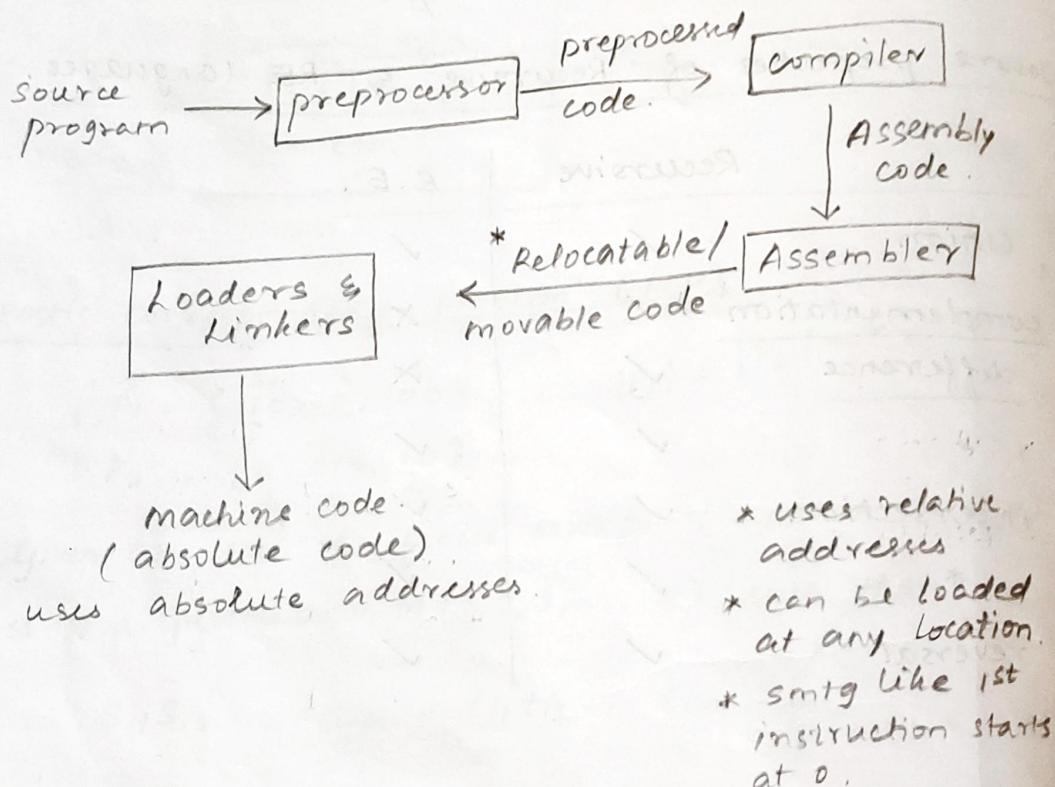
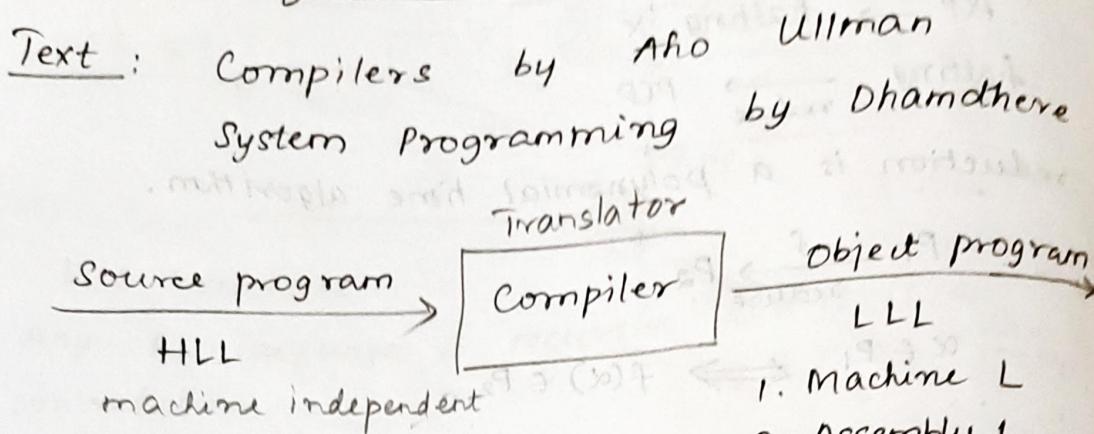
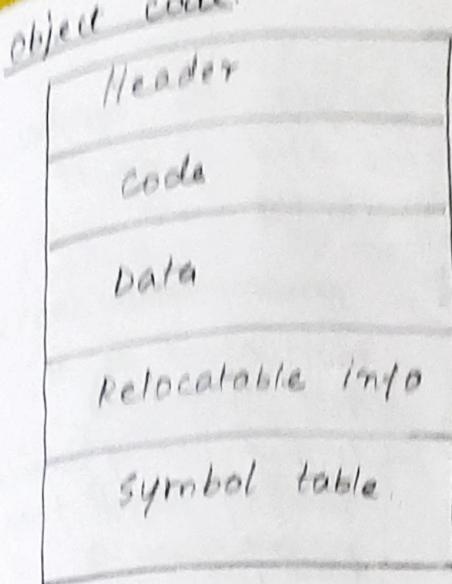
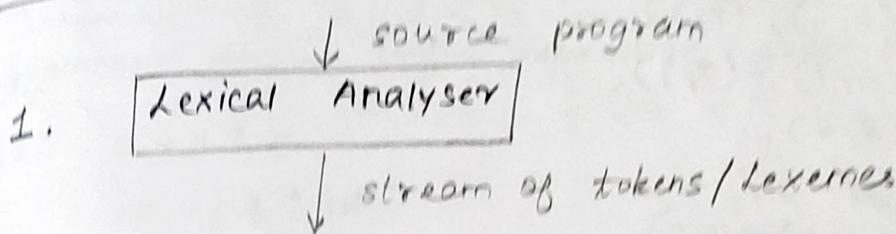


Principles of Compiler design





Modules of a compiler.



Tokens :

- i. keywords
 - ii. identifiers
 - iii. operators
 - iv. constants
 - v. delimiters

Token \rightarrow (type, value).

e.g.

e.g.	Token	Representation
	18	(1, -)
)	and for	(2, -)
)	while	(3, -)
	+	(4, 1)
identifier	$\leftarrow x$	(7, ptr to symbol table)

symbol table

	name	type	address	((8,1)
0	i	int var)	(8,2)
1	10	int const		{	(8,3)
2	a	int var		z	(8,4)
3	b	" var		;	(8,5)
4	l	" const		< =	(6,2)
				=	(6,3)

Stream of tokens.

```
while ( i <= 10 )
{
    a = a + b;
    i = i + 1;
}
```

↓

(3, -),
(8, 1),
(7, 0),
(6, 2),
(9, 1),
(7, 2),
(6, 3), (7, 2),
(4, 1), (7, 3)
(8, 5),
(7, 0),
(6, 3), (7, 0),
(4, 1), (7, 4),
(8, 5),
(8, 4).

lexical analyser
maintains the symbol
table and to do it
efficiently...

symbol table is
maintained as a
self organised
linked list / balanced
binary search tree,
hash / trie.

Two tasks/problems of the lexical analyser.

- i. description of a token. (regular expression)
- ii. recognition of a token. (automaton)

↓ stream of tokens

Parser / Syntax analyser

2.

combines tokens
uses grammar

↓ syntax tree / parse tree
to form a valid syntax.
 $G(N, T, P, S)$.

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E * E$$

$$E \rightarrow E / E$$

$$E \rightarrow id$$

$$N = \{E\}$$

$$T = \{+, -, *, /, id\}$$

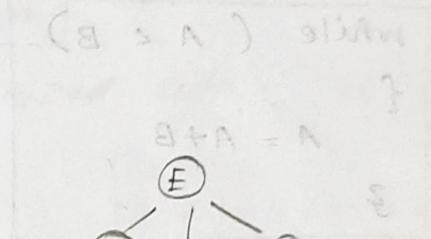
$$S = \{E\}$$

grammar for an infix expression (arithmetic)

ex. $a + b * c$.

