sub.'E'.	$+ id * id \$$	$E' \rightarrow + TE'$
\$ E' T' +	$+ id * id \$$	Reduce
\$ E' T' id	$id * id \$$	$T \rightarrow FT'$
\$ E' T' F	$id * id \$$	$f \rightarrow id$
\$ E' T' id	$id * id \$$	reduce
\$ E' T' *	$* id \$$	<del><math>T' \rightarrow * FT'</math></del>
\$ E' T' F *	$* id \$$	reduce
\$ E' T' F	$id \$$	$f \rightarrow id$
\$ E' T' id	$id \$$	reduce
\$ E' T'	$\$$	$T' \rightarrow E$
\$ E'	$\$$	$E' \rightarrow E$
	$\$$	reduce.

### RECURSIVE DESCENT PARSER:-

A parser that uses the set of recursive procedures to recognize its inputs with no backtracking is called as RECURSIVE DESCENT PARSER.

The two cases that causes backtracking are

1. left factoring (non-deterministic)
2. left recursion (indefinite looping)

The above two problems are overcome in a top down parsers in suitable forms.

LEFT RECURSIVE ( $A \rightarrow A\alpha | \beta$ )  
then 1.  $A \rightarrow \beta A'$   
2.  $A' \rightarrow \alpha A' | \epsilon$

LEFT FACTORING  
 $A \rightarrow q_b$   
 $q_a$   
 $A \rightarrow q_A'$   
 $A' \rightarrow a_{lb}$

### BACKTRACK FREE FORM:

$E \rightarrow TE'$   
 $E' \rightarrow +TE'$   
 $T \rightarrow FT'$   
 $T' \rightarrow *FT | \epsilon$   
 $F \rightarrow id | (E)$

$E \rightarrow TE'$   
procedure  $E()$   
begin  
     $T();$   
    EPRIME();  
end

$T \rightarrow FT'$        $T' \rightarrow *FT | \epsilon$   
procedure  $T()$       procedure  $TPRIME()$   
begin      begin

begin      if ( $\text{inputSymbol} = '*'$ ) then  
     $F();$   
    TPRIME();  
end

if ( $\text{inputSymbol} = '*'$ ) then  
    Advance();  
     $F();$   
    TPRIME();  
end.

$F \rightarrow id | (E)$

procedure  $F()$

begin  
~~RECURSE~~  
RECURSE

if ( $\text{inputSymbol} = 'id'$ ) then

    Advance();

if ( $\text{inputSymbol} = '('$ ) then

$E();$  Advance();

if ( $\text{inputSymbol} = ')'$ ) then

    Advance();

end

$E' \rightarrow +TE'$   
procedure  $EPRIME()$   
begin  
    IF ( $\text{inputSymbol} = '+'$ )  
    then  
        Advance();  
         $T();$   
        EPRIME();  
    end.

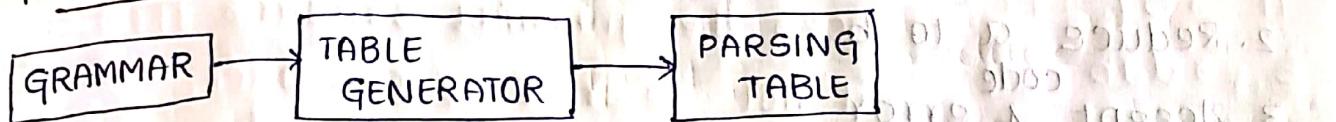
LEFT REC  
bottom scans most FEATURES  
\* LR PCD  
mming LRP  
tech down  
\* LE  
ear GE

LEFT RECURSIVE PARSER: Our aim is to construct efficient bottom up parser for a large class of CFG. The LR parser scans the i/p from left to right and constructs right most derivation in reverse.

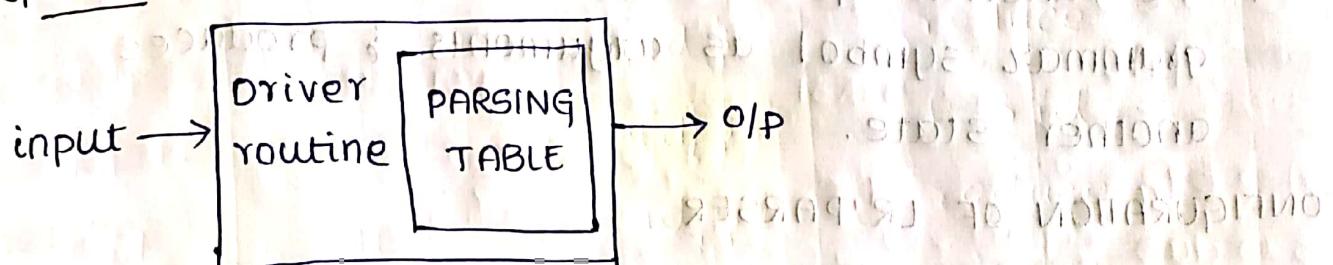
#### FEATURES:

- \* LR parser can be constructed to recognize all programming language constructs for which CFG can be written.
- LR parser method is most general than other parsing techniques also, it dominates the common forms of top down parsing without backtracking.
- \* LR parser can detect the syntactic error at the earliest as it scans the input from left to right.

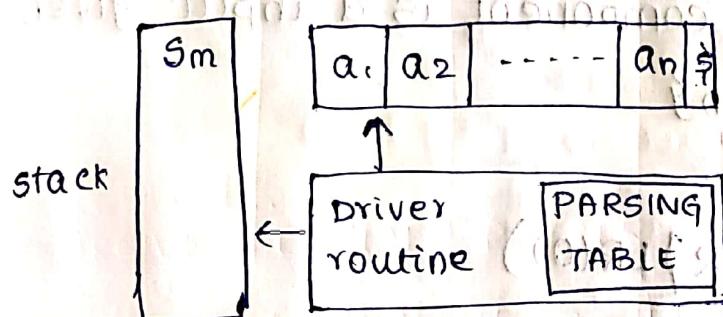
#### GENERATING LR PARSER:



#### Operation:



#### LR Parser:



- \* this parser has the input as stack & parsing table.
- \* the i/p is read one symbol at a time from left to right.

- \* the stack contains the strings of the form,  $s_1, s_2, \dots, s_m, s_m$  where  $s_m$  is the top element in the stack. each  $s_i$  is a grammar symbol and each  $s_i$  is called a state. each state symbol summarizes the information summarized in the stack and is used to guide shift reduce decision.

The parsing table consists of two parts. The parsing action function and 'GOTO' function represented by 'ACTION' & 'GOTO' respectively.

the program driving the LR parsers behaves as follows.

1. It determines  $s_m$ , the state currently on the top of stack and  $a_i$ , the current input symbol.
2. then it consults ACTION OF  $[s_m, a_i]$  the parsing ACTION table entry for the state  $s_m$  and the input entry  $a_i$ .
3. The ACTION  $[s_m, a_i]$  can have one of the 4 values.
  1. Shift the state  $s$ .
  2. Reduce  $a$  to  $b$
  3. Accept A error
4. The function GOTO state takes the state & grammar symbol as arguments & produces another state.

### CONFIGURATION OF LR PARSER

It is a pair whose first component is a stack content and the second component is a input given by  $\{S_0 X_1, S_1 X_2, \dots, X_m S_m\}$

TYPES:-

1. Simple LR parser (SLR) LR(0)
2. Look Ahead Left Recursive parser (LALR/LR(1))
3. Canonical left recursive parser (CLR/LR(2))

### SIMPLE LR PARSER:

LR(0) is a item of grammar  $G$ , is a production of  $G$  with a .(dot) at some position. @ right side

OF the production.

$$E \rightarrow . E + T$$

$$E \rightarrow E. + T$$

$$E \rightarrow E + . T$$

$$E \rightarrow E + T.$$

\* For a given production  $S \rightarrow xyz$  Generate LR(0) itemsets

$$S \rightarrow .xyz$$

$$S \rightarrow x.yz$$

$$S \rightarrow xy.z$$

$$S \rightarrow xyz.$$

NOTE: for a production rule  $A \rightarrow e$ , the LR(0) itemsets is only one  $A \rightarrow .$

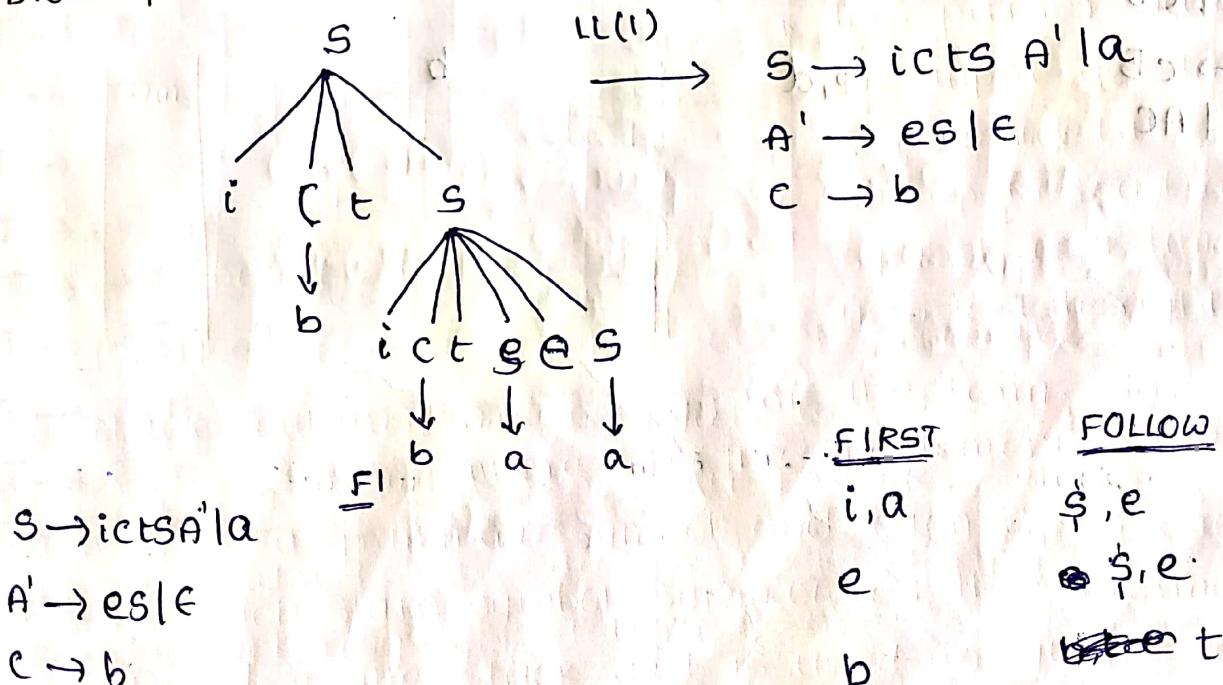
AUGMENTED GRAMMAR: IF  $G$  is a grammar with a start symbol  $S$ , then  $G'$  is the augmented grammar for  $G$ .

$S'$  is the start symbol given by the production

$$S' \rightarrow .S$$

\* The purpose is to indicate when it should stop and announce acceptance, this gives  $S' \rightarrow .S$  if parser reduces.

\* Draw parse tree for "IF-THEN-ELSE STATEMENT"



### PREDICTIVE TABLE:

	i	t	a	b	e	<del>ε</del>
S	$S \rightarrow icta'$		$S \rightarrow a$			
A'						$A' \rightarrow es$
C				$c \rightarrow b$		
						$A' \rightarrow e$

$$\begin{cases} A' \rightarrow es \\ A' \rightarrow e \end{cases}$$

~~here~~

$\rightarrow$  ambiguity

LL(1)

$\therefore$  So, we can't construct the LL(1) grammar.

CONSIDER THE GRAMMAR.

$$A \rightarrow Ba$$

Does the grammar satisfies LL(1),  
if not rewrite & construct.

$$B \rightarrow dab$$

$$| cb$$

$$C \rightarrow cB$$

$$| AC$$

$$d, c$$

$$\$ , c$$

$$A \rightarrow Ba$$

$$B \rightarrow dab$$

$$| cb$$

$$c, d$$

$$C \rightarrow cB$$

$$| AC$$

$$c, d$$



2020.01.20 13:27

CONSIDER THE GRAMMAR

$S \rightarrow E$   
 $E \rightarrow T$   
 $T \rightarrow (E)$   
 $T \rightarrow n$   
 $T \rightarrow +T$   
 $T \rightarrow T+n$

construct LR(1) item sets

$LR(1) = LR(0) + \text{Look ahead}$

	<u>FIRST</u>	<u>FOLLOW</u>	
$S$	$(, n, +$	$\$ )$	
$E$	$(, +, n$	$\$ , )$	
$T$	$+ , n$	$\$ , ), +$	
$LR(1)$	$(S_0)$	$(S_1)$	$(S_2)$
$S \rightarrow .E , \$$	$S \rightarrow E . , \$ \times$	$E \rightarrow (E .), \$$	$E \rightarrow (E .), \$ \times$
$E \rightarrow .T , \$$	$E \rightarrow T . , \$ \times$	$E \rightarrow (E .), )$	$E \rightarrow (E .), ) \times$
$E \rightarrow .T , )$	$E \rightarrow T . , \$ ) \times$	$T \rightarrow +T . , \$ \times$	$T \rightarrow T+n . , \$ \times$
$E \rightarrow .(E) , \$$	$E \rightarrow (.E) , \$$	$T \rightarrow +T . , ( ) \times$	$T \rightarrow T+n . , ( ) \times$
$E \rightarrow .(E) , )$	$E \rightarrow (.E) , )$	$T \rightarrow +T . , + ( ) \times$	$T \rightarrow T+n . , + ( ) \times$
$T \rightarrow .n , \$$	$T \rightarrow n . , \$ \times$	$T \rightarrow T+n , \$$	
$T \rightarrow .n , )$	$T \rightarrow n . , ) \times$	$T \rightarrow T+n , )$	$(S_4)$
$T \rightarrow .n , +$	$T \rightarrow n . , + \times$	$T \rightarrow T+n , +$	NO states
$T \rightarrow .+T , \$$	$T \rightarrow +T , \$$		
$T \rightarrow .+T , )$	$T \rightarrow +T , \$ )$		
$T \rightarrow .+T , +$	$T \rightarrow +T , +$		
$T \rightarrow .T+n , \$$	$T \rightarrow T+n , \$$		
$T \rightarrow .T+n , )$	$T \rightarrow T+n , )$		
$T \rightarrow .T+n , +$	$T \rightarrow T+n , +$		

CANONICAL  
FORM OF  
ITEM SETS

consider the parenthesis grammar.

Goal  $\rightarrow$  list

list  $\rightarrow$  list pair

| pair

pair  $\rightarrow$  (pair)

| ()

what are the LR(1) items for this grammar?  
construct canonical collections of set of LR(1)

follow

Goal  $\rightarrow$  list

list  $\rightarrow$  list pair

| pair

pair  $\rightarrow$  (pair)

| ()

Goal (.

\$

list (.

\$

pair (.

\$

+ .(,

) .(,

+ .(,

) .(,

+ .(,

) .(,

+ .(,

) .(,

+ .(,

) .(,

+ .(,

) .(,

+ .(,

) .(,

+ .(,

) .(,

+ .(,

) .(,

+ .(,

) .(,

+ .(,

) .(,

+ .(,

) .(,

LR(1)

S<sub>0</sub>

Goal  $\rightarrow$  .list, \$

list  $\rightarrow$  .list pair, \$

list  $\rightarrow$  .pair, \$

pair  $\rightarrow$  T : (pair), )

pair  $\rightarrow$  .(pair), )

pair  $\rightarrow$  .(pair), \$

S<sub>1</sub>

(.S<sub>1</sub>)

(.)

(.)

list  $\rightarrow$  list pair, \$

pair  $\rightarrow$  (pair), )

pair  $\rightarrow$  (pair.), )

pair  $\rightarrow$  (pair.), \$

pair  $\rightarrow$  (pair.), \$

Bottom-up

$\rightarrow$  all programming

constructs

LR(1)

procedure closure (I)

begin

Repeat

for each item  $[A \rightarrow \alpha.B\beta.a]$  in I

for each production  $B \rightarrow \gamma$  & each terminal

terminal b in  $\text{FIRST}(\beta a)$  such

that  $[B \rightarrow .\gamma.b]$  is not in I do add

$[B \rightarrow .\gamma.b]$  to I until no more items to

be added to I

return I

end

procedure GOTO (I, x)

begin

Let I be the set of items  $[A \rightarrow \alpha.x.\beta.a]$  s.t

$[A \rightarrow \alpha.x\beta.a]$  is in I

return closure [I]

end

begin

$c = \{\text{closure } ((S' \rightarrow .S,\$))\}$

repeat

for each set of items

I in c & grammar symbol x s.t

I in c & grammar symbol x s.t  
GOTO (I, x) is not empty & already in c

GOTO (I, x) to C

add GOTO (I, x) to C

until no more sets of items

added to C

end.