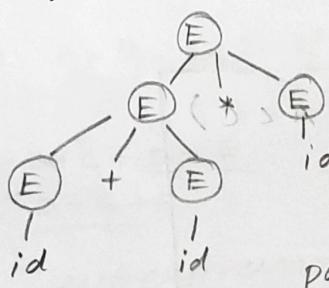
$id + id * id$.

$$\begin{aligned} E &\xrightarrow{1} E + E \\ &\xrightarrow{2} id + E \\ &\xrightarrow{3} id + E * E \\ &\xrightarrow{4} id + id * E \\ &\xrightarrow{5} id + id * id. \end{aligned}$$



since $E \xrightarrow{*} id + id * id$, $id + id * id$ is a valid infix expression.

(or)



parse tree

Derivation can be top-down / bottom-up.

Syntax tree

3.

Intermediate code generator



an intermediate code is Three Address Code (a max. of 3 references).

i. $\underline{A} = \underline{B} + \underline{op} \underline{C}$

ii. $\underline{A} = \underline{op} \underline{B}$

iii. goto \underline{C}

iv. if \underline{A} rel-op \underline{B} goto $\underline{L1}$

while ($A < B$)

{

$A = A + B$

}

1. if $\underline{A} < \underline{B}$ goto $\underline{3}$

2. goto $\underline{6}$

3. $T1 = \underline{A} + \underline{B}$

4. if $\underline{A} = \underline{T1}$

5. goto $\underline{1}$

6.

while ($A < B$ and $A < C$)

{

$A = A + B$

}

1. if $A < B$ goto 3
2. goto 8
3. if $A < C$ goto 5
4. goto 8

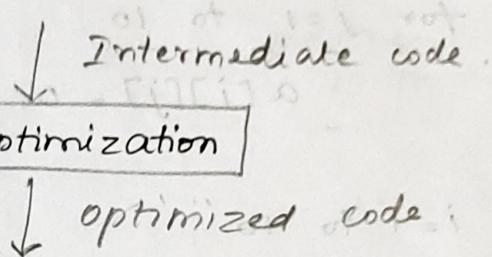
5. $T_1 = A + B$
6. $A = T_1$
7. goto 1

8.

90-10 rule

90% of execution time is spent on 10% of code.

4. Code optimization



eg. ① if $A = B$ goto L2
 goto L3

→

L2 :

1 CMP
1 JZ L2
0.5 JMP L3

* 2.5

1 CMP
1 JNZ L3.
2 ∵ optimized.

* assuming $A = B$ 50% of time.

②

$$A = B + C + \underline{D + E}$$

$$B = \underline{D + E} + F$$

→

$$T_1 = D + E$$

$$A = B + C + T_1$$

$$B = T_1 + F$$

common subexpression

Optimize:

$I = 1;$

while $I \leq 100$ do

{

$a[I] = 0;$

$I = I + 1;$

}

Technique:

Loop unrolling

while $I \leq 50$ do

{ $a[I] = 0;$

$a[J] = 0;$

$I++;$

$J++;$

}

Loop jamming

for $i=1$ to 10

 for $j=1$ to 10.

$a[i][j] = 0$

 redundant?

for $i=1$ to 10

{ for $j=1$ to 10

$a[i][j] = 0$

$a[i][i] = 1$

for $i=1$ to 10

$a[i][i] = 1$

}

other techniques:

index variable elimination

using cheap operations (+ cause adder is
directly available as hardware)

↓ optimized code

Code generation

5.

Basic instruction format

$$X = (A+B) * (C-D)$$

opcode	operation
--------	-----------

ADD R₁, A, B

3 address
instruction.

SUB R₂, C, D

MUL X, R₁, R₂

(General Register
organisation)

MOV R₁, A

ADD R₁, B

2 address
instruction.

MOV R₂, C

SUB R₂, D

MUL R₁, R₂

MOV X, R₁

Load Accumulator

LDA A

1 address
instruction.

ADD B

(using accumulator).

STA T

AC is one of the operands
and holds the result after.

LDA C

(single register
organisation)

SUB D

MUL T

STA X

PUSH A

0 address
instruction.

PUSH B

(stack organisation)

ADD

PUSH C

PUSH D

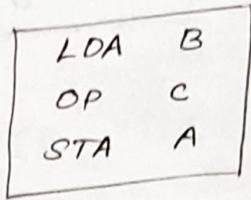
SUB

MUL

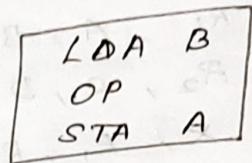
POP X

Three address code \Rightarrow 1 address instruction.

$A = B \text{ OP } C \Rightarrow$



$A = \text{OP } B \Rightarrow$



if $A \geq B$ goto L1 \Rightarrow

LDA	A
CMP	B
JNC	L1
jump no carry	

goto L1 \Rightarrow

JMP	L1
-----	----

Syntax directed Translation Scheme (SDTS)

Translated to code during Syntax analysis.

A ADD

B SUB

T JUM

X ATB

D ADD

G SUB

H JUM

I ATB

A H2VA

B H2VB

C H2VC

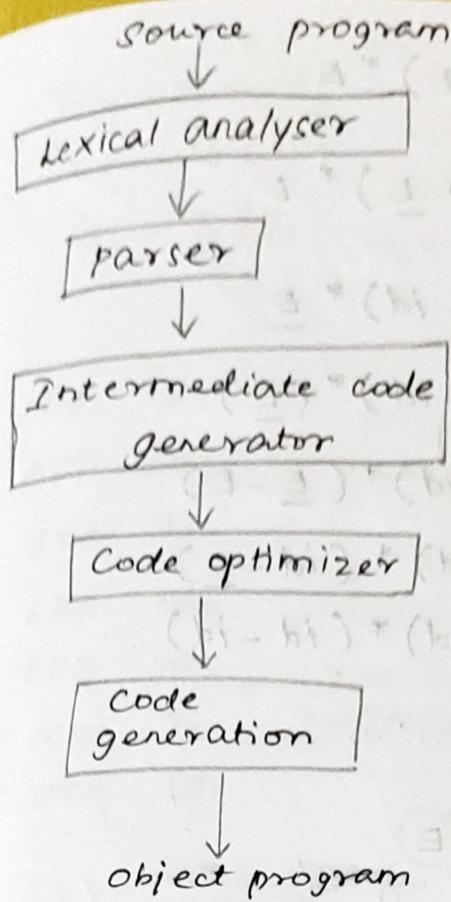
D H2VD

E H2VE

F H2VF

G H2VG

H H2VH



$$x = (A + B) * (C - D)$$

G1

grammer for assignment operation.

start symbol.

$$A \rightarrow id = E$$

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E * E$$

$$E \rightarrow E / E$$

$$E \rightarrow id$$

i. 'id', '=', '(', 'id₂', '+', 'id₃', ')', '*', '(', 'id₄', '-', 'id₅', ')'

ii. $A \Rightarrow id_1 = E$

$$\Rightarrow id_1 = \underline{E} * E$$

$$\Rightarrow id_1 = (\underline{E})^* E$$

$$\Rightarrow id = (\underline{E} + E)^{*E}$$

$$\Rightarrow id = (id + E)^* E.$$

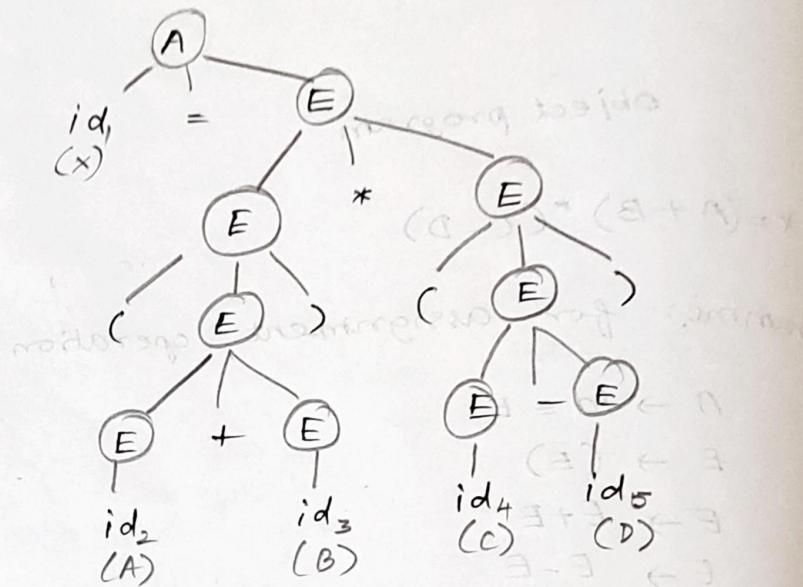
$$\Rightarrow id = (id + id)^* E$$

$$\Rightarrow id = (id + id)^* (E)$$

$$\Rightarrow id = (id + id)^* (E - E)$$

$$\Rightarrow id = (id + id) * (id - \underline{id})$$

$$\Rightarrow id = (id + id) * (id - id)$$



$$iii. \quad \overline{AB} = A - B$$

$$T_1 = id_2 + id_3$$

$$T_2 = id_4 - id_5$$

$$T_3 = T_1 * T_2$$

$$id_1 = T_3$$

$$T_1 = id_2 + id_3 \rightarrow \text{Instruction 1}$$

$$T_2 = id_4 - id_5$$

$$id_1 = T_1 * T_2$$

- v. 1 LDA id_2 }
2 ADD id_3 } instruction 1
3 STA T_1 }
4 LDA id_4
5 SVB $id_5 \xrightarrow{id_5}$ 6 STA T_2 redundant
6 MUL $T_2 \xrightarrow{T_2}$ LDATA
7 STA id_1