### AUGMENTED GRAMMAR:

	<u>FIRST</u>	<u>FOLLOW</u>
1. $S' \rightarrow S$	c, d	\$
$S \rightarrow CCG$	c, d	\$
$C \rightarrow CC/d$	c, d	c/d

ITEMSETS:

$$I_0 : S' \rightarrow .S, \$$$

$$S \rightarrow .CCG, \$$$

$$C \rightarrow .CG, c/d$$

$$G \rightarrow .d, c/d$$

$$\boxed{GOTO(I_0, S) : S' \rightarrow S., \$ \rightarrow I_1}$$

$$GOTO(I_0, G) : S \rightarrow C.G, \$$$

$$C \rightarrow .CG, \text{ etc } \$ \} I_2$$

$$G \rightarrow .d, \text{ etc } \$$$

$$GOTO(I_0, C) : C \rightarrow c.C, c/d \} I_3$$

$$C \rightarrow .d, c/d$$

$$C \rightarrow .c.G, c/d$$

$$\boxed{GOTO(I_0, d) : C \rightarrow d., c/d \quad I_4}$$

$$\boxed{GOTO(I_2, G) : S \rightarrow CG., \$ \quad I_5}$$

$$GOTO(I_2, C) : C \rightarrow c.G, \text{ etc } \$$$

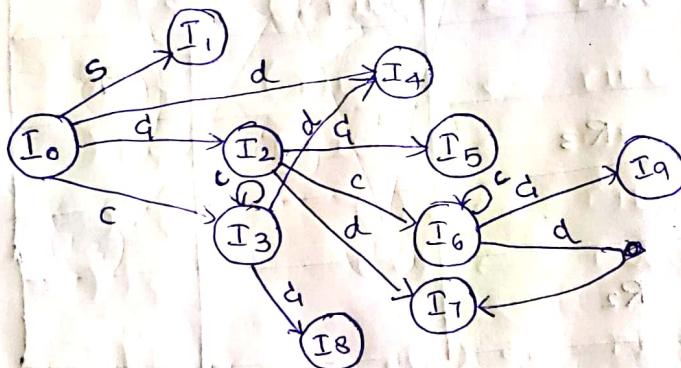
$$C \rightarrow .d, \text{ etc } \$ \quad I_6$$

$$C \rightarrow .c.G, \text{ etc } \$$$

$$\boxed{GOTO(I_2, d) : C \rightarrow d., \cancel{\text{etc } \$} \quad I_7}$$

$$\boxed{GOTO(I_3, G) : C \rightarrow CG., c/d \quad I_8}$$

$\text{GOTO}(I_3, c) : \begin{array}{l} G \rightarrow c \cdot G, c/d \\ G \rightarrow \cdot c G, c/d \\ G \rightarrow \cdot d, c/d \end{array}$   $I_3$   
 $\text{GOTO}(I_3, d) : \boxed{G \rightarrow d \cdot, c/d I_4}$   
 $\text{GOTO}(I_6, \cdot) : \boxed{G \rightarrow c G \cdot, \$ I_9}$   
 $\text{GOTO}(I_6, c) : \begin{array}{l} G \rightarrow c \cdot G, \$ \\ G \rightarrow \cdot c G, \$ \\ G \rightarrow \cdot d, \$ \end{array} I_6$   
 $\text{GOTO}(I_6, d) : \boxed{G \rightarrow d \cdot, \$ I_7}$



### LR(1) PARSING TABLE:

1. If  $[A \rightarrow \alpha \cdot a\beta, b]$  is in  $I_i$ ; &  $\text{GOTO}(I_i, a) = I_j$   
then set ACTION  $[i, a]$  to shift j
2. If  $[A \rightarrow \alpha \cdot, a]$  is in  $I_i$ ; then set ACTION  $[i, a]$  to  
Reduce  $A \rightarrow \alpha$
3. If  $[S' \rightarrow S \cdot, \$]$  is in  $I_i$ , then set ACTION  $[i, \$]$  to  
Accept
4. If  $\text{GOTO}(I_i, A) = I_j$  then  $\text{GOTO}(i, A) = j$  [Non-terminal]

Action : ~~non~~ Terminal,  
~~non~~ Terminal  $\rightarrow$  non-Terminal  
 $\text{GOTO}$  : nonTerminal  $\rightarrow$  non-Terminal

Cons  
collie  
che

	ACTION			S	C
	C	d	\$		
0					2
1	$s_3$	$s_4$		Accept	
2					5
3	$s_3$	$s_4$			8
4	$R_3$	$R_3$			
5			$R_1$		
6	$s_6$	$s_1$			9
7			$R_3$		
8	$R_2$	$R_2$			
9			$R_2$		

$$\begin{aligned} S &\rightarrow Gd \\ G &\rightarrow Cd \\ G &\rightarrow d \end{aligned}$$

\*check if the string cdd can be parsed

word      stack      action

c       $\$_0$        $S_3$        $S_3 \rightarrow S_3$        $S_3 \rightarrow S_3$

d       $\$_0 C_3$        $S_4$

d       $\$_0 C_3 d_4$        $R_3$        $d \rightarrow d$

d       $\$_0 C_3 C_8$        $R_2$        $d \rightarrow Cd$

d       $\$_0 C_2$       GOTO 2       $i1 = [A, 1]$  or op

$\$$        $\$_0 C_2 d_7$

$\$$        $\$_0 C_2 C_5$

$\$$        $\$_0 S_1$       Accept

Construct the parenthesis grammar & canonical collection of set of items. derive the table and check it is LR(1). And.

FIRST

Goal $\rightarrow$ List	, ( )	( )	\$   ( )
List $\rightarrow$ List pair	, ( )	( )	\$   ( )
pair $\rightarrow$ (pair)	, ( )	( )	\$   ( )

I <sub>0</sub> : Goal $\rightarrow$ .List	,	\$	
List $\rightarrow$ . List pair	,	\$   ( )	
List $\rightarrow$ . Pair	,	\$   ( )	
Pair $\rightarrow$ . (pair)	,	\$   ( )	
Pair $\rightarrow$ . ( )	,	\$   ( )	

GOTO(I <sub>0</sub> , List) :	Goal $\rightarrow$ List .	,	\$
	List $\rightarrow$ List . pair	,	\$   ( )
	List $\rightarrow$ . pair	,	\$   ( )
	pair $\rightarrow$ . (pair)	,	\$   ( )
	pair $\rightarrow$ . ( )	,	\$   ( )

GOTO(I <sub>0</sub> , pair) :	List $\rightarrow$ pair .	,	\$   ( )
GOTO(I <sub>0</sub> , () :	pair $\rightarrow$ . (pair)	,	\$   ( )
	pair $\rightarrow$ . (pair)	,	\$   ( )
	pair $\rightarrow$ . ( )	,	\$   ( )
	pair $\rightarrow$ (.)	,	\$   ( )

GOTO(I <sub>1</sub> , pair) :	<del>pair <math>\rightarrow</math> . pair</del>	,	\$   ( )
	List $\rightarrow$ List pair .	,	\$   ( )
	pair $\rightarrow$ ( . pair)	,	\$   ( )
	pair $\rightarrow$ . (pair)	,	\$   ( )
	pair $\rightarrow$ . ( )	,	\$   ( )

	(
0	S <sub>3</sub>
1	S <sub>5</sub>
2	R <sub>2</sub>
3	S <sub>3</sub>
4	R
5	S
6	
7	R
8	
9	
10	
11	

pair  $\rightarrow (\cdot)$ , \$ | ( ) , \$ | ( ) I<sub>6</sub>

GOTO(I<sub>3</sub>, pair) : pair  $\rightarrow (\text{pair} \cdot)$ , \$ | ( )

GOTO(I<sub>3</sub>, ( )) : pair  $\rightarrow (\cdot \text{pair})$ , \$ | ( ) I<sub>3</sub>

pair  $\rightarrow \cdot(\text{pair})$ , \$ | ( )

pair  $\rightarrow \cdot(\cdot)$ , \$ | ( )

pair  $\rightarrow (\cdot \cdot)$ , \$ | ( )

GOTO(I<sub>3</sub>, )) : pair  $\rightarrow (\cdot \cdot)$ , \$ | ( ) I<sub>7</sub>

GOTO(I<sub>5</sub>, pair) : pair  $\rightarrow (\text{pair} \cdot)$ , \$ | ( ) I<sub>8</sub>

GOTO(I<sub>5</sub>, ( )) : pair  $\rightarrow (\cdot \text{pair})$ , \$ | ( )

pair  $\rightarrow \cdot(\text{pair})$ , \$ | ( ) I<sub>5</sub>

pair  $\rightarrow \cdot(\cdot)$ , \$ | ( )

pair  $\rightarrow (\cdot \cdot)$ , \$ | ( )

GOTO(I<sub>5</sub>, )) : pair  $\rightarrow (\cdot \cdot)$ , \$ | ( ) I<sub>9</sub>

GOTO(I<sub>6</sub>, )) : pair  $\rightarrow (\text{pair} \cdot)$ , \$ | ( ) I<sub>10</sub>

GOTO(I<sub>8</sub>, )) : pair  $\rightarrow (\text{pair} \cdot)$ , \$ | ( ) I<sub>11</sub>

Goal  $\rightarrow$  list

R<sub>1</sub> - List  $\rightarrow$  List pair

R<sub>2</sub> - pair  $\rightarrow$  pair

R<sub>3</sub> - pair  $\rightarrow$  (pair)

R<sub>4</sub> - ( )  $\rightarrow$  ( )

R<sub>5</sub> - ( )  $\rightarrow$  ( )

R<sub>6</sub> - ( )  $\rightarrow$  ( )

R<sub>7</sub> - ( )  $\rightarrow$  ( )

R<sub>8</sub> - ( )  $\rightarrow$  ( )

R<sub>9</sub> - ( )  $\rightarrow$  ( )

2020.01.20 13:28

	ACTION	\$	list	go TO pair
0	$S_3$	-	1	2
1	$S_5$	Accept		4
2	$R_2$	$R_2$		
3	$S_3$	$S_7$	-	6
4	$R_1$		$R_1$	
5	$S_5$	$S_9$		8
6		$S_{10}$		
7	$R_4$	$R_4$	$R_4$	
8		$S_{11}$		
9	$R_4$		$R_4$	
10	$R_3$	$R_3$	$R_3$	
11	$R_3$	$R_3$	$R_3$	

check if the string () can be parsed.

word      stack      action

(	$\$_0$	$S_3$
)	$\$_0 (3$	$S_3$
)	$\$_0 (3 (3$	$S_7$
)	$\$_0 (3 (3 (7$	$R_4$ pair $\rightarrow ()$
)	$\$_0 (3 (3 (7$	$S_{10}$
\$	$\$_0 (3 pair_6 )_{10}$	$R_3$ pair $\rightarrow (pair)$
\$	$\$_0 pair_2$	$R_2$ list $\rightarrow pair$
\$	$\$_0 list_1$	Accept

## UNIT - II

### CONTEXT SENSITIVE ANALYSIS:

→ context of

1) Variable (data type)

2) expression (type compatibility)

3) procedure calls

compiler builds a large knowledge base from programs.

↓  
about details in the program

↓ to do it, it should know

1. How values are represented?

2. How values flow between  
variable?

3. understand structure

4. program interaction with external  
files.

### knowledge base:

ex: if var 'x' is used in program, it checks

1. what kind of values are stored.

2. memory requirements ⇒ (overflow can be avoided)

3. If 'x' is a procedure, checks args & return type

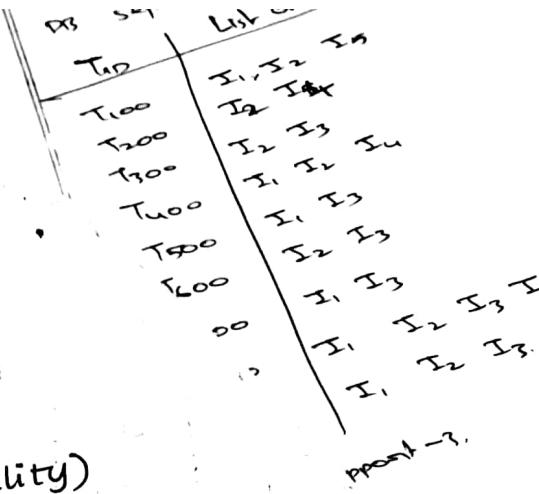
4. memory allocation

5. How long it is to be stored

compile  
(or)  
runtime.

### Type System:

→ The set of types in programming language  
along with rules.



TID	List
100	I <sub>1</sub> , I <sub>2</sub>
200	I <sub>2</sub> I <sub>1</sub>
00	I <sub>2</sub> I <sub>1</sub>
00	I <sub>1</sub> , I <sub>2</sub>
00	I <sub>1</sub> , I <sub>2</sub>
0	I <sub>2</sub>
0	I <sub>1</sub> , I <sub>2</sub>
I <sub>1</sub>	I <sub>1</sub>
I <sub>1</sub>	I <sub>1</sub>

Need for type system.

1. ensure runtime safety } avoid compile time errors.  
ex:  $a+b$

ex:  $a+b$

+ int r d c

int i r d c

real  $\gamma$   $\tau$   $d$   $c$

double d      d      d . not possible

complex C      C.      n.p      C

two types

strongly  
typed  
(C, C++, Java)

weakly typed

cpython

## 2. Expressiveness

→ precise meaning

## addition

ex: + (or) overloading

3. Generate better code

$$a + b = a + b$$

int      int      int

int real real

int ~~real~~

**TYPE CHECKING:** compiler analyse program and assign types to each expression also ensure context - this is called type checking.

Type checking: 2 ways

Checking:- 2 ways  
attribute grammar framework  
directed

- attribute grammar
- ASDT [adhoc Syntax Directed Translation]

attribute grammar frameworks:

→ attribute grammar on CFG.



Set of rules on CFG



each rule specifies an attribute in terms of another.

I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub>  
I<sub>2</sub>, I<sub>4</sub>  
I<sub>3</sub>  
I<sub>2</sub>, I<sub>4</sub>

ex: CFG for Signed Binary no's (SBN)

$$\{ S, T, NT, P \}$$

$$S = \{ \text{Number} \}$$

$$T = \{ -, +, 1, 0 \}$$

$$NT = \{ \text{Number}, \text{list}, \text{Bit}, \text{sign} \}$$

$$P = \{ S \rightarrow \text{Number}$$

$$\text{Number} \rightarrow \text{sign List}$$

$$\text{Sign} \rightarrow - | +$$

$$\text{List} \rightarrow \text{List bit} | \text{bit}$$

$$\text{bit} \rightarrow 0 | 1$$

Attribute grammar to each non-terminal

<u>Symbol</u>	<u>attributes</u>
Number	value
Sign	negative
List	position value
bit	position value

Attributes grammar for SBN

productions

$$1. \text{ Number} \rightarrow \text{Sign list}$$

attribute rules

$$\text{list.position} \leftarrow 0$$

IF sign. negative then

$$\text{number.val} \leftarrow -\text{list.val}$$

else

$$\text{number.val} \leftarrow \text{list.val}$$

1. sign  $\leftarrow +$
2. sign  $\leftarrow -$
3. list  $\leftarrow$  bit

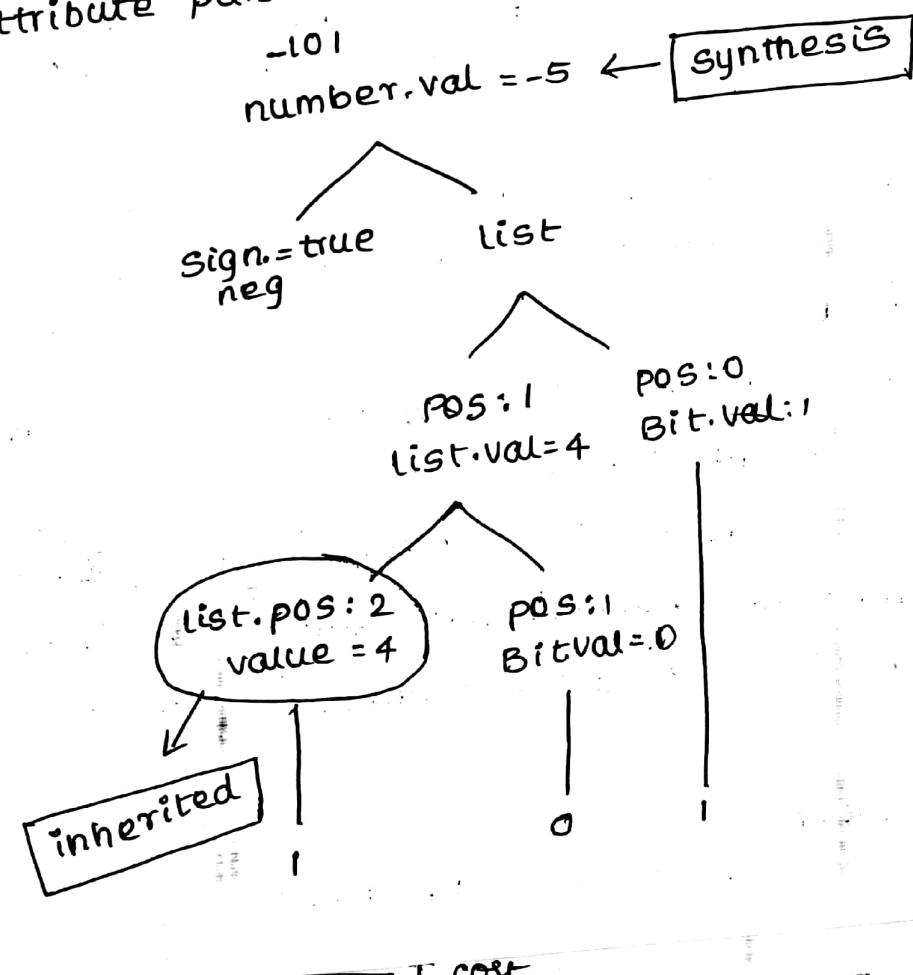
5. list  $\leftarrow$  list bit

6. Bit  $\leftarrow 0$
7. Bit  $\leftarrow 1$

## 2M TWO TYPES OF ATTRIBUTES

- \* synthesized  $\rightarrow$  only on the own node from children
- \* inherited  $\rightarrow$  depended on its own node, parent, sibling, children.

ex:- attribute parse tree for -101



To evaluate the attributed parse tree. the sentence of the grammar attributes are instantiated. for each node in the parse tree. each rule defines the set of dependancies and depends on each argument in the rule.