$$S = P * N * R / 100$$

To g_1 add

i. $'id'_1 = 'id'_2 * 'id'_3 * 'id'_4 / '1 '0 '0'$

$E \rightarrow I$
 $I \rightarrow I.\text{digit} | \text{digit}$
 $\text{digit} \rightarrow 0|1|2|3\dots|9$

ii. $A \Rightarrow id_1 = E$

Grammar for
integer const, I

$$\Rightarrow id_1 = \underline{E} * E$$

$$\Rightarrow id_1 = id_2 * \underline{E}$$

$$\Rightarrow id_1 = id_2 * id_3 * \underline{E}$$

$$\Rightarrow id_1 = id_2 * id_3 * \underline{id}_4 / E$$

$$\Rightarrow id_1 = id_2 * id_3 * id_4 / \underline{E}$$

$$\Rightarrow id_1 = id_2 * id_3 * id_4 / \underline{I}.$$
 digit

$$\Rightarrow id_1 = id_2 * id_3 * id_4 / \underline{I}.d$$

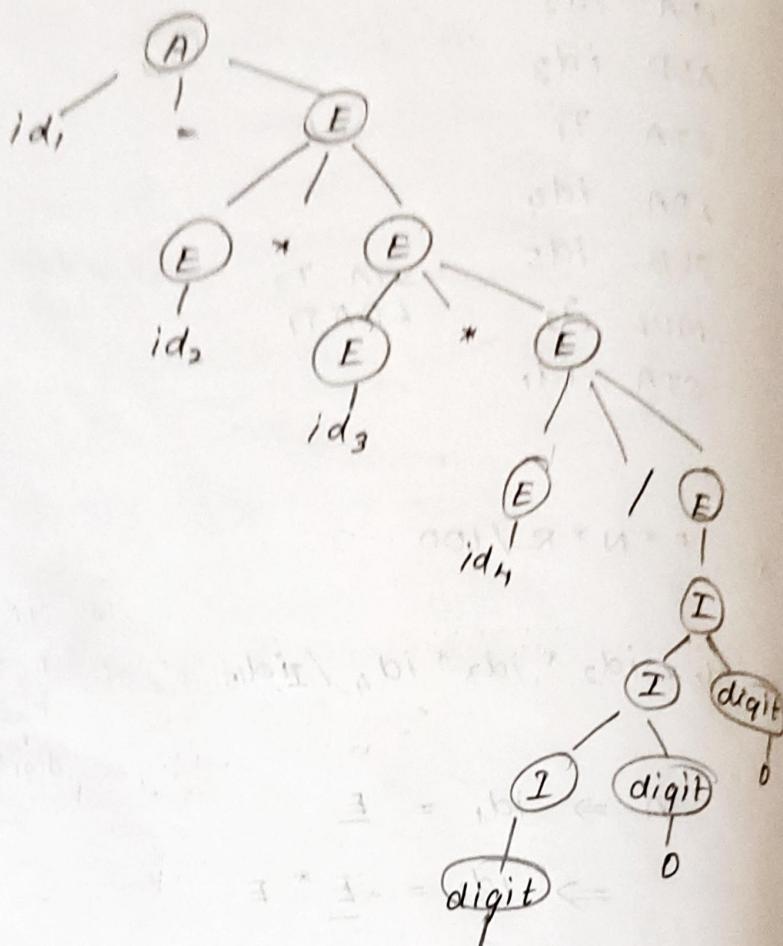
...

I.dd
ddd

$$\Rightarrow id_1 = id_2 * id_3 * id_4 / 100$$

$$\Rightarrow \quad \quad \quad / 100$$

$$\Rightarrow id_1 = id_2 * id_3 * id_4 / 100$$



iii.

$$T_1 = P * N$$

$$T_2 = T_1 * R$$

$$T_3 = T_2 / 100$$

$$S = T_3$$

$$3 * 3 * 3 / 100 = 100$$

$$3 * 3 * 3 / 100 = 100$$

$$3 * 3 * 3 / 100 = 100$$

$$3 * 3 * 3 / 100 = 100$$

If then else grammar.

condition statement

$S \rightarrow i \ C \ t \ s \ e \ s$

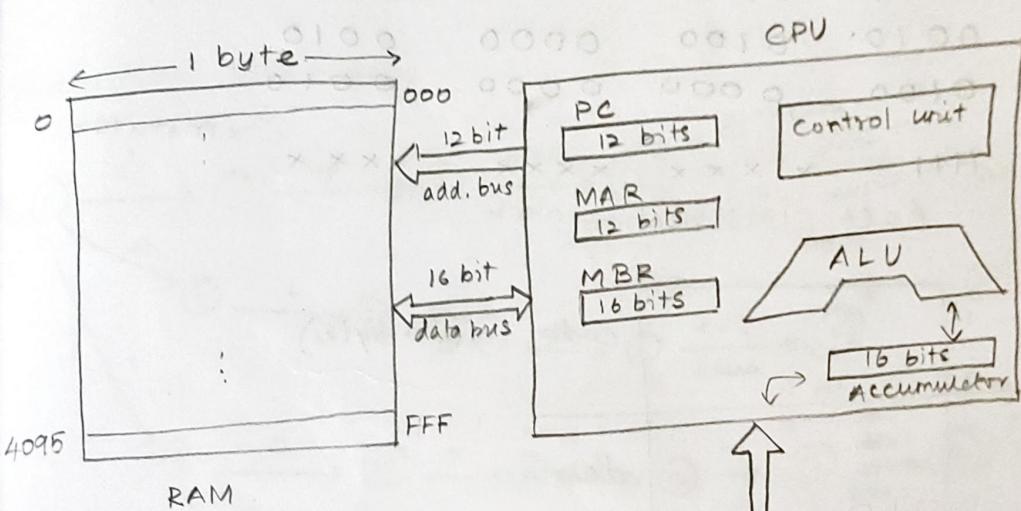
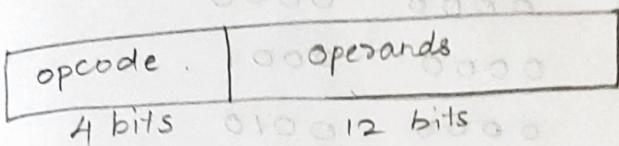
if then else

$e \rightarrow id \ relop \ id$

$S \rightarrow A$

if $a > b$ then $a = a + 2$ else $b = b * 2$. H.W.

system with 16 bit instruction



```

int a, b, c; // make entries in S.T.
Read a, b; // read from keyboard.
c = a+b;
Print c;
    
```

Symbol Table

name	type	address (in hexadec)
a	int	400
b	int	401
c	int	402

Keyboard

1	0000	0000	0000	0001
2	0001	0100	0000	0000
3	0000	0000	0000	0001
4	0001	0100	0000	0000
5	0010	0100	0000	0001
6	0011	0100	0000	0010
7	0001	0100	0000	0010
8	0010	0100	0000	0010
9	0100	0000	0000	0010
10	1111	-	xxxx	xxxx

monitor

halt instruction

0	i,	code	(20 bytes)
1	i.o.		
9	a b c	data	

opcodes :

- 0000 - Fetch a word from I/P device specified
- 0001 - Transfer AC to memory address specified.
- 0010 - Load the AC from the " " .

0011 - Add AC to memory specified
0100 - Transfer AC to O/P device.