Write the attribute grammar for signed decimal (SDN)

// FORMALITIES IN PRE-PAGE

Number → Sign List

Sign → + / -

List → bit / List bit . bit

bit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

\* A compiler writer may create attribute grammar & design set of rules can be created to do it.

EVALUATION METHODS:



### 1. Dynamic

- depends on structure of attributed parse tree
- rule fixed as soon as operands are available

### 2. Oblivious

- independent on order
- designer selects method for evaluation

### 3. Rule based

- it depends on static analysis and on grammar.

### PROBLEMS WITH THE RULE BASED ATTRIBUTE GRAMMAR:

- refer to its own results either directly or indirectly
- such attribute grammar are called CIRCULAR ATTRIBUTE GRAMMAR. ( $\rightarrow$  fail to get evaluated  $\xrightarrow{\text{so}}$  removed.)

### METHODS TO REMOVE:

#### 1. Avoidance

- $\rightarrow$  restrict on CFG to avoid.

use evaluation → to identify circularity.

## Ex: 2 EXPRESSION TYPES.

Inter expression type

$$E_0 \rightarrow E_1 + T_1$$

$$| E_1 - T_1$$

$$| T$$

$$T_0 \rightarrow T_1 * F$$

$$| T_1 \% F$$

$$| F$$

$$F \rightarrow (E)$$

$$| id$$

$$| num$$

function attribute rules

$$E_0 \cdot type \rightarrow F^+ (E_1 \cdot type, T_1 \cdot type)$$

$$\rightarrow F^- ( " )$$

$$\rightarrow T \cdot type$$

$$T_0 \cdot type \rightarrow F^* (T_1 \cdot type, F \cdot type)$$

$$\rightarrow F \% ( " )$$

$$\rightarrow F \cdot type$$

$$F \cdot type = (E \cdot type)$$

$$\rightarrow id$$

already def.

$$\rightarrow num$$

F<sup>+</sup>, F<sup>-</sup>, F<sup>\*</sup>, F% → for evaluating.

Eg: Execution time estimator

To estimate execution time of sequence of

\* To estimate execution time of assignment statements

assignment by adding 3 new productions

\* we generate by adding 3 new productions

Black → Black Assignment

| Assignment

Assignment → id = Expr is added to classical

expr grammar:

production

$$Black_0 \rightarrow Black, Assign$$

$$| Assign$$

$$Assign \rightarrow id = Expr$$

attribute rules

$$Black_0 \cdot cost \leftarrow Black \cdot cost +$$

$$Assign \cdot cost$$

$$\leftarrow Assign \cdot cost$$

$$Assign \cdot cost \rightarrow cost(store) +$$

$$E \cdot cost$$

$$E_0 \cdot cost = E_1 \cdot cost + T_1 \cdot cost +$$

$$cost(add)$$

$$E_1 \cdot cost + T_1 \cdot cost + cost(sub)$$

$$T \cdot cost$$

$$E_0 \rightarrow E + T$$

$$| E - T$$

$$| T$$

Top

T<sub>100</sub>T<sub>200</sub>T<sub>300</sub>T<sub>400</sub>T<sub>500</sub>T<sub>600</sub>T<sub>700</sub>T<sub>800</sub>T<sub>900</sub>T<sub>1000</sub>T<sub>1100</sub>T<sub>1200</sub>T<sub>1300</sub>T<sub>1400</sub>T<sub>1500</sub>T<sub>1600</sub>T<sub>1700</sub>T<sub>1800</sub>T<sub>1900</sub>T<sub>2000</sub>T<sub>2100</sub>T<sub>2200</sub>T<sub>2300</sub>T<sub>2400</sub>T<sub>2500</sub>T<sub>2600</sub>T<sub>2700</sub>T<sub>2800</sub>T<sub>2900</sub>T<sub>3000</sub>T<sub>3100</sub>T<sub>3200</sub>T<sub>3300</sub>T<sub>3400</sub>T<sub>3500</sub>T<sub>3600</sub>T<sub>3700</sub>T<sub>3800</sub>T<sub>3900</sub>T<sub>4000</sub>T<sub>4100</sub>T<sub>4200</sub>T<sub>4300</sub>T<sub>4400</sub>T<sub>4500</sub>T<sub>4600</sub>T<sub>4700</sub>T<sub>4800</sub>T<sub>4900</sub>T<sub>5000</sub>T<sub>5100</sub>T<sub>5200</sub>T<sub>5300</sub>T<sub>5400</sub>T<sub>5500</sub>T<sub>5600</sub>T<sub>5700</sub>T<sub>5800</sub>T<sub>5900</sub>T<sub>6000</sub>T<sub>6100</sub>T<sub>6200</sub>T<sub>6300</sub>T<sub>6400</sub>T<sub>6500</sub>T<sub>6600</sub>T<sub>6700</sub>T<sub>6800</sub>T<sub>6900</sub>T<sub>7000</sub>T<sub>7100</sub>T<sub>7200</sub>T<sub>7300</sub>T<sub>7400</sub>T<sub>7500</sub>T<sub>7600</sub>T<sub>7700</sub>T<sub>7800</sub>T<sub>7900</sub>T<sub>8000</sub>T<sub>8100</sub>T<sub>8200</sub>T<sub>8300</sub>T<sub>8400</sub>T<sub>8500</sub>T<sub>8600</sub>T<sub>8700</sub>T<sub>8800</sub>T<sub>8900</sub>T<sub>9000</sub>T<sub>9100</sub>T<sub>9200</sub>T<sub>9300</sub>T<sub>9400</sub>T<sub>9500</sub>T<sub>9600</sub>T<sub>9700</sub>T<sub>9800</sub>T<sub>9900</sub>T<sub>10000</sub>T<sub>10100</sub>T<sub>10200</sub>T<sub>10300</sub>T<sub>10400</sub>T<sub>10500</sub>T<sub>10600</sub>T<sub>10700</sub>T<sub>10800</sub>T<sub>10900</sub>T<sub>11000</sub>T<sub>11100</sub>T<sub>11200</sub>T<sub>11300</sub>T<sub>11400</sub>T<sub>11500</sub>T<sub>11600</sub>T<sub>11700</sub>T<sub>11800</sub>T<sub>11900</sub>

min-supp

$$\begin{aligned}
 T_0 &\rightarrow T_1 * F \\
 T_1 &\div F \\
 \text{IF} \\
 F &\rightarrow (E) \\
 \text{id} \\
 \text{inum} \\
 \text{consider } x=y \text{ for the above} \\
 \text{the load operation, activation record} \\
 \text{load AI} &\quad \text{Yarp} @ y \Rightarrow \gamma_y \\
 \text{add} &\quad \gamma_y, \gamma_x \Rightarrow \gamma_x \\
 \text{store AI} &\quad \gamma_x \Rightarrow \text{Yarp}, @ x \\
 \text{To appropriate above behaviour, for } y \text{ only one load} \\
 \text{operation inside a block then to include} \\
 \text{if (ident has not been loaded) then} \\
 F &\rightarrow \text{cost(load)} \\
 \text{else} &\quad f.\text{cost} \leftarrow 0
 \end{aligned}$$

To implement the above we include before & after to each node.

Before  $\rightarrow$  contains names of id in the block-  
after  $\rightarrow$  before + any id that would be loaded.

After modification:

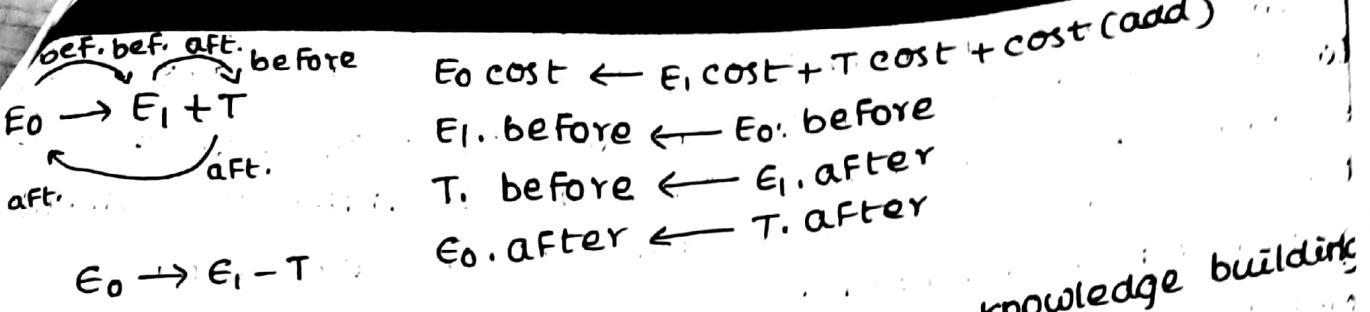
F  $\rightarrow$  (E)

F  $\rightarrow$  num

F  $\rightarrow$  id

### Attribute

$f.\text{cost} \leftarrow e.\text{cost}$   
 $e.\text{before} \leftarrow f.\text{before}$   
 $f.\text{before} \leftarrow e.\text{after}$   
 $f.\text{cost} \leftarrow \text{cost(load num)}$   
 $f.\text{after} \leftarrow f.\text{before}$   
 IF (idname  $\notin$  f.before) then  
 $f.\text{cost} \leftarrow \text{cost(load)}$   
 $f.\text{after} \leftarrow f.\text{before} \cup \text{id.name}$   
 else  
 $f.\text{cost} \leftarrow 0$   
 $f.\text{after} \leftarrow f.\text{before}$ .



PROBLEMS OF ATTRIBUTE GRAMMAR (AGF)  $\Rightarrow$  CSA knowledge building

- \* Handling Non local infor. Hence, rules are included as requirement.
  - \* It constructs parse tree for evaluation but it is a known fact that compiler never builds a parse tree as it constructs only a abstract syntax.
  - \* To locate the rules it needs to traverse through the parse tree to find the appropriate one.
- Adhoc Syntax Directed Translation (ASDT) (or) (ADST): The problems with AGF is the increase in complexity & space requirement to overcome the above problem ASDT does the following:

1. associates a snippet of code with each production.
2. At each reduction of the rule the corresponding snippet runs.

IMPLEMENTING ASDT: To implement ASDT, the parser must include the mechanisms to parse value from their definitions in one action to the users & in another to provide convenient meaning allowing for the actions for the parser. Some of the mechanisms for handling are

- i. COMMUNICATE BIN THE ACTION: In ASDT, the parser doesn't construct a parse tree so, the values are to be stored in an internal stack. usually, the parse tree stores as a pair of info. given by <symbol, state>, for each grammar rule & now for implementing ASDT, the parser includes values & hence the triple info. is given by <value, symbol state> is stored for every grammar symbol. To manage the stack the values are pushed or popped up.

for instance a reduction given by  $A \rightarrow B$  leads to locations popped from the top of stack. To solve this storage, the parser may omit actual grammar symbols but the grammar symbol provides better error reporting & debugging.

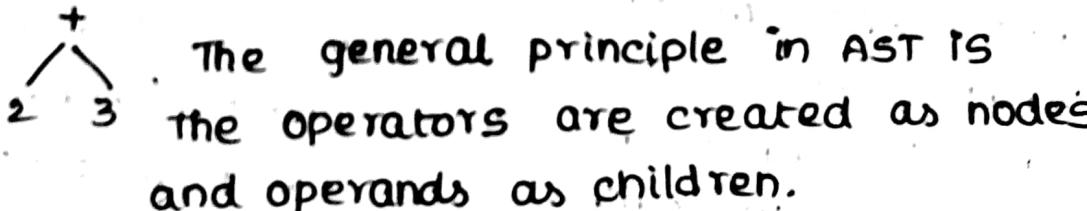
NAMING VALUES: To simplify the use of stack based values, a convenient notation is required for naming them. The symbol  $\$S$  refers to the current location for the current production of the right hand symbol. Hence, the symbol  $\$1, \$2, \dots, \$n$  refers to 1<sup>st</sup>, 2<sup>nd</sup> & n<sup>th</sup> location in RHS for eg: production code snippet

number  $\rightarrow$  signlist  $\$\$ \leftarrow \$1 * \$2$   
 sign  $\rightarrow + \quad \$\$ \leftarrow 1$   
 sign  $\rightarrow - \quad \$\$ \leftarrow -1$   
 list  $\rightarrow$  bit  $\$\$ \leftarrow \$1$

list $\rightarrow$ list bit	$\$\$ = 2 * \$1$
	$+ \$2$
Bit $\rightarrow 0$	$\$\$ \leftarrow 0$
Bit $\rightarrow 1$	$\$\$ \leftarrow 1$

ACTIONS: To perform the actions in the middle of the compiler writer can transform the grammar so that it performs a reduction at each point when required. It performs an abstract syntax tree. The compiler front end must build an IR of the program for the use in middle & back end part of compilation. AST is a common form of tree structured IR. A compiler uses the series of routines to construct AST & one such routine called MakeNode where  $0 \leq i \leq 3$ . The routine takes the 1<sup>st</sup> argument a constant that uniquely identifies a grammar symbol that node will represent the remaining arguments represent a subtree for eg:- MakeNode<sub>0</sub>(Num), constructs a leaf node. If, MakeNode<sub>2</sub>(plus, MakeNode<sub>0</sub>(num), MakeNode<sub>0</sub>(num)). this builds an AST rooted with "+" & with 2 children.

i.e.,  $2+3$



\* pattern  
 DB S  
 T<sub>100</sub>  
 T<sub>200</sub>  
 T<sub>300</sub>  
 T<sub>400</sub>  
 T<sub>500</sub>  
 T<sub>600</sub>  
 T<sub>700</sub>  
 T<sub>800</sub>  
 T<sub>900</sub>  
 min-s

### PRODUCTION

$$E \rightarrow E + T$$

(Graphical form of IR) ASDT  
 ↑ to create  
 ↑ 2 nodes are created  
 $\$\$ \leftarrow \text{make node}_2(\text{plus}, \$1, \$2)$   
 $\$\$.\text{type} \leftarrow F+(\$1.\text{type}, \$2.\text{type})$   
 $\$\$ \leftarrow \text{make node}_2(\text{sub}, \$1, \$2)$   
 $\$\$.\text{type} \leftarrow F-(\$1.\text{type}, \$2.\text{type})$   
 $\$\$ \leftarrow \$2$

$$| E - T$$

| T

$$T \rightarrow T * F$$

"

$\$2, \$3$

$b_1, b_2$

| F

$\$\$ \leftarrow \$3$

$$F \rightarrow (E)$$

$\$\$ \leftarrow (\$1)$

↑ the node is created.

| id

$\$\$ \leftarrow \text{make node}_2(\text{id})$

| num

$\$\$ \leftarrow \text{make node}_2(\text{num})$

### ILOC : Intermediate Language OF Compiler:

- for an optimized compiler

- IR

- uses 4 supporting routines

RENS

Address → takes var names as its arguments.

Emit →

Next register →

value →

to the variable

ADDRESS: Takes 'a' variable name as its argument and returns a register number. (temporary memory in processor). It ensures that the register contains the variable address (memory address)

EMIT: It handles the details of creating a concrete representation for the various ILOC Operations. It can even format and print to the file.

NEXTREGISTER: Returns a new register number  
 simple implementation to increment a global count  
VALUE: Takes the number as its argument and returns the register that contains the numbers passed as its argument. If required, it generates the code to move that number into a register.

### PRODUCTION

$$E \rightarrow E + T$$

$$| E - T$$

$$T \rightarrow T * F$$

$$| T / F$$

$$| F$$

$$F \rightarrow (E)$$

$$| id$$

$$| num$$

$T_{200}$	$I_2 I_3$
$T_{300}$	$I_1 I_2 I_4$
$T$	$I_1 I_2 I_3$

$I_2 I_3$   
 $I_3$   
 $I_2 I_3 I_4$   
 $I_2 I_3 I_5$

$I_2 I_3 I_4$   
 $I_2 I_3 I_5$

$I_2 I_3 I_4 I_5$

## INTERMEDIATE REPRESENTATION:

IR is required to represent the code for further optimization across various phases of the compiler. Hence, IR must be

1. expressive to record all useful facts in a program
2. In expensive
3. concise
4. easily examinable

## CLASSIFICATION:-

IR has a structural organization and can be classified into three.

1. Graphical IR
2. Linear IR
3. Hybrid IR.

GRAPHICAL IR: Encodes the compiler knowledge in the form of graph. The algorithms express in terms of nodes and edges (or) list and trees.

Eg:- parse tree, abstract syntax tree.

LINEAR IR:- resembles pseudo code for some machine languages. The algorithms iterate over simple linear sequences of operators.

Eg :- ILOC, 3 address code

HYBRID IR:- this combines the elements of both linear and graphic IR in an attempt to minimize drawbacks of both the IR's.

## GRAPHICAL IR:-

- consists of nodes & edges
- difference lies in structure & in graph.

# 1. Syntax Related trees

1. parse tree
  2. AST - Abstract Syntax tree
  3. DAG - Directed acyclic graph
- imp
2. Graph
    - \* control flow graph [CFG]
    - \* Dependency graph

## Tree

→ correspond to Syntax of source code

### 1. parse tree

(syntax tree)

corresponds to derivation from i/p program.

Ex: Draw parse tree for  $x * 2 + x * 2 * y$

$$E \rightarrow E + T$$

$$| \quad E - T$$

$$| \quad T$$

$$T \rightarrow T * F$$

$$| \quad T \div F$$

$$| \quad F$$

$$F \rightarrow id$$

$$| \quad num$$

$$| (E)$$



## AST:

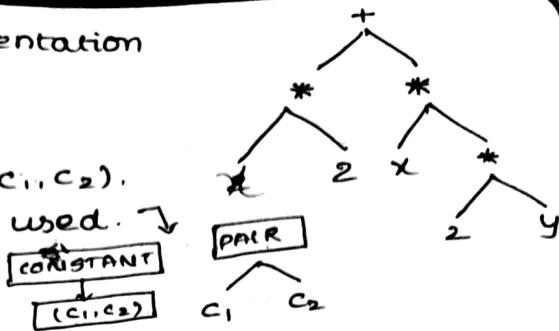
- \* Retains basis of parse tree but eliminates extra nodes
- \* precedence & meaning of expression remains (no extra nodes whereas in parse tree, we have extra nodes)

near source level representation  
used in many compilers

Ex:      in      FORTRAN

a complex constant ( $c_1, c_2$ ).

2 full AST may be used. →



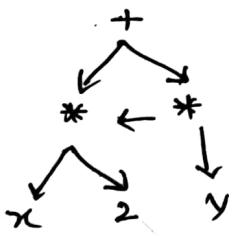
$I_1$ ,  $I_2$   
 $I_1$ ,  $I_3$   
 $I_2$ ,  $I_3$   
 $I_1$ ,  $I_3$   
 $I_1$ ,  $I_2$   
 $I_1$ ,  $I_2$   
 appear -?

DAG:

- contradiction of AST
  - Avoids duplicates
  - In DAG, a node can have multiple parents but identical subtrees reused.  
hence a more compact representation (reused)

$$x * 2 + x * 2 * y$$

DAG



- better code
- But compiler has to prove that  $x * 2$  doesn't change.  
(with directions more importantly)

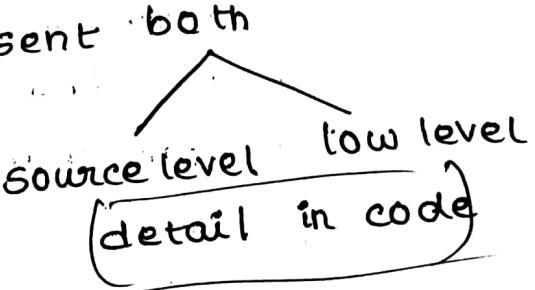
## ADVANTAGES:-

- ADVANTAGES:

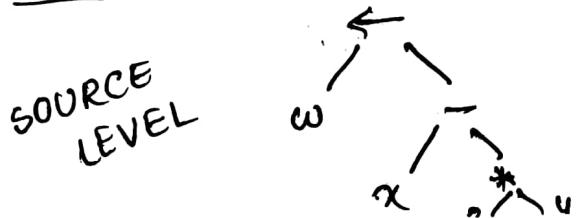
  1. Smaller memory requirements
  2. To expose potential redundancy.
  3. Better for the optimization.

## LEVELS OF DISTRACTION:-

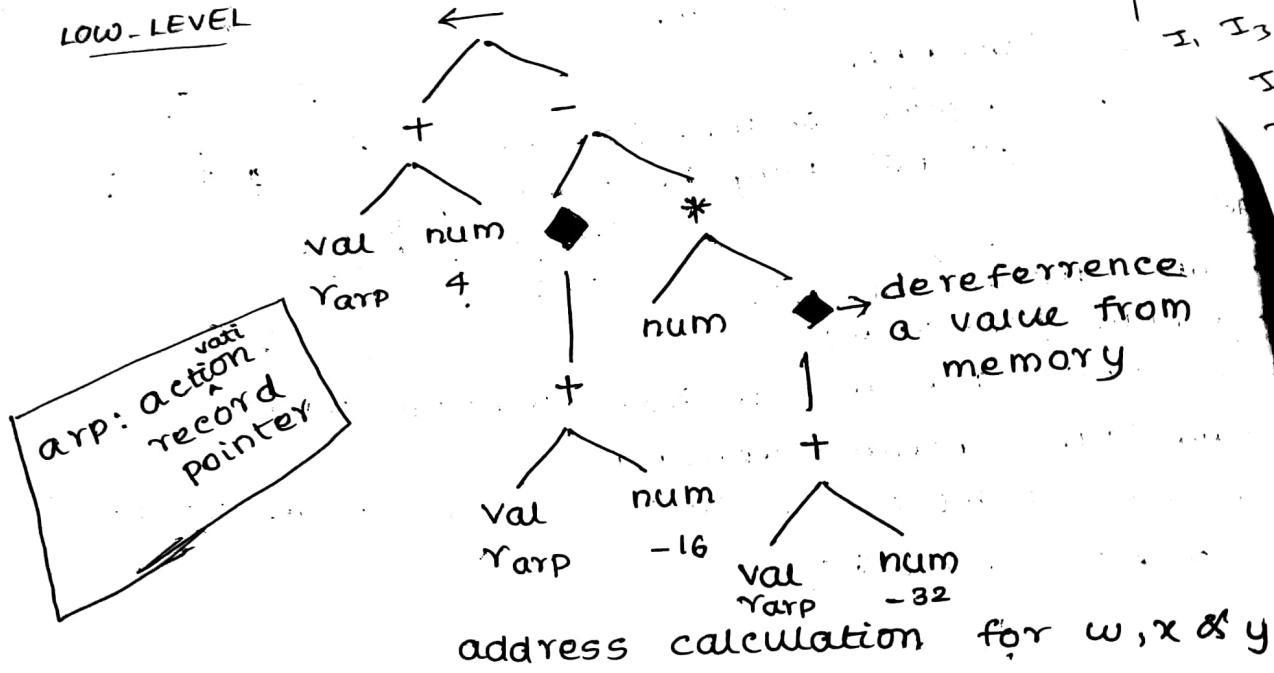
Trees can be used to represent both



e.g.:  $\omega \leftarrow x - 2 * y$



LOW-LEVEL



T<sub>300</sub>  
10

I<sub>1</sub> I<sub>2</sub> I<sub>3</sub>  
I<sub>1</sub> I<sub>2</sub> I<sub>3</sub>  
I<sub>2</sub> I<sub>3</sub>  
I<sub>1</sub> I<sub>3</sub>

I<sub>2</sub> I<sub>3</sub> I<sub>4</sub>  
I<sub>2</sub> I<sub>3</sub>

1. Draw AST for

1.  $2 * 7 + 3$
2.  $2 * (7 + 3)$
3.  $((1 + 2 + 3 + 4) + 5)$
4.  $7 + 3 * (10 | (12 | (3 + 1) - 1))$

5. Low level AST

$$c = a + b * d$$

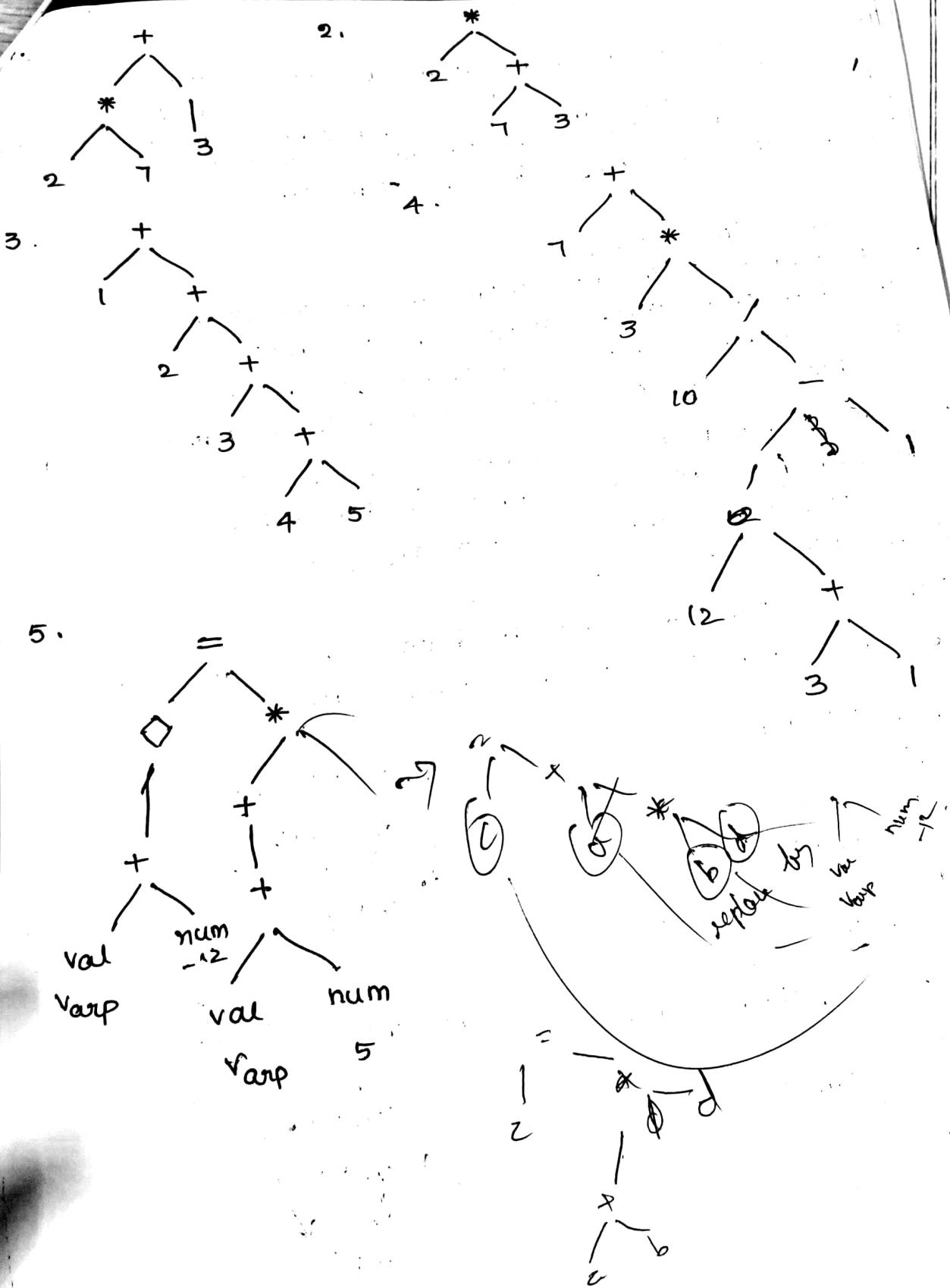
a  $\rightarrow$  5 from Yarp

b  $\rightarrow$  ptr from Yarp @ 6

d  $\rightarrow$  7 from Yarp

c  $\rightarrow$  ptr from Yarp @ -12

	List
T100	I <sub>1</sub>
T200	I <sub>2</sub>
T300	I <sub>2</sub>
T400	I <sub>1</sub>
T500	I <sub>1</sub>
T600	I <sub>2</sub>
T700	I <sub>1</sub>
T800	I <sub>1</sub>
T900	I <sub>1</sub>
	min-support



2020.03.02 14:40

NOTE:  
Trees provide natural representation for the grammatical structure of course code, rigid structure makes them less useful for representing other properties of P. Hence, to model the other properties compiler uses graphs that are more general than trees.

### CFG:

- Control Flow Graph
- It models flow of control in the programs, it is the directed graph given by  $G = (N, E)$  where  $N$  = is a node that corresponds to the each basic block. (set of operations within a loop)
- \* A basic block is a sequence of operations that are always executed together. each edge 'E' corresponds to the possible transfer of control from block  $N_i$  to  $N_j$ .
- Thus, CFG provides the graphical representation of possible runtime control flow path.

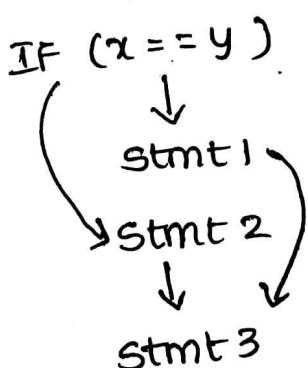
### CODE:

1.  $\text{while } (i < 100)$   
begin  
     $x = x + 10;$   
     $y = y + x;$   
end  
     $z = x + y$

2.  $\text{if } (x == y) \text{ then stmt1}$   
else  
    stmt2  
stmt3

### CFG:

$\text{while } (i < 100)$   
     $x = x + 10;$   
     $y = y + x;$   
     $z = x + y$



A CFG for conditional statement is acyclic & for looping construct, it is cyclic

\* compiler uses CFG along with other IR, such as AST, DAG etc.

- \* CFG may be useful as it supports
1. Code optimization
  2. Instruction scheduling
  3. Global register allocation

### DEPENDANCY GRAPH:

- These are useful for determining an evaluation model for attribute instances in a given parse tree, while unallocated parse trees shows the attribute values. The dependency graph emphasizes the flow of information among the attribute instances. An edge from one attribute instance to another means the value of the 1<sup>st</sup> is required to compute the second.

Edges express the constraints implied by the semantic use. compilers use the Data dependency graph to encode the flow of values, from the point where the values created (definition points) to the points where it is used (used points). Through the dependency graph, the real constraints are represented as sequence that it doesn't capture the program control flow fully.

### pgm:

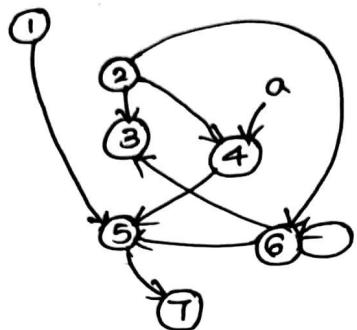
- ①  $x \leftarrow 0$
- ②  $i \leftarrow 1$
- ③ while ( $i < 100$ )
- ④ IF ( $a[i] > 0$ )

- ⑤ then  $x \leftarrow x + a[i]$
- ⑥  $i \leftarrow i + 1$
- ⑦ print(x)

T <sub>100</sub>	I <sub>1</sub> , I <sub>2</sub>	I <sub>5</sub>
T <sub>200</sub>	I <sub>2</sub>	I <sub>4</sub>
T <sub>300</sub>	I <sub>2</sub>	I <sub>3</sub>
	I <sub>1</sub> , I <sub>2</sub>	I <sub>4</sub>
	I <sub>1</sub> , I <sub>3</sub>	
	I <sub>2</sub>	I <sub>3</sub>
	I <sub>3</sub>	
	I <sub>2</sub> , I <sub>3</sub>	I <sub>5</sub>
	I <sub>2</sub>	I <sub>3</sub>

2020-03-02 14:41:00

DEPENDANCY TREE

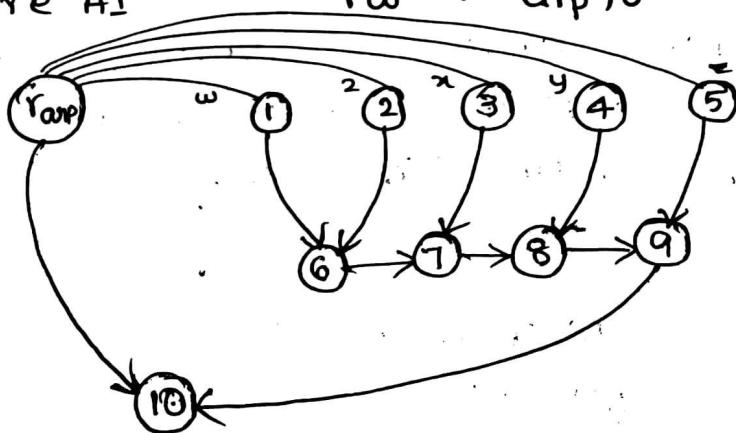


$I_2 \rightarrow I_4$   
 $I_2 \rightarrow I_3$   
 $I_1 \rightarrow I_2 \rightarrow I_4$   
 $I_3$   
 $I_3$   
 $I_3$

2020-03-02 14:40

ILOC BLOCK:

1. LOAD AI       $r_{arp} @ w \Rightarrow r_w$
2. LOAD #1       $2 \Rightarrow r_z$
3. LOAD AI       $r_{arp}, @ x \rightarrow r_x$
4. LOAD AI       $r_{arp}, @ y \rightarrow r_y$
5. LOAD AI       $r_{arp}, @ z \rightarrow r_z$
6. mult       $r_w, r_z \rightarrow r_w$
7. mult       $r_w, r_x \rightarrow r_w$
8. mult       $r_w, r_y \rightarrow r_w$
9. mult       $r_w, r_z \rightarrow r_w$
10. STORE AI       $r_w \rightarrow r_{arp}, 0$



### LINEAR IR:-

- Requirements of instructions that execute in their order of appearance.