Lexical Analyzer :

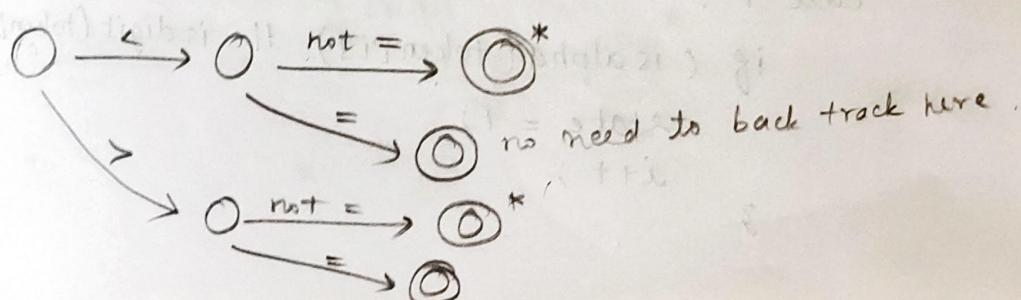
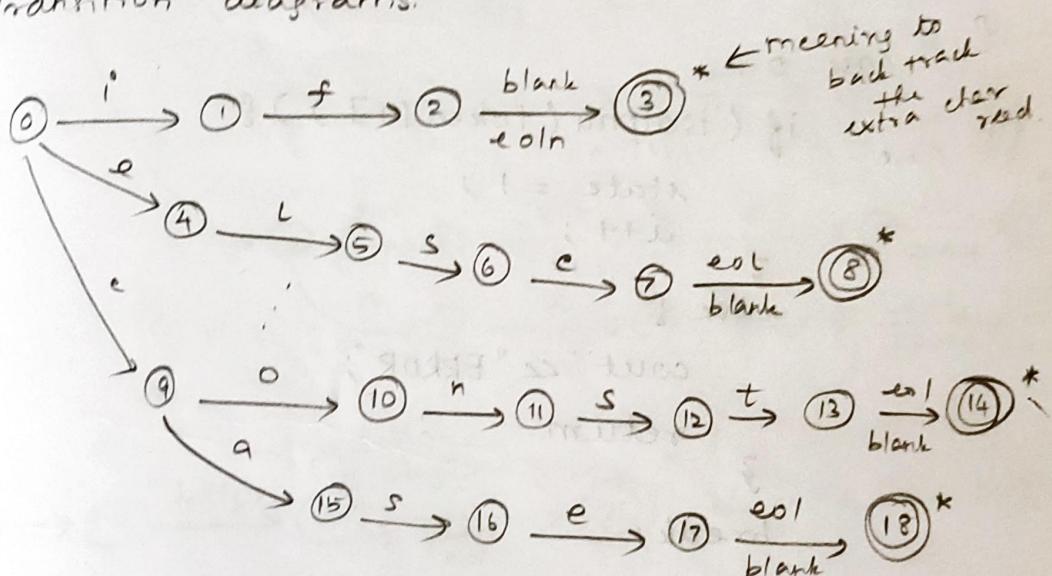
an effort

assume the following.

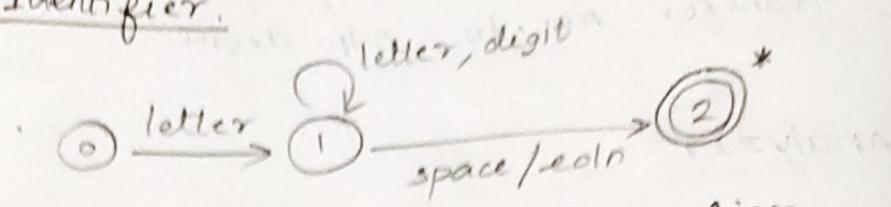
i. keywords : if, else, break, volatile,
void, const, case.

ii. relop : <
<=
>
>=
!=
==

Transition diagrams.



Identifier.



Program to read an identifier

```
int state, i;
char token[50];
cin >> token;
```

state = 0;
i = 0.

-switch (state)

```
while(1) {
    switch (state)
    {
        case 0:
            if (isalpha(token[i])) {
                state = 1;
                i++;
            }
            else {
                cout << "ERROR";
                return ;
            }
            break;

        case 1:
            if (isalpha(token[i]) || isdigit(token[i]))
                state = 1;
            i++;
    }
}
```

```

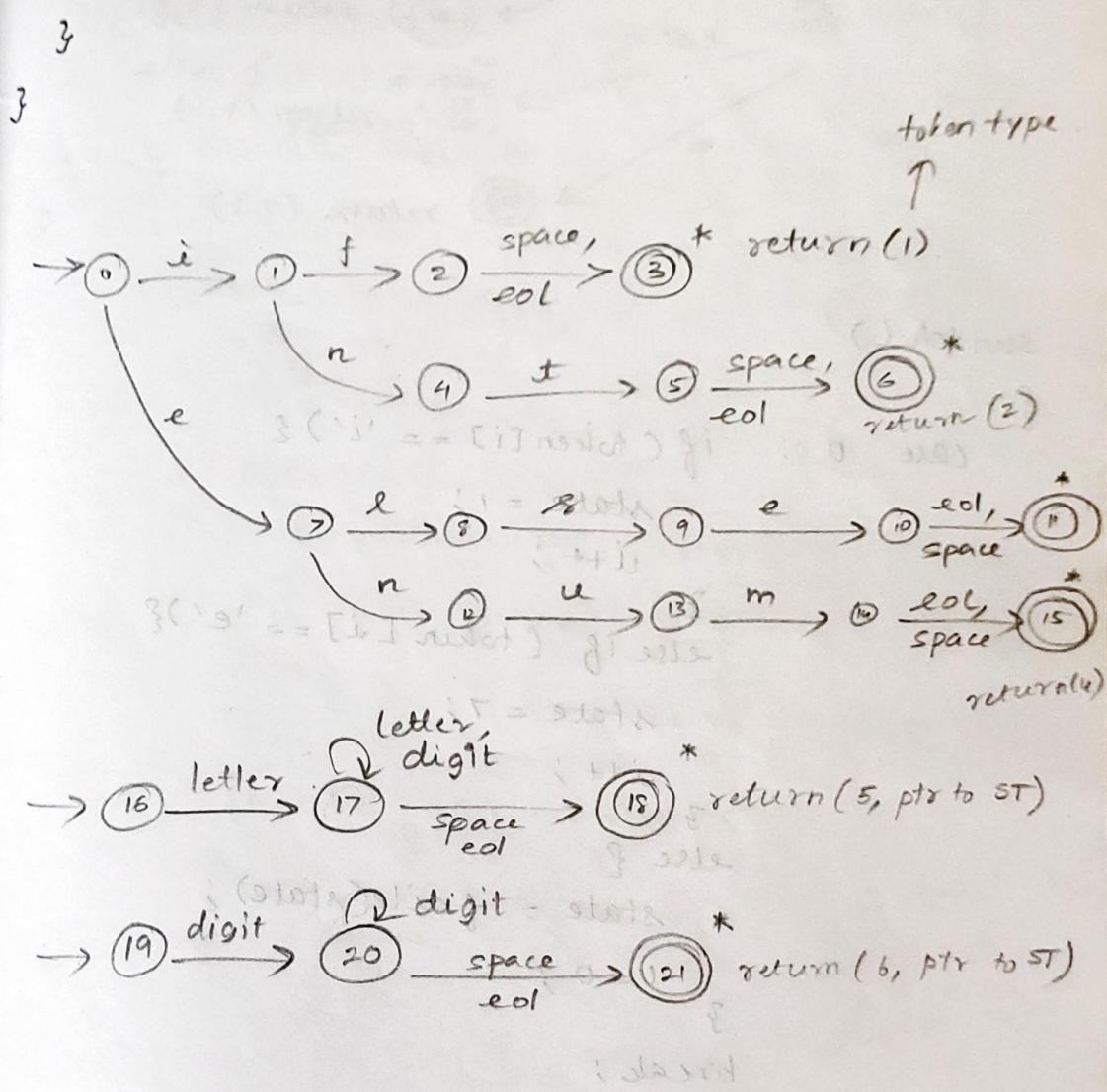
else if (token[i] == ' ' || token[i] == '\0') {
    state = 2;
    i++;
}

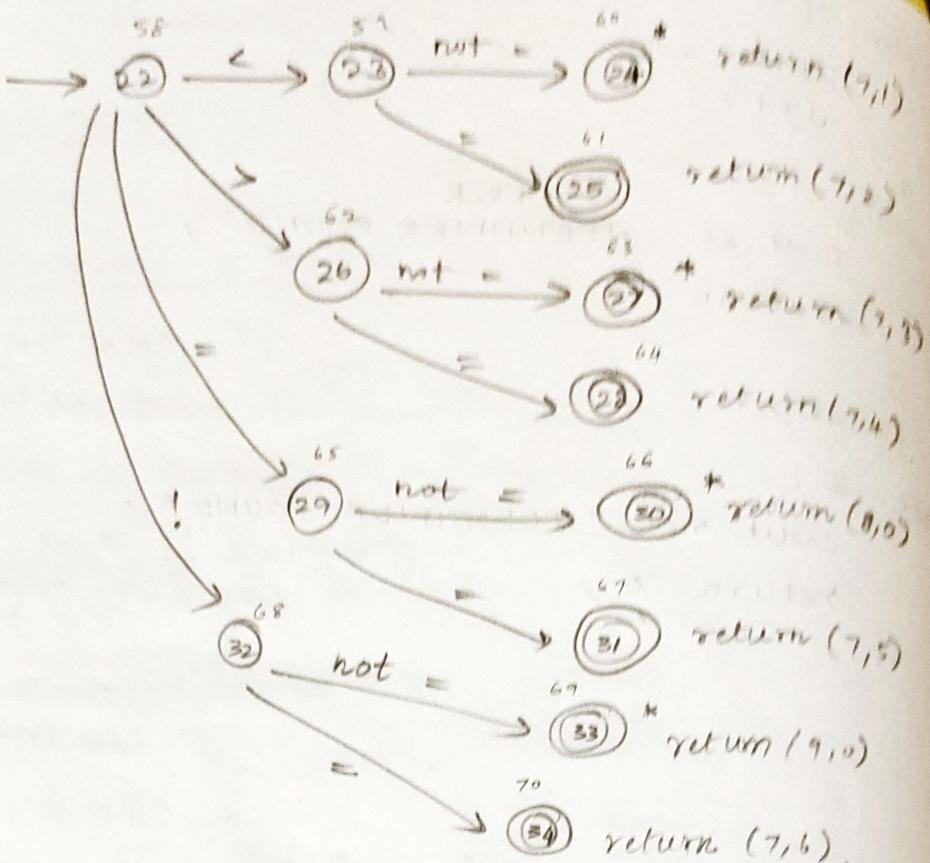
else {
    cout << "ERROR";
    cout << " IDENTIFIER FOUND ";
    return 1;
}

break;

case 2:
    cout << " IDENTIFIER FOUND ";
    cout << " // retract if there ";
    return 0;
    are more tokens
}