- One addresscode  $\xrightarrow{\text{I/P}}$  accumulator
- two addresscode  $\xrightarrow{\text{I/P}}$  obsolete
- three addresscode
- Stack machine code

\* ILOC is called linear IR.

### STACK MACHINE CODE:-

- compact representation (size)
- Operations take operands from stack & push results on to the stack.

creates implicit name space.

ex:  $f(x - 2 * y)$

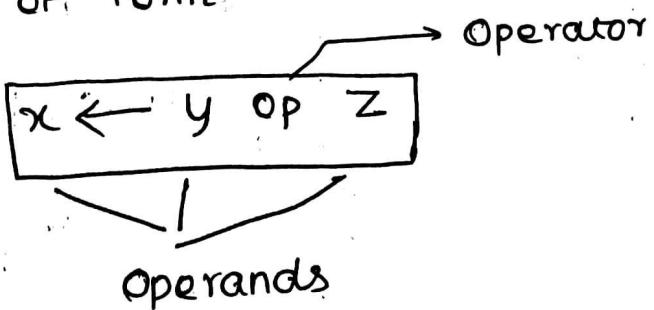
```

LOAD Y      PUSH Y
L.          PUSH 2
            MULTIPLY
            PUSH X
            SUBTRACT
  
```

### TAC:-

### THREE ADDRESSCODE:

Operations of form



### ADVANTAGES:-

- Reuse names & values

- NAMES & OPERATOR  
 ADDRESSING IMPLEMENTS 3 address operator  
 WELL FOR EXECUTION.  
 TO OPTIMIZE AS IT CONTAINS LOW LEVEL OPERATION  
 IMPLEMENTED AS A SET OF QUADRUPLES.

$x = *a$

TAC

$t_1 \leftarrow a$   
 $t_2 \leftarrow y$   
 $t_3 \leftarrow t_1 * t_2$   
 $t_4 \leftarrow x$   
 $t_5 \leftarrow t_4 - t_3$

$x = *a;$   
 $a = \&b;$   
 $t_1 \leftarrow a$   
 $t_2 \leftarrow *t_1$        $t_4 \leftarrow b$   
 $t_3 \leftarrow x$        $t_5 \leftarrow \&t_4$   
 $t_3 \leftarrow t_2$        $t_1 \leftarrow t_5$   
 registers

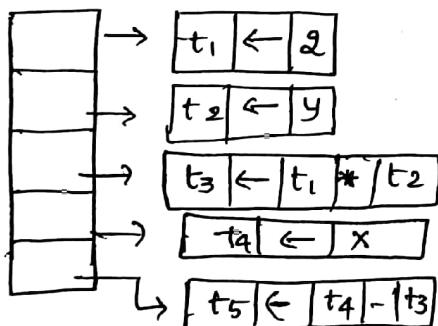
### TAC IMPLEMENTATION:

#### 1. Simple array

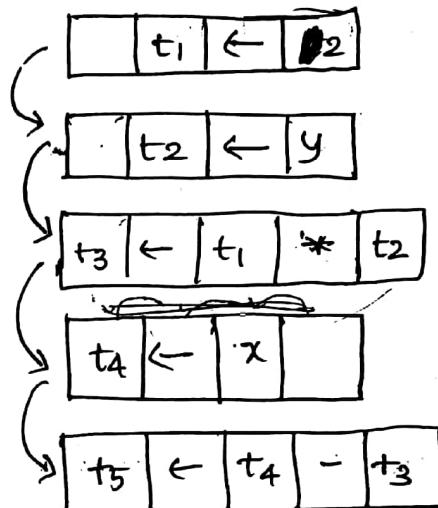
target op arg1 arg2

$t_1$	$\leftarrow$	2	} most of space is not utilised & can overcome by following methods
$t_2$	$\leftarrow$	y	
$t_3$	*	$t_1$	
$t_4$	$\leftarrow$	x	
$t_5$	-	$t_4$	$t_3$

#### 2. Array of pointer



#### 3. Linked List (usually followed)



\* compact 3 names & 1 operator  
 \* modern processors implement 3 address operators  
 model - Fixes well for execution.

\* easier to optimize as it contains low level operations  
 \* represented as a set of quadruples.

$$x = 2 * y$$

TAC:

$$\begin{aligned} t_1 &\leftarrow a \\ t_2 &\leftarrow y \\ t_3 &\leftarrow t_1 * t_2 \\ t_4 &\leftarrow x \\ t_5 &\leftarrow t_4 - t_3 \end{aligned}$$

$$\begin{aligned} x &= * a; \\ a &= \& b; \end{aligned}$$

$$\begin{aligned} t_1 &\leftarrow a \\ t_2 &\leftarrow * t_1 \\ t_3 &\leftarrow x \\ t_3 &\leftarrow t_2 \\ &\text{Registers} \end{aligned}$$

$$\begin{aligned} t_4 &\leftarrow b \\ t_5 &\leftarrow \& t_4 \\ t_1 &\leftarrow t_5 \end{aligned}$$

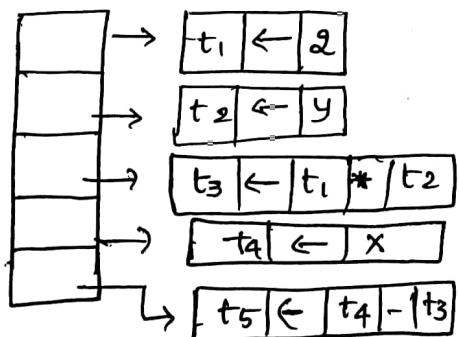
TAC IMPLEMENTATION:

1. Simple array

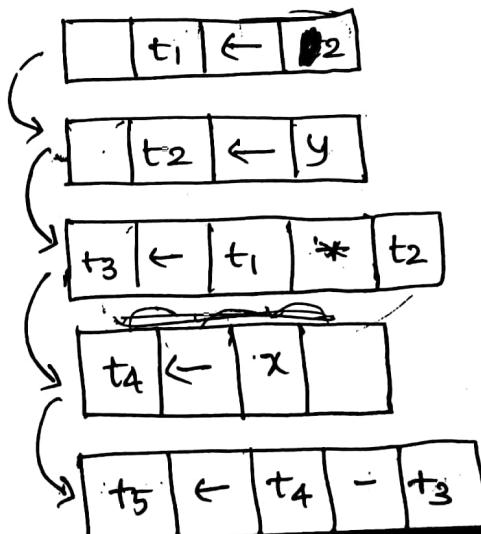
target op arg1 arg2

$$\begin{array}{llll} t_1 & \leftarrow & 2 & \} \text{ most of space} \\ t_2 & \leftarrow & y & \text{is not utilised} \\ t_3 & * & t_1 & \& \text{can overcome} \\ & & & \text{by following} \\ t_4 & \leftarrow & x & \text{methods} \\ t_5 & - & t_4 & t_3 \end{array}$$

2. Array of pointer



3. Linked List (usually followed)



2020.03.02 14:41

T100

T200

T300

T400

T500

T600

T700

T800

T900

T1000

### MAPPING VALUES TO NAMES:

#### Source code

$$\begin{aligned} a &\leftarrow b+c \\ b &\leftarrow a-d \\ c &\leftarrow b+c \\ d &\leftarrow a-d \end{aligned}$$

#### 2 methods:

1. Using source names
2. Using value names

2MQ

#### MEMORY MODELS:

1. Register to Register
2. Memory to Memory

#### SINGLE STATIC ASSIGNMENT FORM (SSA FORM):

In this form, values are assigned a static name, both data flow and control flow can be represented by  $\phi$  (Psi) function. This function corresponds uniquely to specific determined points in the code as each name has one operation. It is called as SINGLE STATIC ASSIGNMENT FORM.

#### CODE

$$\begin{aligned} x &\leftarrow \\ y &\leftarrow \\ \text{while } (x < 100) \\ &\quad x \leftarrow x+1 \\ &\quad y \leftarrow y+x \end{aligned}$$

#### SSA FORM

$$\begin{aligned} x_0 &\leftarrow \\ y_0 &\leftarrow \\ \text{if } (x_0 \geq 100) \text{ go to next} \\ \text{loop: } x_1 &\leftarrow \phi(x_0, x_2) \\ y_1 &\leftarrow \phi(y_0, y_2) \\ x_2 &\leftarrow x_1 + 1 \\ y_2 &\leftarrow y_1 + x_1 \\ \text{if } (x_2 < 100) \text{ go to loop} \\ \text{next: } x_3 &\leftarrow \phi(x_0, x_2) \\ y_3 &\leftarrow \phi(y_0, y_2) \end{aligned}$$

- \* Two considerations for SSA, each definition is distinguishable.
- 2. each use refers to single definition so, compiler inserts a  $\phi$  function. when different controller merge & it renames variables to create SSA.  
 $\phi$  functions can be inserted where multiple distinct values reach for the block

ex:-

$$\begin{aligned} y &= x - a \\ k &= y - t \\ y &= m - k \\ k &= y - k \end{aligned}$$

$$\begin{aligned} y_1 &= x_1 - a, \\ k_1 &= y_1 - t, \\ y_2 &= m_1 - k_1, \\ k_2 &= y_2 - k_1 \end{aligned}$$

imp

SYMBOL TABLE: - It is central repository of the facts to provide efficient access to it, compiler discovers and stores a wide variety of information from the program as given below.

1. Variables - stores datatype, name & scope
2. arrays - Name, dimension ; datatype, lower bound & upperbound.
3. functions & procedures - Argument list, function name, return type
4. records & structures - Fields & associated types

USES:

- provides centralized uniform access with the single implementation
- localizes information from the source code.
- avoids the expense of searching from IR to find fact

### IMPLEMENTATION: Better method - HashTable.

- HASH TABLE: The key issue efficiency as a compiler access symbol table very frequently. It provides a constant time for storing and retrieval using hash.

- The hash function maps the facts to small integers and these integers are used to index the table, with hash-symbol table; the compiler stores all information that it derives about the name 'n' into the table slot  $H(n)$ .

- Building up a single table uses two routines,

1. LOOKUP(name): returns the record stored at  $H(name)$ , if it exists and it returns a value indicating that the name was found (or) not

2. INSERT(name, record): stores the information in the hash table @ $H(name)$ , it may expand the table to accommodate new class

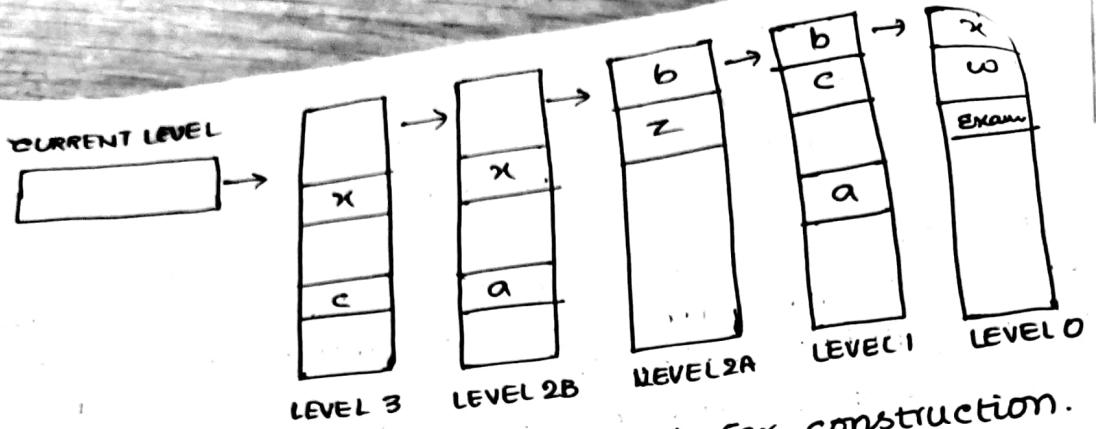
### SCOPE (LEVEL)

```

static int w; } # LEVEL 0
int x;
void example (int a, int b) } # LEVEL 1
int c;
{
    int b,z; } # LEVEL 2a
    :
}
{
    int a,x; } # LEVEL 2b
    :
}
int c,x;
b=a+b+c+w; } # LEVEL 3

```

$$\rightarrow \Psi(y_0, y_2)$$



Two more routines are used for construction.

1. Initialize scope
2. Finalize scope

#### MANAGEMENT OF SPACE:-

- \* Separate table for each record
- \* Selector table for field names that are identified.
- \* Unified table

#### NOTE:-

MEMO FUNCTION:- A function that avoid recompilation of prior results by storing in a hash table under a key built with its arguments.

PURE FUNCTION:- A function always returning the same results.

to optimize code compiler uses various strategies & knowledge.

### 3. OPTIMIZATION

#### STORING & ACCESSING ARRAYS:

- collection of elements of same data type.
- requirements coniguous memory location
- One dimension
  - ↳ vector (6)

1	2	4	5	6	7
0	1	2	3	4	5

- array element required referred by address/index
- third element  $v[2] = 4$ .
- requirement offset of  $i^{\text{th}}$  element

#### FORMULA:

$$(i - \text{low}) * w \rightarrow \text{length in bytes}$$

for int,  
 $\text{low} = 0$   
 $w = 2$

offset of 3

$$(3 - 0) * 2 = 6$$

#### 3 ADDRESS CODE FOR ACCESSING ELEMENT:

load I  $@v \rightarrow y_0$

sub I  $y_{i,3} \rightarrow y_1$

Mult  $y_1, 2 \rightarrow y_2$

add  $y_0, y_2 \rightarrow y_3$

load  $y_3 \rightarrow y_4$

#### MULTIDIMENSION ARRAY:

##### 3 schemes

1. Row major order
2. Column major order
3. Indirection vectors

## ROW MAJOR ORDER:

- adjacent elements of single row occupies contiguous location.

Eg:

1,1	1,2	1,3	1,4
2,1	2,2	2,3	2,4

row major order is given by

1,1	1,2	1,3	1,4	2,1	2,2	2,3	2,4
-----	-----	-----	-----	-----	-----	-----	-----

loop for above,

for  $i \leftarrow 1$  to 2

    for  $j \leftarrow 1$  to 4

$A[i,j] \leftarrow A[i,j] + 1$

## ADDRESS CALCULATION:

$$A[1, \dots, 2, 1, \dots, 4]$$

where  $low_1 = 1$  &  $high_2 = 4$

$A[i,j]$

row i & follow that offset

for element j

each row 4 element

$$len_k = high_k - low_k + 1$$

$$= 4 - 1 + 1 = 4$$

follow base address  $A[2,3]$   
(<sup>row index</sup>)

$$(i - low_1) * len * w \rightarrow \text{float}$$

$$(2 - 1) * 4 * 4 = 16.$$

Address computation ( $A[2,3]$ )

$$\text{row} = @A + 16$$

ACOL row

$$@A + 16 + (j - low_2) * w$$

$$\Rightarrow @A + 16 + (3 - 1) * 4 \Rightarrow @A + 16 + 8 \Rightarrow @A + 24.$$

1,1	1,2	1,3	1,4	2,1	2,2	2,3	2,4
0	4	8	12	16	20	24	

COLUMN MAJOR ORDER:

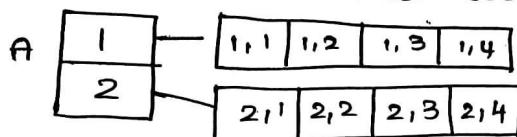
→ elements of array mapped into adjacent locations of single column.

1,1	2,1	1,2	2,2	1,3	2,3	1,4	2,4
-----	-----	-----	-----	-----	-----	-----	-----

→ non sequential access.

### INDIRECTION VECTORS:

→ multi dimensional arrays into vectors.



elements address as vectors.