```





switch ()

{

```

case 0 : if (token[i] == 'i') {
            state = 1;
            i++;
        }
        else if (token[i] == 'e') {
            state = 7;
            i++;
        }
        else {
            state = fail(state);
            i = 0;
        }
        break;

```

case 1 :

```
int fail(int s)
```

```
{ if (s >= 22) {
```

```
    cout << "NOT a TOKEN";
```

```
    return ;
```

```
3 } else if (s >= 19) {
```

```
    return 22;
```

```
3 } else if (s >= 16) {
```

```
    return 19;
```

```
3 }
```

```
else {
```

```
    return 16;
```

```
3 }
```

```
3 } else {
```

Regular expression over a language L .

Regular expression is a finite tool
that can represent a finite/infinite language.

- i. ϵ is a Regular expression representing $\{\epsilon\}$.
- ii. for $a \in L$, ' a ' is a regular expression representing $\{a\}$.
- iii. if R and S are regular expressions,
representing L_R and L_S respectively.
 - a) $R + S$ is a regex, denotes $L_R \cup L_S$,
 - b) $R.S$ is a regex, denotes $L_R L_S$,
 - c) R^* is a regex, denotes L_R^* .
o or more 'a's.

Find language of a^* .

→ 'a' is a reg ex. representing $L_1 = \{a\}$.

→ a^* is a reg ex. representing L_1^*

$$L_1^* = \bigcup_{i=0}^{\infty} L_1$$

$$= \{\epsilon\} \cup \{a\} \cup \{aa\} \cup \{aaa\} \cup \dots$$

$$= \{\epsilon, a, aa, aaa, aaaa, \dots\}$$

Priority. ($*$, $,$, $,$, 1).
Note:

$(a/b)^*$.

- a is a regex representing $L_1 = \{a\}$
→ b is " " " " $L_2 = \{b\}$
→ a/b is " " " " $L_3 = L_1 \cup L_2$
= $\{a, b\}$.

→ $(a/b)^*$ is " " " " L_3^*

$$\begin{aligned}L_3^* &= \bigcup_{k=0}^{\infty} L_3^k \\&= \{\epsilon\} \cup \{a, b\} \cup \{aa, ab, ba, bb, \dots\} \\&= \{\epsilon, a, b, aa, ab, ba, bb, \dots\}.\end{aligned}$$

$a(a/b)^*$?

$$\begin{aligned}L &= \{az, \epsilon, a, b, aa, ab, ba, bb, \dots\} \\&= \{a, aa, ab, aaa, aab, aba, abb, \dots\}.\end{aligned}$$

Letter = $a | b | c | \dots | z$

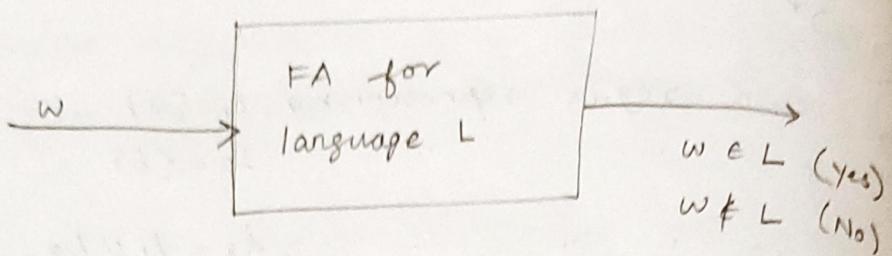
digit = $0 | 1 | 2 | \dots | 9$

id = letter. (letter | digit)*

const = digit. digit*

keywords = if | else | switch | case

Finite Automaton



- i. Non Deterministic Finite - state Automaton.
 - allow ϵ transitions.
 - allow many transition from a state on taking the same input symbol.
- ii. Deterministic Finite - state Automaton.
 - one and only one path for a string.

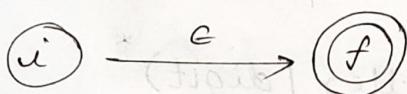
Regular expression to NFA :

over Σ

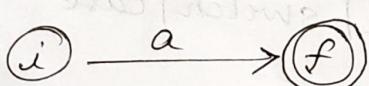
Input : RegEx representing Language L

Output : NDFA for the language L .

1. For RegEx , e.

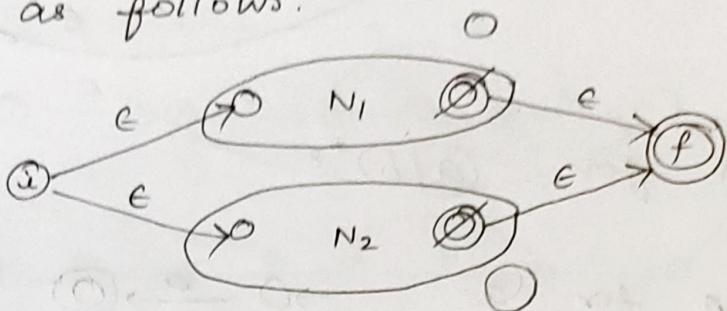


2. For RegEx 'a' over Σ .

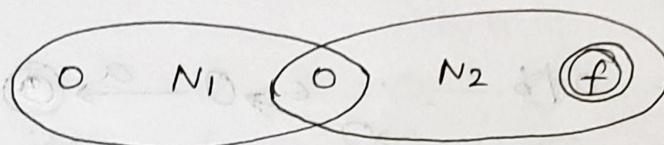


3. if there are regular expressions R_1 and R_2 representing automata N_1 and N_2 respectively.