1. extra storage

2. Initialization code representation

(Forecasting)

CODE OPTIMIZATION: It is to discover information about the runtime behaviours at compile time.

Two types of compilers:

1. Debugging compilers
2. Optimizing compilers

### OPTIMIZING CONSIDERATION:

1. Results should be same / safety concern (meaning shouldn't change)

2. profitability (NO:of iterations in pgm. should be reduced) & (Removing duplicates)

3. ~~Risk~~ Identifies RISK

[Registered Demand should be minimum  
address calculations]

### OPPORTUNITY:

1. Reduce the override of abstraction.

eg: array address calculations.

2. matching the code to the system resources decrease memory allocation.
  3. Takes advantage of special functions.

### SCOPE OF OPTIMIZATION OR CATEGORIES:

1. Local methods
    - Basic blocks
  2. Regional methods
    - Extended basic blocks larger than single block
  3. Intraprocedural optimization
    - procedures
  4. Interprocedural optimization
    - whole program

## LOCAL METHODS:-

- Optimization within basic blocks
  - local value numbering [LVN]
    - ↳ algorithm uses hashing to map names, constants & expressions into distinct value numbers.
  - IF "hash key" already found then
    - replace
  - else
    - insert a new value

$$\begin{array}{l}
 \text{Ex: } * \\
 \begin{array}{rcl}
 a & \leftarrow & b + c \\
 & & \begin{array}{c} 1 \\ 3 \\ 5 \end{array} \\
 & & \begin{array}{c} 2 \\ 4 \end{array} \\
 b & \leftarrow & a - d \\
 & & \begin{array}{c} 3 \\ 5 \end{array} \\
 & & \begin{array}{c} 2 \\ 4 \end{array} \\
 c & \leftarrow & b + c \\
 & & \begin{array}{c} 6 \\ 5 \end{array} \\
 & & \begin{array}{c} 3 \\ 4 \end{array} \\
 d & \leftarrow & a - d \\
 & & \begin{array}{c} 7 \\ 3 \end{array}
 \end{array}
 \end{array}$$

REWRITTEN  
CODE:- \*

$$\begin{aligned} a &\leftarrow b + c \\ b &\leftarrow a - d \\ c &\leftarrow b + c \\ d &\leftarrow b \end{aligned}$$

$$\begin{aligned} a^3 &\leftarrow x+y^2 \\ b^4 &\leftarrow x+y^2 \\ a^5 &\leftarrow 17 \\ c^6 &\leftarrow x+y^2 \end{aligned}$$

$$\begin{array}{l} a^{\bullet} \leftarrow x + y \\ b^{\bullet} \leftarrow a \\ a \leftarrow 17 \\ c \leftarrow b \end{array}$$

# LOCAL NUMBERING TECHNIQUE //

LINE SUBSTITUTION :

→ code substituted

```
inline void sum (int a, int b)
{
    cout << a+b;
}
```

3. CONSTANT FOLDING:

$$\textcircled{1} \quad x = 3$$

$$y = x+2$$

$$\textcircled{2} \quad \text{return } (2+3)$$

optimized as  
return (5);

# compiler only  
optimizes the  
code.

Folding:

$$x = 3$$

$$y = 5$$

EXTENSIONS: Algebraic identities  $\Rightarrow$  for constant folding

$$1. \quad a * b = b * a$$

$$2. \quad a + b = b + a$$

$$3. \quad A - 0 = A \leftarrow \text{(constant folding)}$$

$$4. \quad A - A = 0$$

$$5. \quad a * 1 = a$$

$$6. \quad a * 0 = 0$$

$$7. \quad a \text{ and } a = a$$

$$8. \quad a \text{ or } a = a$$

ILOCAL METHODS:

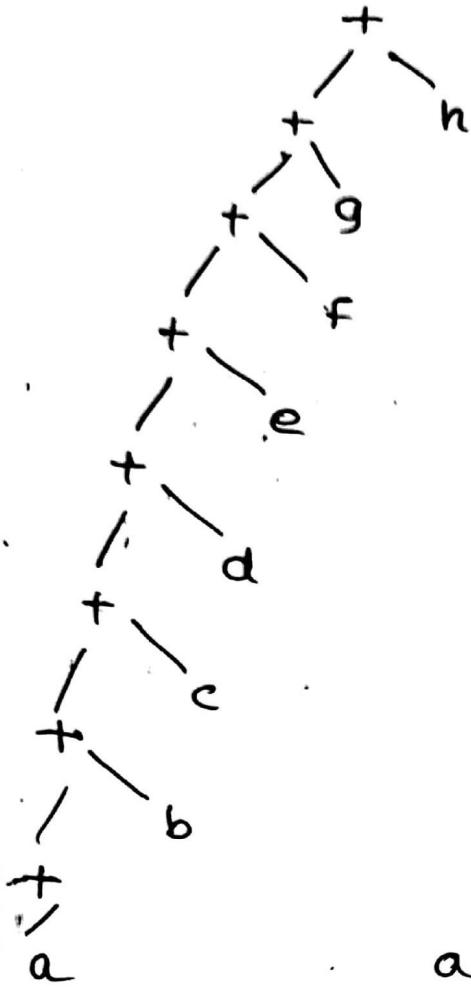
1. LVN

2. Tree height balancing

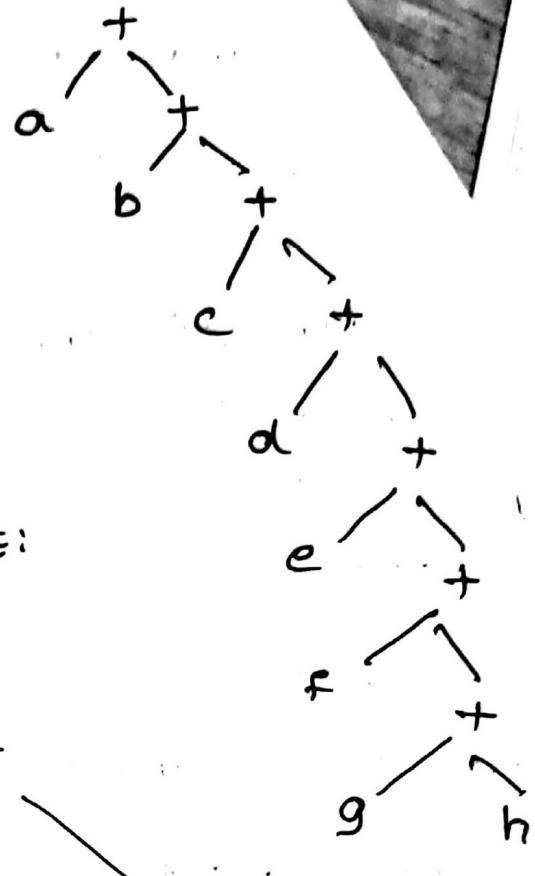
TREE HEIGHT BALANCING:

Eg:- a + b + c + d + e + f + g + h.

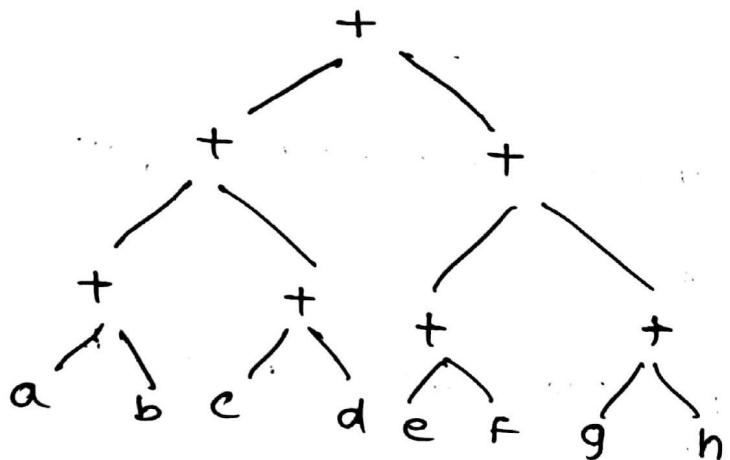
### Left associative tree



### right associate



### BALANCED TREE:



### Balanced tree

#### unit-0

- 1)  $t_1 \leftarrow a+b$
- 2)  $t_3 \leftarrow e+f$
- 3)  $t_5 \leftarrow t_1+t_2$
- 4)  $t_7 \leftarrow t_5+t_6$

#### unit-1 $\rightarrow$ processors

- 1)  $t_2 \leftarrow c+d$
- 2)  $t_4 \leftarrow g+h$
- 3)  $t_6 \leftarrow t_3+t_4$

4) units both processors used

### left association

#### unit 0

- 1)  $t_1 \leftarrow a+b$
- 2)  $t_2 \leftarrow t_1+c$
- 3)  $t_3 \leftarrow t_2+d$
- 4)  $t_4 \leftarrow t_3+e$
- 5)  $t_5 \leftarrow t_4+f$
- 6)  $t_6 \leftarrow t_5+g$
- 7)  $t_7 \leftarrow t_6+h$

#### unit 1

7 units  
1 processor idle

- Compiler identifies candidate expression in blocks. All the expression in the candidate tree must be commutative and associative. Identitative each name that labels the interior node must be used exactly once.
- \* For each candidate tree, the algorithm finds all its operands assigns each a rank and enters them into a priority queue ordered in ascending order. From this queue, the algorithm builds a balance tree.

\* Several factors for identifying the candidate tree

1. Trees cannot be larger than the block
2. The rewritten code cannot change the observable values
3. Trees cannot extend backwards.

- The algorithm uses the following methods for the tree finding phase:
- \* The set  $\text{uses}(T_i)$  contains the index in the block for each use of  $T_i$ . If  $T_i$  is used after the block then  $\text{uses}(T_i)$  should contain '2' additional entries.

i.e., Integer > no:of operations in the block.

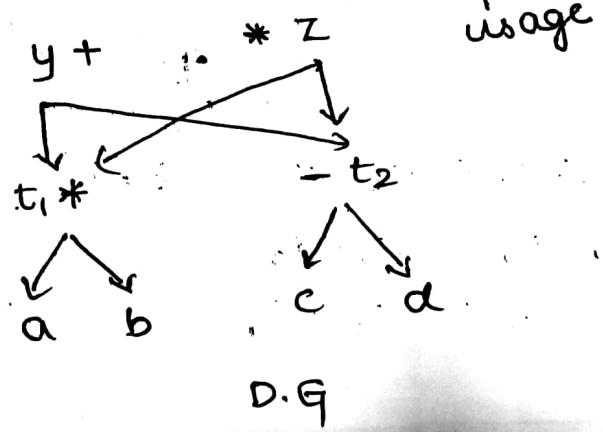
### TREE MHT BALANCING :-

- \* The algorithm constructs maximum signed candidate tree using DG (Dependancy Graph)

Eg:

$$\begin{aligned} t_1 &\leftarrow a * b \\ t_2 &\leftarrow c + d \\ y &\leftarrow t_1 + t_2 \\ z &\leftarrow t_1 * t_2 \end{aligned}$$

code



2020.03.02 14:42

## SEVERAL FACTORS FOR TREE ID:

1. trees cannot be larger than block.
2. rewritten code cannot change its value
3. trees cannot extend backward.

## ALGORITHM TREE FOR FINDING PHASE:

Each  $T_i$  def in the block.

- (i)  $uses(T_i)$  contains the index in the block. For each use ( $T_i$ ).

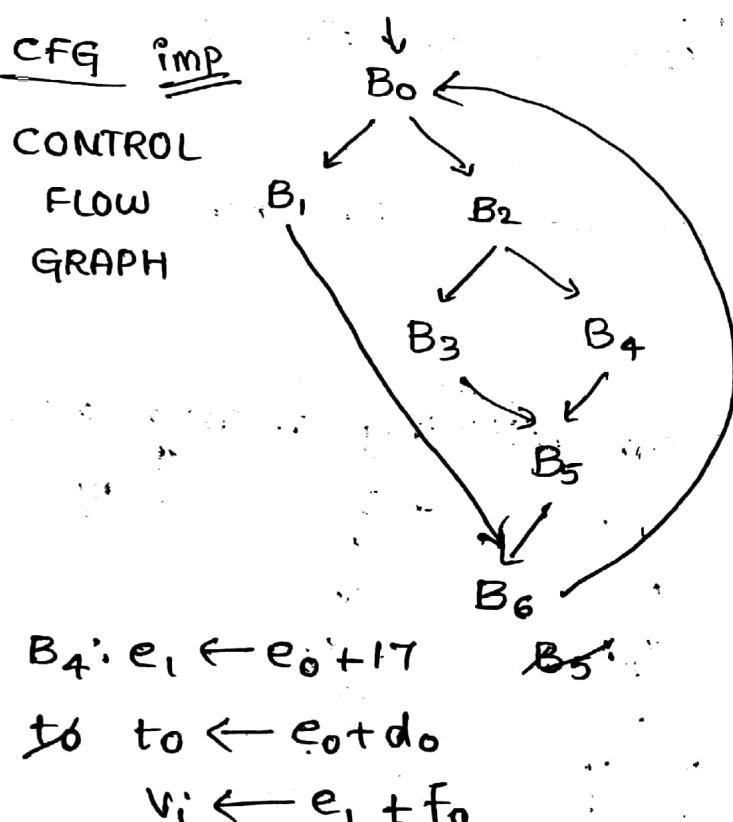
If  $T_i$  used after  $\downarrow$  should contain 2 additional entries greater than no: of operations in the block.

## II REGIONAL OPTIMIZATION:

Larger than single block.

points 2 methods

- ① Super local value no (larger region) → used for redundant equation elimination in more than one block.
- ② loop unrolling.



$B_5: e_2 \leftarrow \phi(e_0, e_1)$   
 $v_2 \leftarrow \phi(v_0, v_1)$   
 $v_0 \leftarrow a_0 + b_0$   
 $w_0 \leftarrow l_0 + d_0$   
 $x_0 \leftarrow e_2 + f_0$   
 $\rightarrow B_6$

$B_6: r_2 \leftarrow \phi(r_0, r_1)$   
 $y_0 \leftarrow a_0 + b_0$   
 $z_0 \leftarrow c_0 + d_0$   
 $\rightarrow B_0$

# for loop functions we use the  $\phi$  fn's in code loop

### Algorithm include

1.  $B_5, B_6 \rightarrow$  not important
  2.  $B_0, B_1, B_2, B_3, B_4 \rightarrow$  important
- In  $B_0$ :  ~~$r_0$~~  is redundant  
 $B_1$ :  $r_0$  is redundant  
 $B_2$ :  $q_0, r_1$  is redundant

## 2) LOOP UNROLLING

→ unroll inner loop → "loop fusion".  
 → unroll outer loop → "Jamming".

### UNROLL INNER LOOP:

```

do 60 i=1, n2
    do 50 i=1, n1
        y[i] = y[i] + x[j] * m(i,j)
    50 continue
60 continue
    
```

### UNROLL BY 3:

```

do 60 i=1, n2
    do 50 i=1, n3
        nex = i mod 3
        IF (nex .ge. 1) then
            do 49 i=1, nex
                y[i] = y[i] + x[j] * m(i,j)
            49 continue
        do 50 i=nex+1, n, 3
            y[i] = y[i] + x[j] + m(i,j)
            y[i+1] = y[i+1] + x[j] + m(i,j)
        50 continue
    50 continue
60 continue
    
```

$$y[i+2] = y[i+2] + x[j] * m[i+2, j]$$

50 continue

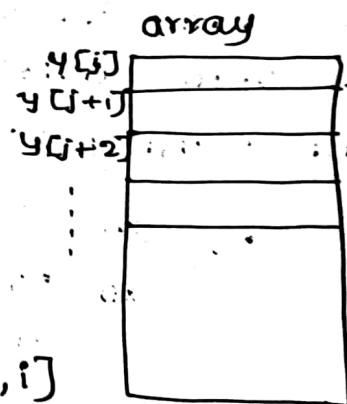
60 continue

#### EQUIVALENT CODE:

```
for (i=1; i<n2; i++)
{
    for (j=1; j<i; j++)
    {
        y[j] = y[j] + x[i] * m[j, i]
    }
}
```

#### UNROLL BY 3:

```
for (i=1; i<n2; i++)
{
    for (j=1; j<i; j+=3)
    {
        y[j] = y[j] + x[i] + m[j, i]
        y[j+1] = y[j+1] + x[i] + m[j+1, i]
        y[j+2] = y[j+2] + x[i] + m[j+2, i]
    }
}
```



→ Reduce test  
& branch  
sequence...

#### LOOP UNROLLING:

##### Outer loop(fusion)

```
eg: DO I=1,N
    DO J=1,N
        DO K=1,N
            A[i,j] = A[i,j] + B[i,k] * C[k,j]
        end DO
    end DO
end DO.
```

ROLL OUTER LOOP BY THE FACTOR OF 2:

```
I=1,N,2  
DO J=1,N  
DO K=1,N  
A[I,J] = A[I,J]*B.....  
end DO  
end DO  
DO j=1,N  
DO K=1,N  
A[i+j] = A[i+j] + B[i+j,K] * C[K,j]  
end DO  
end DO
```

USING INNERLOOP (JAMMING) UNROLL BY A FACTOR OF 3:-

### GLOBAL OPTIMIZATION:

- an entire procedure
- As loops are included requires global analysis before modifying.