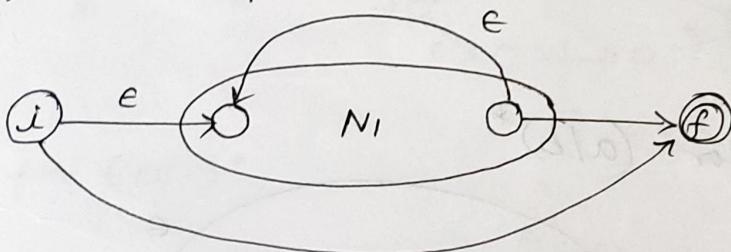
a) For $R_1 | R_2$, NFA is constructed as follows.



b) For $R_1 \cdot R_2$



c) for R_1^*

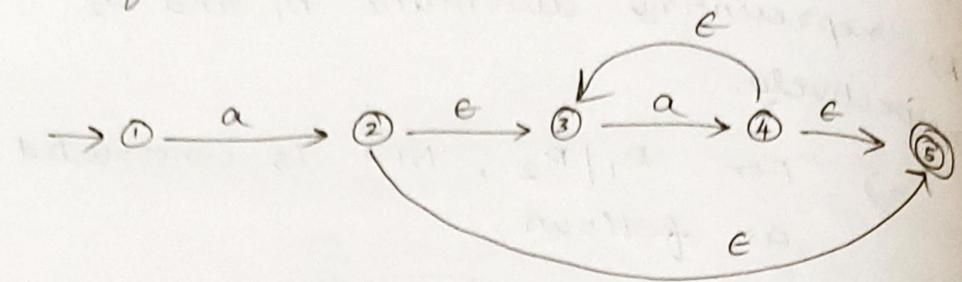


Consider aa^*

NFA for $a \rightarrow q_1 \xrightarrow{a} q_f \quad N_1$

NFA for $a^* \rightarrow q_1 \xrightarrow{\epsilon} q_1 \xrightarrow{a} q_2 \xrightarrow{\epsilon} q_2 \xrightarrow{a} q_2 \xrightarrow{\epsilon} q_f \quad N_2$

NFA for $a a^*$

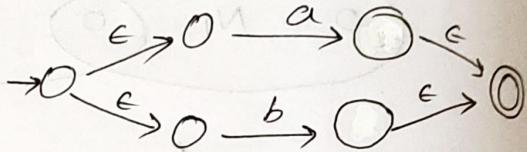


NFA for $(a/b)^*$?

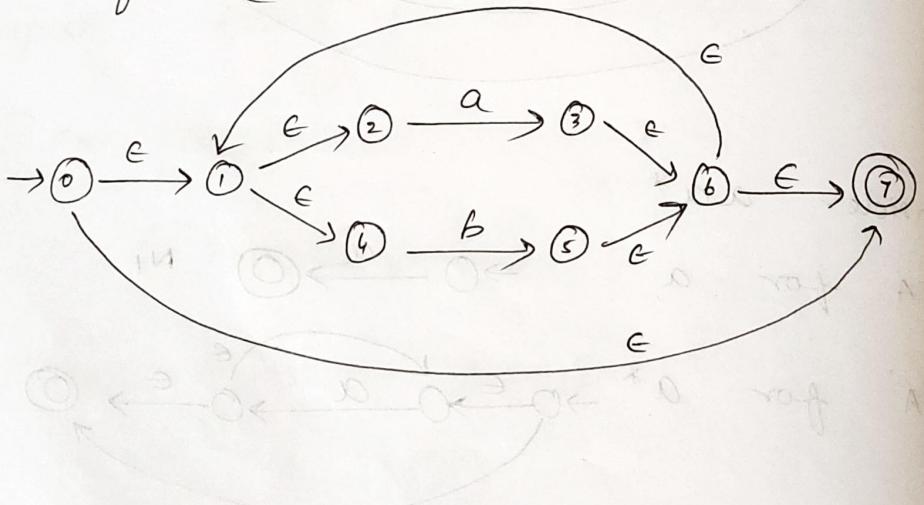
NFA for $a \rightarrow 0 \xrightarrow{a} 0$

" for $b \rightarrow 0 \xrightarrow{b} 0$

for a/b



NFA for $(a/b)^*$

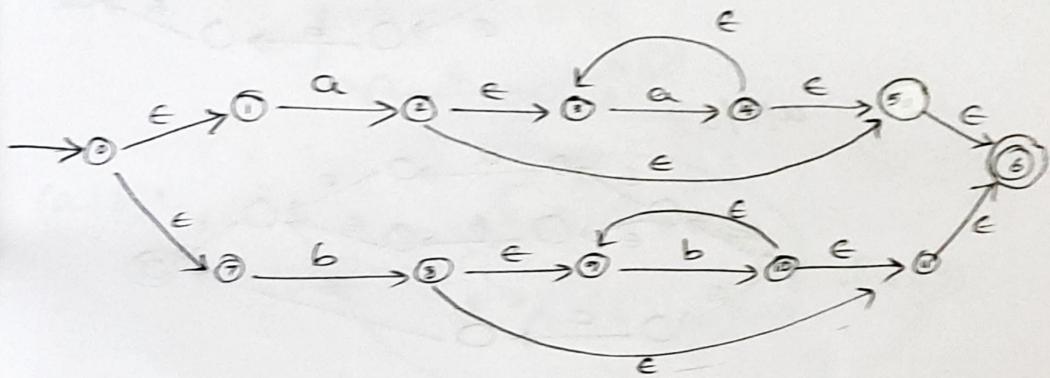


i. aa^*/bb^*

ii. $(0/1)^*(0/1)$

iii. $(a/b/c/d)^*$

i.

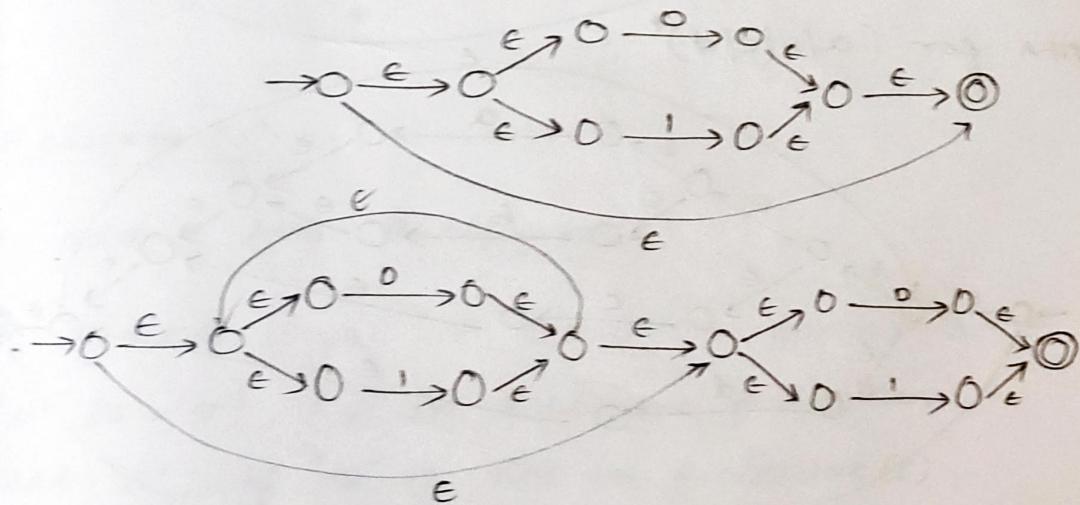


ii. NFA for 0 $\xrightarrow{0} \circ$

" for 1 $\xrightarrow{1} \circ$

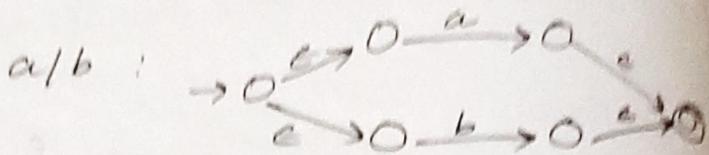
" for 0/1 $\xrightarrow{\epsilon} \circ \xrightarrow{0} \circ \xrightarrow{\epsilon} \circ$
 $\xrightarrow{\epsilon} \circ \xrightarrow{1} \circ \xrightarrow{\epsilon} \circ$