STEPS:  
1. finding uninitialized var with live information  
2. Global code placement.

FINDING UNINITIALIZED VARIABLE WITH LIVE INFORMATION:-

In a procedure B until a variable v is assigned a value, it is uninitialized.

uninitialized variable indicates a logical error. Hence, a programmer should be informed about its existence.

uninitialized variables can be identified by using liveliness.

LIVEVAR  $use(v)$   
 A variable 'v' is live at a point P iff there  
 exist a path in CFG from P to  $use(v)$  along  
 which v is not redefined.

We encode a set of liveout(B) that contains all  
 the variables that are live on exit from B.

$$\text{liveout}(n) = \bigcup_{m \in \text{succ}(n)} (v \in \text{VAR}(m) \cup (\text{liveout}(m) \cap \text{VARKILL}(m)))$$

$(v \in \text{VAR}) \Rightarrow$  Upward exposed variable in m, i.e., those  
 $(v \in \text{VAR}) \rightarrow$  variables that are used in the 'm' before  
 any redefinition in m.

$\text{VARKILL}(m) \Rightarrow$  set of variables not defined in m.

To complete the liveout sets the compiler uses 3  
 step algorithm:

- (1) Build a CFG (Control Flow Graph)
- (2) Gather initial information by computing UEVAR and  
 VARKILL set for each block B.
- (3) Solve the equation for  $\text{LIVEOUT}(B)$  for each block.

### ITERATIVE LIVE ANALYSIS:

// block b & operations  $x \leftarrow y \text{ op } z$

for each block b

init(b)

init(b)

$\text{UEVAR}(b) \leftarrow \emptyset$

$\text{VARKILL}(b) \leftarrow \emptyset$

for  $i \leftarrow 1 \text{ to } K$

if  $y \notin \text{VARKILL}(b)$

then add y to  $\text{UEVAR}(b)$

LIVEVAR: A variable 'v' is live at a point  $P$  iff there exist a path in CFG from  $P$  to  $use(v)$  along which  $v$  is not redefined.

We encode a set of liveout(b) that contains the variables that are live on exit from  $B$ :

$$\text{liveout}(n) = \bigcup_{m \in \text{succ}(n)} (v \in \text{VAR}(m)) \cup (\text{liveout}(m) \cap \overline{\text{VARKILL}}(m))$$

(UEVAR)  $\Rightarrow$  Upward exposed variable in  $m$ , i.e., those variables that are used in the ' $m$ ' before any redefinition in  $m$ .

VARKILL(m)  $\Rightarrow$  set of variables not defined in  $m$ .

To complete the liveout sets the compiler uses 3 step algorithm:

- (1) Build a CFG (control flow graph)
- (2) Gather initial information by computing UEVAR and VARKILL set for each block  $B$ .
- (3) Solve the equation for LIVEOUT(B) for each block.

### ITERATIVE LIVE ANALYSIS:-

// block b & operations  $x \leftarrow y \text{ op. } z$

for each block  $b$

init( $b$ )

init( $b$ )

    UEVAR( $b$ )  $\leftarrow \emptyset$

    VARKILL( $b$ )  $\leftarrow \emptyset$

    for  $i \leftarrow 1$  to  $K$

        if  $y \notin \text{VARKILL}(b)$

            then add  $y$  to UEVAR( $b$ )

if  $z \notin \text{VARKILL}(b)$   
 then add  $z$  to  $\text{UEVAR}(b)$   
 Add  $x$  to  $\text{VARKILL}(b)$

### a) GATHERING INITIAL INFORMATION

#### SOLVING THE INITIAL EQUATION :-

ii. CFG with  $V$  blocks  $0$  to  $n-1$ .

for  $i \leftarrow 0$  to  $n-1$

$\text{liveout}(i) \leftarrow \emptyset$

$\text{changed} \leftarrow \text{true}$

while ( $\text{changed}$ )

$\text{changed} \leftarrow \text{false}$

for  $i \leftarrow 0$  to  $n-1$

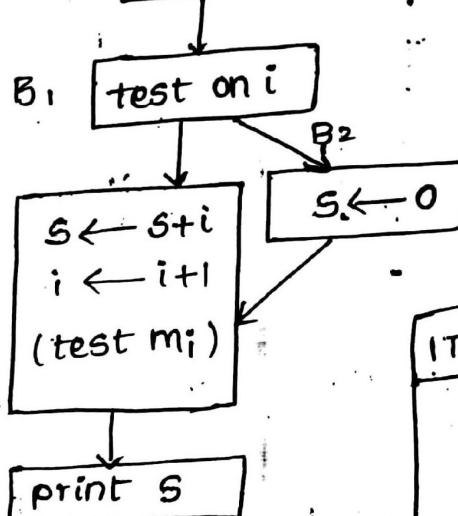
recompute  $\text{LIVEOUT}(i)$

IF ( $\text{liveout}(i)$  changed) then

$\text{changed} \leftarrow \text{true}$ . Initial var.

UEVAR - used  
 VARKILL - assigned

eg:  $B_0$   $i \leftarrow 1$



BLOCK	UEVAR	VARKILL
$B_0$	$\emptyset$	$i$
$B_1$	$i$	$\emptyset$
$B_2$	$\emptyset$	$S$
$B_3$	$S, i$	$S, i$
$B_4$	$S$	$\emptyset$

ITERATION	$B_0$	$B_1$	$B_2$	$B_3$	$B_4$
0	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
1	$i$	$i, S$	$S, i$	$S, i$	$\emptyset$
2	$S, i$	$i, S$	$i, S$	$i, S$	$\emptyset$
3	$S, i$	$i, S$	$i, S$	$i, S$	$\emptyset$

COMPARE  
SUCCESSOR BLOCK

LIVEVAR USES:-

1. used for global register allocation
2. improve SSA construction
3. can discover useless store operation.

GLOBAL OPTIMIZATION:

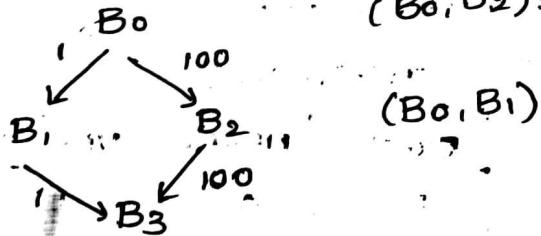
2 techniques

- finally uninitialized var using Liveout
- Global code placement.
- Branching cost considered
- Fall through branch - either taken falls through or next instruction.
- taken.

→ each branch has either falls through (or) taken.

→ compiler can choose which block lives on it.

eg:

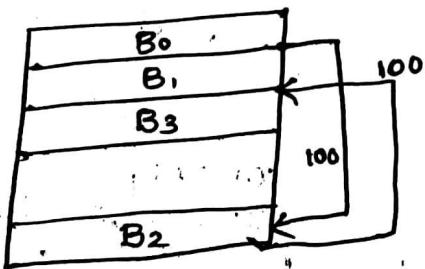


$(B_0, B_2)$  → 100 times more often than

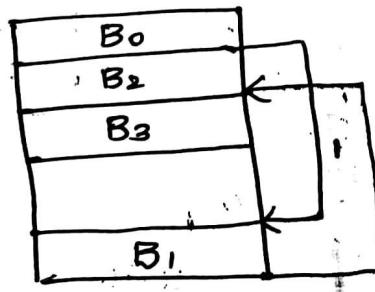
$(B_0, B_1)$

Two layouts:

## 1. slow layout



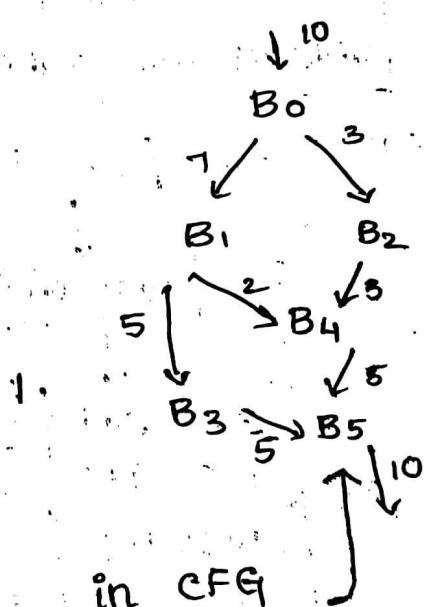
## Fast layout

Placement of code:

- constructing chains as hotpaths in CFG
  - layout of code
- "COMPILER USES CFG PATHS THAT ARE MOST FREQUENTLY EXECUTED EDGES CALLED HOTPATHS".

$x = y + 2 \rightarrow x = 4 \text{ or } 2$   
 $y = 1 \text{ or } 2$   
 $z = 2 \text{ or } 1$

1. CRG
2. id VEVAL & VARKIL
3. compute Liveout formula



using hot paths:

- edges  
 for each block b  
 make generate chain d for b  
 $P \leftarrow \emptyset$   
 for each CFG edge  $\langle x, y \rangle$ ,  $x \neq y$ ,  
 IF 'x' is tail of chain a &  
 'y' is head of chain b then  
 $t \leftarrow \text{priority}(a)$   
 append b to a  
 $\text{priority}(a) \leftarrow \min(t, \text{priority}(b), P++)$   
 GLOBAL OPTIMIZATION

$(B_0, B_1)$	$(B_3, B_5)$	$(B_4, B_5)$	$(B_1, B_2)$
$(B_0, B_2)$	$(B_2, B_4)$	$(B_1, B_4)$	
$B$	$B$	$B$	$P$
$(B_0)_e (B_1)_e (B_3)_e (B_5)_e (B_4)_e (B_2)_e (B_1)_e (B_2)_e$			

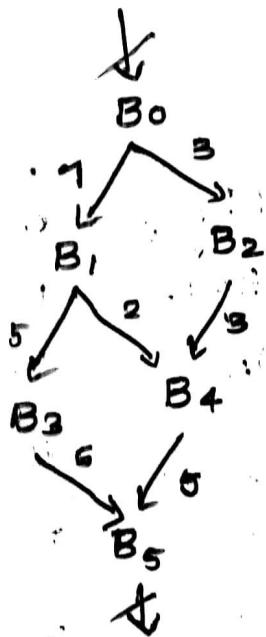
GLOBAL OPTIMIZATION: -

- GLOBAL OPTIMIZATION:-

2 tools → Finding uninitialized using liveout  
→ code placement (to minimize jumps)  
using notpaths

(most frequent executable blocks)

- 1) degenerate block
  - 2) find chains from decreasing order.



edges	chains	P
$\{B_0, B_1\}$	$B_0, B_1, B_2, B_3, B_4, B_5$	0
	$(B_0B_1B_3B_5)_6 (B_2B_4)_5$	

## CODE LAYOUT ALGORITHM

CODE LAYOUT ALGORITHM:  
chain headed by CFG entry node.

$t \leftarrow \text{chain header}, \text{arity}(t)\}$

$t \leftarrow \text{chain head}$   
 $\text{worklist} \leftarrow \{ t, \text{priority}(t) \}$

while (worklist ≠ 0)

remove a chain  $C \in t$   
lowest priority from worklist  
for each block  $x_i$  in  $C$  in chain order  
place  $x_i$  at the end of execu. code  
for each edge  $(x, y)$  where  $y$  is unplaced.

$t \leftarrow$  chain containing  $(x, y)$   
IF  $(t, \text{priority}(t)) \notin \text{worklist}$  then  
 $\text{worklist} \leftarrow \text{worklist} \cup t, \text{priority}(t)$   
(utilizes branching controls)

### CODE PLACEMENT

worklist      code layout

Step

$\{B_0, B_1, B_3, B_5\}_0$

-

$\{B_2, B_4, y_3\}_3$

$B_0, B_1, B_3, B_5$

2

$B_0, B_1, B_3, B_5, B_2, B_4$

$\emptyset$

\* Edges

chains

$(B_0)_E \quad (B_1)_E \quad (B_2)_E \quad (B_3)_E \quad (B_4)_E \quad (B_5)_E$

$(B_0, B_1)_0$

$(B_0, B_1)_0 \quad (B_2)_E \quad (B_3)_E \quad (B_4)_E \quad (B_5)_E$

$(B_3, B_5)_0$

$(B_0, B_1)_0 \quad (B_2)_E \quad (B_3, B_5), (B_4)_E$

$(B_1, B_3)_0$

$(B_0, B_1, B_3, B_5)_0 \quad (B_2)_E \quad (B_4)_E$

$(B_4, B_5)_0$

$(B_0, B_1, B_3, B_5)_0 \quad (B_2)_E \quad (B_4)_E$

$(B_0, B_2)_0$

$(B_0, B_1, B_3, B_5)_0 \quad (B_2)_E \quad (B_4)_E$

$(B_1, B_4)_0$

$(B_0, B_1, B_3, B_5)_0 \quad (B_2)_E \quad (B_4)_E$

$(B_2, B_4)_0$

$(B_0, B_1, B_3, B_5)_0 \quad (B_2, B_4)_0$

Edges(B<sub>0</sub>, B<sub>1</sub>)(B<sub>4</sub>, B<sub>5</sub>)(B<sub>1</sub>, B<sub>3</sub>)(B<sub>3</sub>, B<sub>5</sub>)(B<sub>0</sub>, B<sub>2</sub>)(B<sub>2</sub>, B<sub>4</sub>)(B<sub>1</sub>, B<sub>4</sub>)chains(B<sub>0</sub>)<sub>E</sub> (B<sub>1</sub>)<sub>E</sub> (B<sub>2</sub>)<sub>E</sub> (B<sub>3</sub>)<sub>E</sub> (B<sub>4</sub>)<sub>E</sub> (B<sub>5</sub>)<sub>E</sub> P(B<sub>0</sub>, B<sub>1</sub>)<sub>D</sub> (B<sub>2</sub>)<sub>E</sub> (B<sub>3</sub>)<sub>E</sub> (B<sub>4</sub>)<sub>E</sub> (B<sub>5</sub>)<sub>E</sub> 0(B<sub>0</sub>, B<sub>1</sub>)<sub>D</sub> (B<sub>2</sub>)<sub>E</sub> (B<sub>3</sub>)<sub>E</sub> (B<sub>4</sub>, B<sub>5</sub>)<sub>E</sub> 1(B<sub>0</sub>, B<sub>1</sub>)<sub>D</sub> (B<sub>2</sub>)<sub>E</sub> (B<sub>3</sub>)<sub>E</sub> (B<sub>4</sub>, B<sub>5</sub>)<sub>E</sub> 2(B<sub>0</sub>, B<sub>1</sub>, B<sub>3</sub>)<sub>E</sub> (B<sub>2</sub>)<sub>E</sub> (B<sub>4</sub>, B<sub>5</sub>)<sub>E</sub> 3(B<sub>0</sub>, B<sub>1</sub>, B<sub>3</sub>)<sub>D</sub> (B<sub>2</sub>)<sub>E</sub> (B<sub>4</sub>, B<sub>5</sub>)<sub>E</sub> 3(B<sub>0</sub>, B<sub>1</sub>, B<sub>3</sub>, B<sub>2</sub>)<sub>E</sub> (B<sub>4</sub>, B<sub>5</sub>)<sub>E</sub> 4(B<sub>0</sub>, B<sub>1</sub>, B<sub>3</sub>, B<sub>2</sub>)<sub>D</sub> (B<sub>4</sub>, B<sub>5</sub>)<sub>E</sub> 4(B<sub>0</sub>, B<sub>1</sub>, B<sub>3</sub>, B<sub>2</sub>)<sub>D</sub> (B<sub>2</sub>, B<sub>4</sub>, B<sub>5</sub>)<sub>E</sub> 4INTERPROCEDURAL OPTIMIZATION:-

- \* Having more no:of procedure in the program has the advantage where it limits the length of the code (decreases in size).

Disadvantages.

- \* limits the compiler's ability to understand, what happens inside a call.
- Eg:- call by reference.
  - \* as call by reference changes the actual arguments, it cannot use it for further optimization.
  - \* each call requires a pre & post overheads (content should be stored in stack & PC is in different location)

In order to reduce the execution time, PC is involved & its content changes

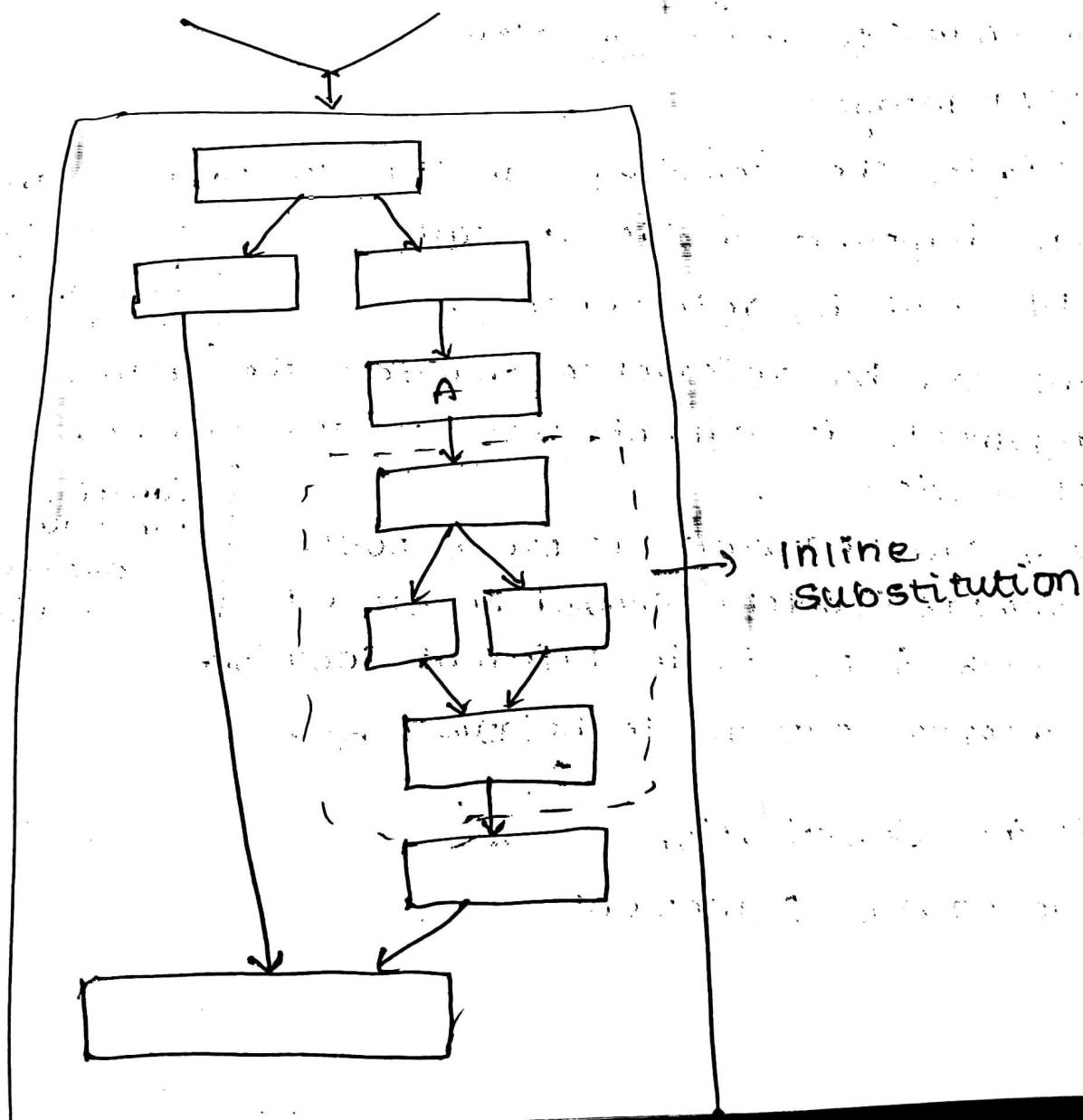
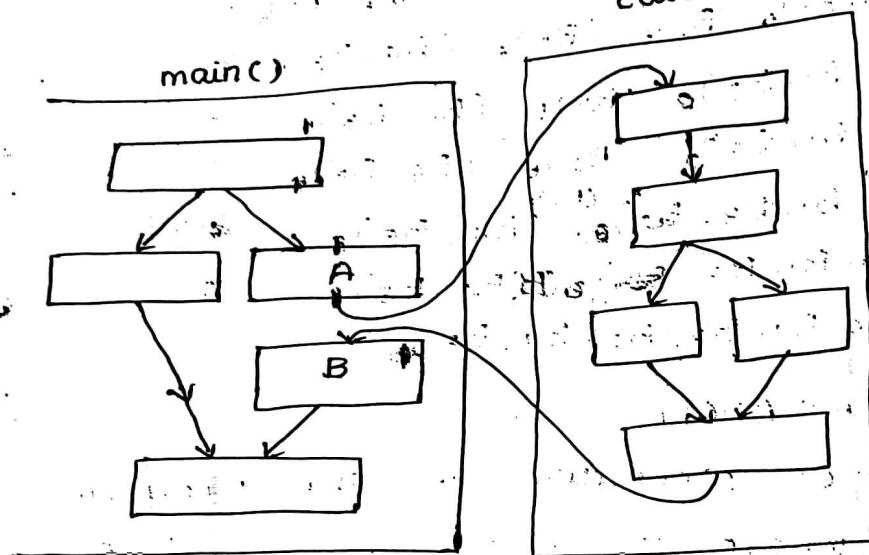
function  
procedure →  
overheads  
involved

Two interprocedural Techniques are:

1. Inline Substitution
2. procedural placement.

## INLINE SUBSTITUTION

- \* size is small
- \* less no of calls to function
- calc()



Empirical studies show that the decision on caller, the callee and the call site, the criteria's include

1. The callee size, if the callee size is smaller to the procedure linkage (precall, postcall, epilog & prolog) Inlining the callee should reduce the code size and also executes with fewer operation

2. If the overall size of any procedure Be lesser

3. Dynamic call count: Improve on frequently executed call site than infrequently executed call site

4. The constant valued actual parameters, and will definitely improve
5. other parameter values includes a static call count. The parameters count looping & meeting depth & fraction of execution time.

PROCEDURE PLACEMENT: Rearranging procedures

using callgraph. A callgraph as a node for each procedure and an edge  $(x,y)$  for each call from  $x$  to  $y$ . procedure placement algorithm has more freedom than block layout and depends on cache

- 1. Analysis - Operates on call graph
- 2. Transformation

\* The two phases are

ANALYSIS: Repeatedly selects two nodes, from callgraph and combines them. The order of the combination is driven by execution frequency data either measured or estimated.

TRANSFORMATION: Dependent on the order of combination (or) layout.

ALGORITHM: It operates over the programs' callgraph iteratively by considering edges in the order of their estimated execution frequency.

eg:-

1. Builds a callgraph

2. Assigns each edges an execution frequency into

3. combines all edges between two nodes into a single edge.

INITIALIZE: // Build the call multigraph  $G$ .

for each edge  $(x,y) \in G$  // add weight

IF  $(x=y)$  then

    delete  $(x,y)$  from  $G$

else

$wt((x,y))$  // execute frequently

for each node  $x \in G$

$list(x) \leftarrow \{x\}$

// init list

IF multiple edges exists from  $x$  to  $y$  then

    combine them & their actions.

for each edge  $(x,y) \in G$

? put each edge

    Enqueue  $((0, (x,y), wt(x,y)))$  } into  $Q$

? into  $Q$

// iterative reduction of graph.

while  $Q$  is not empty // take highest

$(x,y) \leftarrow Dequeue(Q)$  // priority else

    for each edge  $(y,z) \in G$

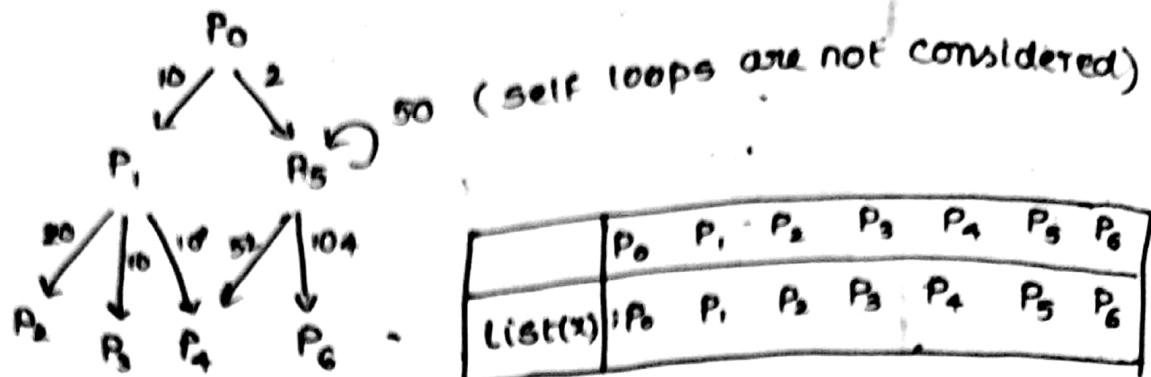
        Resource  $((y,z), x)$  // move source from  $y \rightarrow x$

    for each edge  $(y,z) \in G$  // move target

        Retangent  $((z,z), x)$  "from  $y$  to  $z$

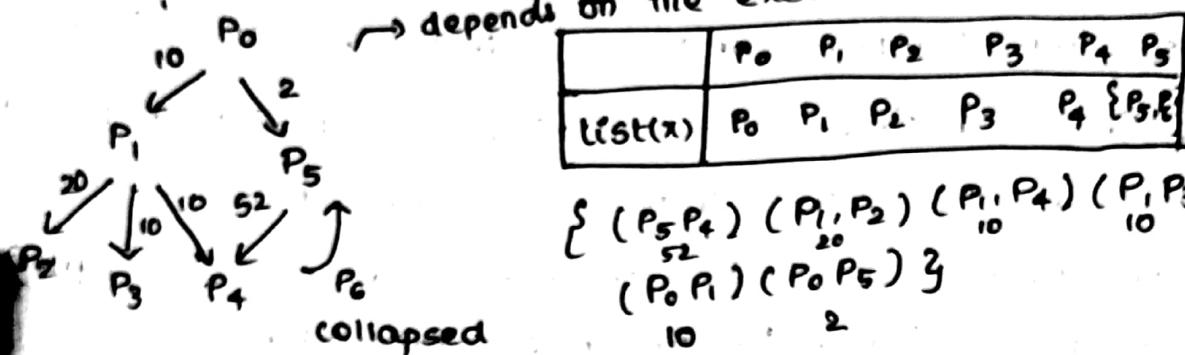
    append  $list(y)$  to  $list$ " // update

; delete  $y$  & its edges from

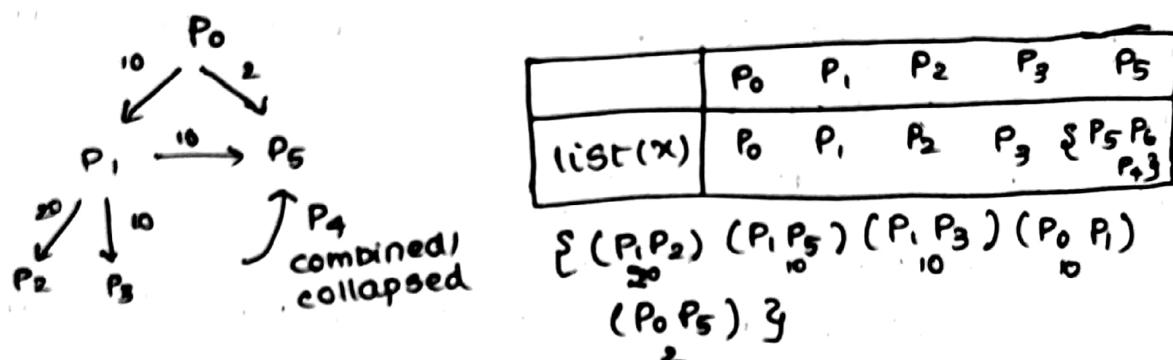


$\{ (P_5, P_6) (P_5, P_4) (P_1, P_2) (P_1, P_4) (P_1, P_3) (P_0, P_1) (P_0, P_5) \}$

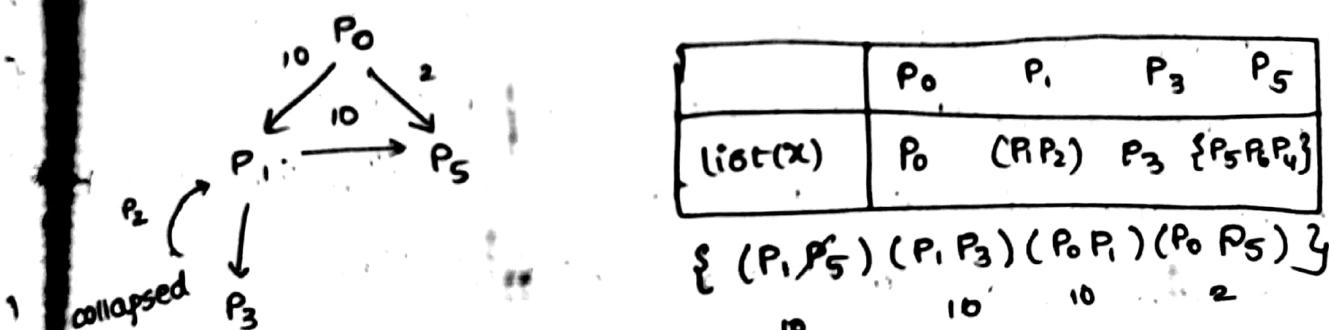
depends on the execution frequency.



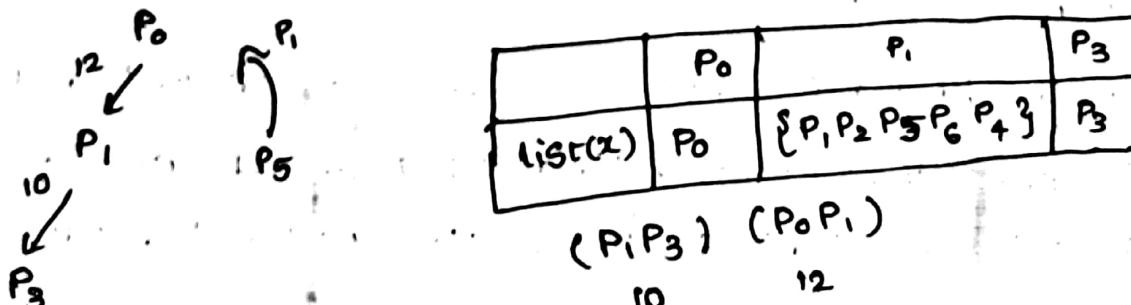
$\{ (P_5, P_4) (P_1, P_2) (P_1, P_4) (P_1, P_3)$   
 $(P_0, P_1) (P_0, P_5) \}$



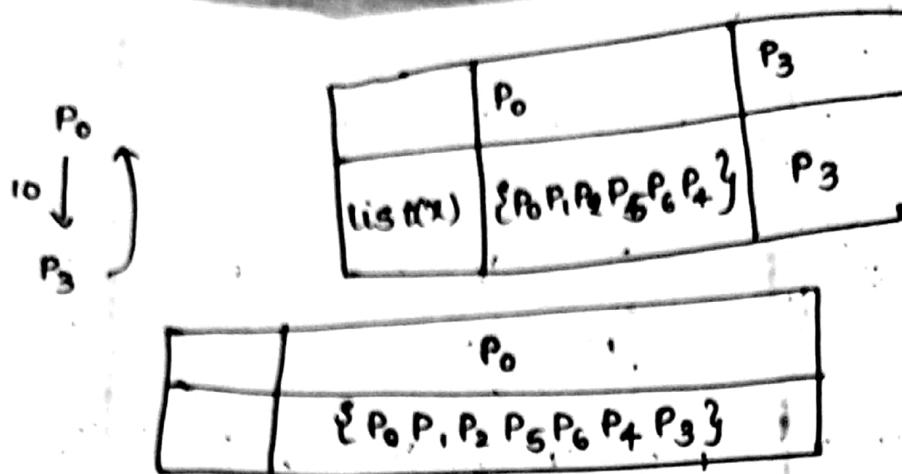
$\{ (P_1, P_2) (P_1, P_5) (P_1, P_3) (P_0, P_1)$   
 $(P_0, P_5) \}$



$\{ (P_1, P_5) (P_1, P_3) (P_0, P_1) (P_0, P_5) \}$



$(P_1, P_3) (P_0, P_1)$



- \* **DATA FLOW ANALYSIS:** It is a classical technique for compile time program analysis and it allows to reason the run time flow of values in the program.
- \* Before improving the code, the compiler must locate points in the program where modifying the code is likely to improve and safe.
- \* compiler uses static analysis methods for run time flow of values.
- \* Two steps in static analysis are:
  1. control flow analysis
  2. Data Flow analysis followed by SSA form.

By performing static analysis, the compiler will know exactly what will happen for instance,

- \* either immediately display results (or) if it gets inputs from users and takes modest time to display results.

**ITERATIVE DATA FLOW ANALYSIS:** It is a collection of techniques for compile time reasoning on the run time flow of values to locate opportunities for optimization and to prove the safety of specific optimization.

list of  
 T100  
 T1000  
 T10000  
 T100000  
 T1000000

transformation.  
THREE STEPS: FOR ITERATIVE ANALYSIS:

Build CFG

Gather initial information about blocks.

Solve equation to produce DOM (Dominance)

Dominance is the shape and structure of CFG is

the idea behind [dominance class]. dominators.

$$\text{DOM}(n) = \{n\} \cup (\bigcap_{m \in \text{pred}_S(n)} \text{DOM}(m))$$