NFA for $(0/1)^*$

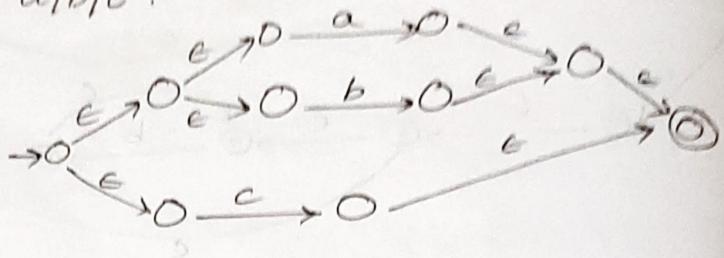


iii. NFA for a : $\rightarrow O \xrightarrow{a} O$

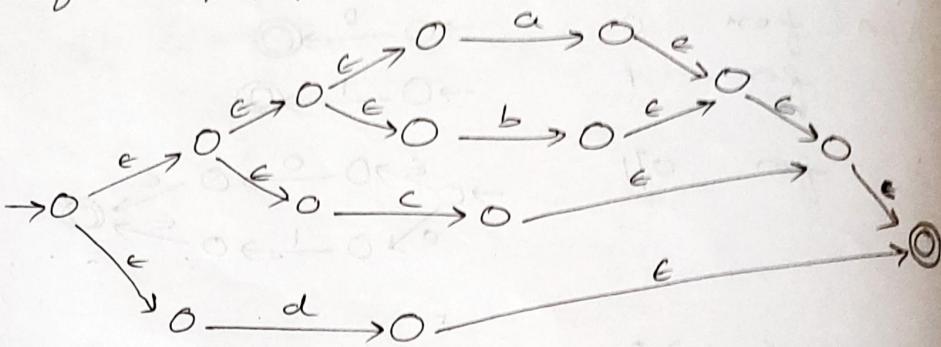
b : $\rightarrow O \xrightarrow{b} O$



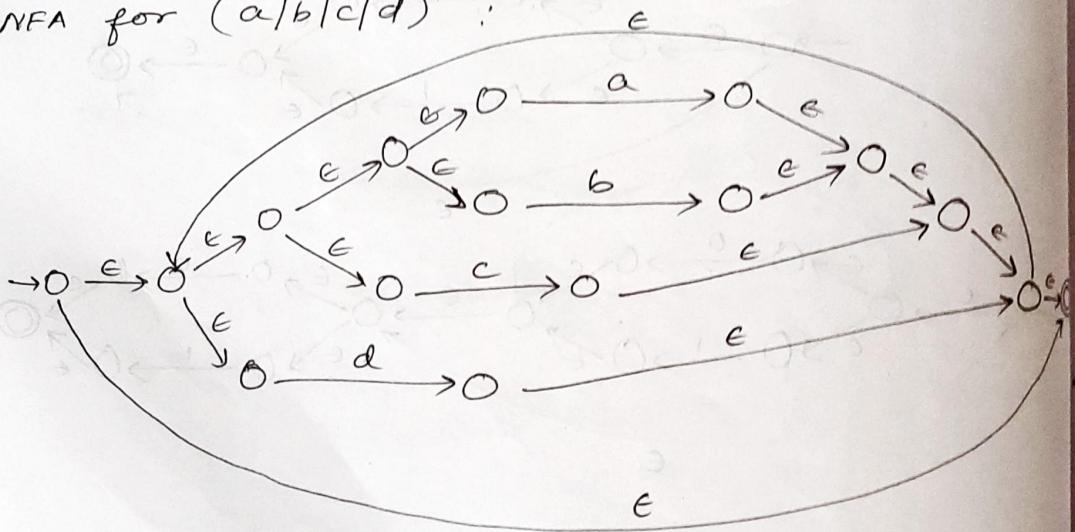
$a/b/c$:



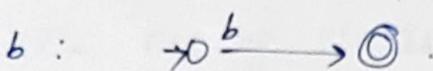
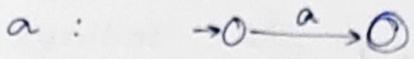
NFA for $a/b/c/d$:



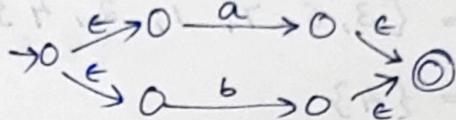
NFA for $(a/b/c/d)^*$:



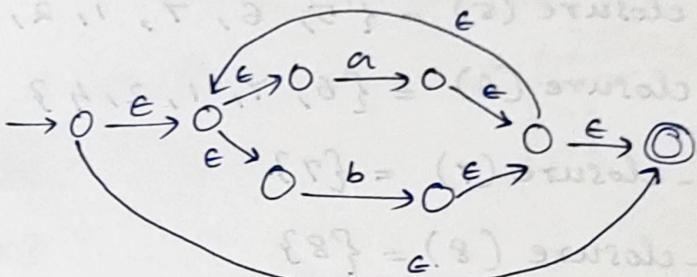
$(a/b)^*abb$



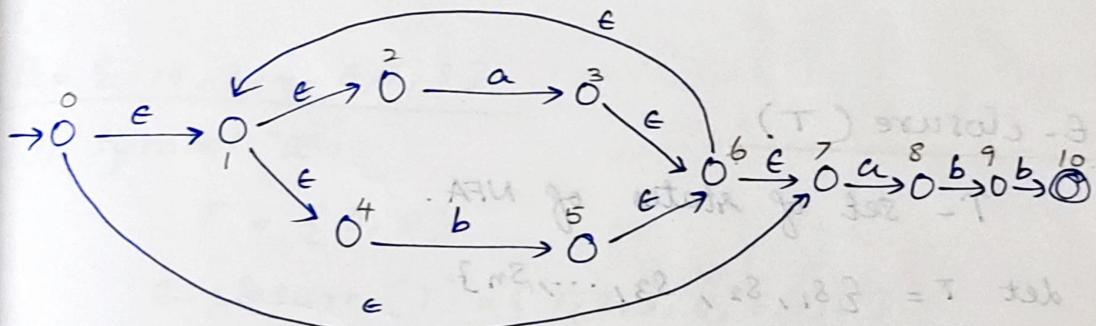
a/b :



$(a/b)^*$:



$(a/b)^*abb$:



ϵ closure (s)

1. add s into ϵ closure(s).

2. If there is an ϵ transition from s to t , if 't' is in ϵ closure of s, then add 't', if 't' is not in ϵ -closure(s).

3. repeat step 2 till no more states can be added.

$$\epsilon\text{-closure}(0) = \{0, 1, 7, 2, 4\}$$

$$\epsilon\text{-closure}(1) = \{1, 2, 4\}$$

$$\epsilon\text{-closure}(2) = \{2\}$$

$$\epsilon\text{-closure}(3) = \{3, 6, 7, 1, 2, 4\}$$

$$\epsilon\text{-closure}(4) = \{4\}$$

$$\epsilon\text{-closure}(5) = \{5, 6, 7, 1, 2, 4\}$$

$$\epsilon\text{-closure}(6) = \{6, 7, 1, 2, 4\}$$

$$\epsilon\text{-closure}(7) = \{7\}$$

$$\epsilon\text{-closure}(8) = \{8\}$$

$$\epsilon\text{-closure}(9) = \{9\}$$

$$\epsilon\text{-closure}(10) = \{10\}$$

$\epsilon\text{-closure}(T)$

T - set of states of NFA.

$$\text{det } T = \{s_1, s_2, s_3, \dots, s_n\}$$

$$\epsilon\text{-closure}(T) = \epsilon\text{-closure}(s_1) \cup \epsilon\text{-closure}(s_2) \cup \dots$$

(2) symbols \rightarrow start 2 lines
 $\epsilon\text{-closure}(s_n)$.

$$\epsilon\text{-closure}(\{1, 9\}) = \{1, 2, 4, 9\}$$

(2) symbols \rightarrow in top 2 lines
if no state draw on hit so get target

and no state draw on hit so get target
below

Initial of NFA
 Consider $A = \{0, 1, 2, 4, 7\} = \epsilon\text{-closure}(0)$,
 i/p symbol 'a'. The initial of DFA

$T = \text{set of states that can be reached from } A \text{ through 'a'}$.

$$= \{3, 8\}$$

$\epsilon\text{-closure}(T) = \{1, 2, 4, 3, 6, 7, 8\}$ new state
 $= B$. $\begin{array}{l|l} \text{'a' transitions} & \\ \hline 1 & \\ 2 \rightarrow 3 & \\ 8 & \\ 7 \rightarrow 8 & \\ 9 & \\ 9 \rightarrow 10 & \end{array}$

i/p symbol 'b': $\{1, 2, 4, 3, 6, 7, 8\} = (T)\bar{S}$

$$T = \{5\}$$

$\epsilon\text{-closure}(T) = \{1, 2, 4, 5, 6, 7\}$
 $= C$. $\{1, 2, 3, 4, 5, 6, 7\} = T$

$B = \{1, 3, 4, 6, 7, 8\}$ $\{1, 2, 3, 4, 5, 6, 7\} = T$

i/p symbol 'a':

$$T = \{3, 8\}$$

$\epsilon\text{-closure}(T) = \{1, 2, 3, 4, 6, 7, 8\}$
 $= B$. $\{1, 2, 3, 4, 5, 6, 7\} = T$

i/p symbol 'b':

$$T = \{5, 9\}$$

$\epsilon\text{-closure}(T) = \{1, 2, 4, 5, 6, 7, 9\}$
 $= D$. $\{1, 2, 4, 5, 6, 7, 9\} = T$

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

i/p symbol 'a':

$$T = \{3, 8\}$$

$$\bar{\epsilon}(T) = \{1, 2, 3, 4, 6, 7, 8\} = B$$

i/p symbol 'b': $\{1, 2, 4, 5, 6, 7, 9\} = T$

$$T = \{5\}$$

$$\bar{e}(T) = \{1, 2, 4, 5, 6, 7, 3\} \\ = C.$$

$$D = \{1, 2, 4, 5, 6, 7, 9\}$$

$\{1, 2, 4, 5, 6, 7, 9\} = (T) \text{ states}$

i/p symbol 'a':

$$T = \{3, 8\}$$

$$\bar{e}(T) = \{1, 2, 3, 4, 6, 7, 8\}$$

$$= B.$$

i/p symbol 'b':

$$T = \{5, 10\}$$

$$\bar{e}(T) = \{1, 2, 4, 5, 6, 7, 10\}$$

$$= E.$$

\therefore final

final state from NFA

$$E = \{1, 2, 4, 5, 6, 7, 10\}$$

i/p symbol 'a':

$$T = \{3, 8\}$$

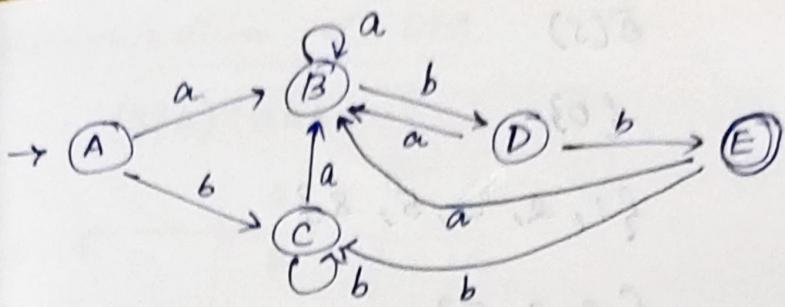
$$\bar{e}(T) = B.$$

i/p symbol 'b':

$$T = \{5\}$$

$$\bar{e}(T) = C.$$

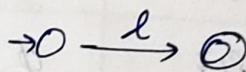
$$S = \{8, 5, 2, 1, 8, 1, 3\} = (T) \bar{S}$$



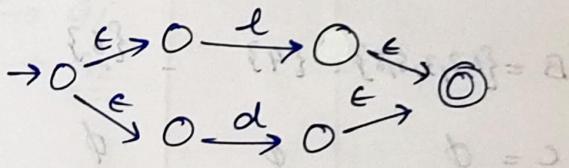
	a	b
A	B	C
B	B	D
C	B	C
D	B	E
* E	B	C

construct DFA for $l \cdot (l/d)^*$

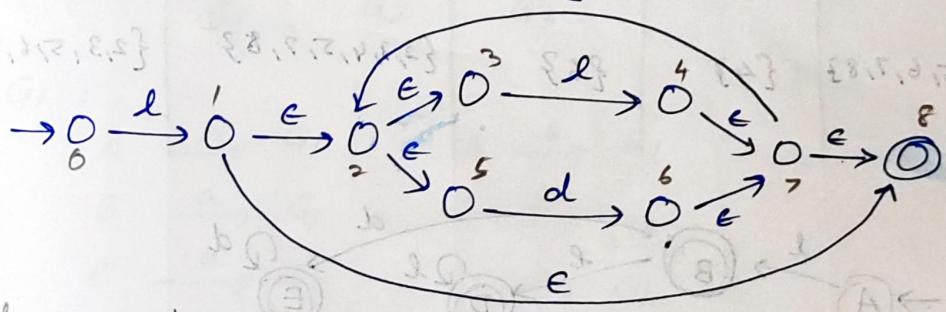
NFA for l :



l/d :



$l \cdot (l/d)^*$



l :
 $0 \rightarrow 1$
 $3 \rightarrow 4$
 d :
 $5 \rightarrow 6$
 ϵ :
 $2 \rightarrow 3$
 $2 \rightarrow 4$
 $4 \rightarrow 5$
 $4 \rightarrow 6$
 $5 \rightarrow 7$
 $5 \rightarrow 8$
 $6 \rightarrow 7$
 $6 \rightarrow 8$
 $7 \rightarrow 8$

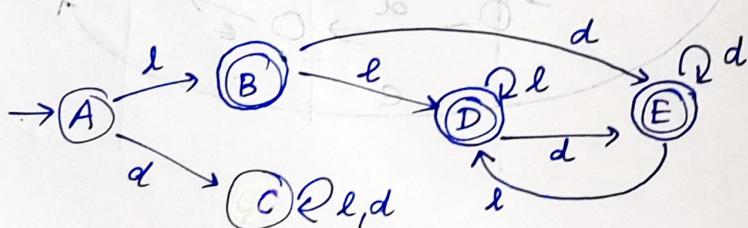
$l \cdot l = A$

$l \cdot d = B$

$l \cdot \epsilon = C$

s	$\bar{e}(s)$
0	$\{0\}$
1	$\{1, 2, 3, 5, 8\}$
2	$\{2, 3, 5\}$
3	$\{3\}$
4	$\{4, 2, 3, 5, 7, 8\}$
5	$\{5\}$
6	$\{6, 2, 3, 5, 7, 8\}$
7	$\{2, 3, 5, 7, 8\}$
8	$\{8\}$

	T	$*(\text{b12}).2$	$\bar{e}(T)$	
$A = \{0\}$	$\{1\}$	\emptyset	$\{1, 2, 3, 5, 8\}$	\emptyset
$*B = \{1, 2, 3, 5, 8\}$	$\{4\}$	$\{6\}$	$\{2, 3, 4, 5, 7, 8\}$	$\{2, 3, 5, 6, 7\}$
$C = \emptyset$	\emptyset	\emptyset	\emptyset	\emptyset
$D = \{2, 3, 4, 5, 7, 8\}$	$\{4\}$	$\{6\}$	$\{2, 3, 4, 5, 7, 8\}$	$\{2, 3, 5, 6, 7\}$
$*E = \{2, 3, 5, 6, 7, 8\}$	$\{4\}$	$\{6\}$	$\{2, 3, 4, 5, 7, 8\}$	$\{2, 3, 5, 6, 7\}$



Minimization of DFA

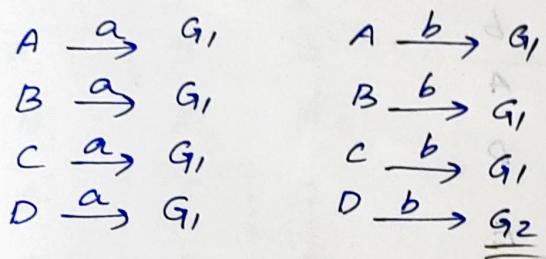
$(a/b)^*abb$.

	a	b
$\rightarrow A$	B	C
B	B	D
C	B	C
D	B	E
* E	B	C

$$\Pi_0 = \{(\text{Non final}) (\text{Final})\}$$

$$= \{ (ABCD) (E) \}$$

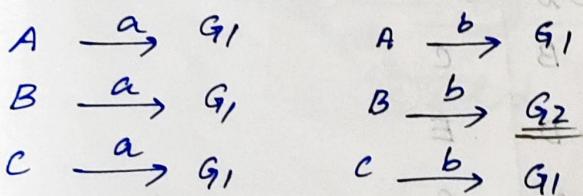
G_1 :



$$\Pi_1 = \{ (ABC) (D) (E) \}$$

$(^*(ab)) G_1 \cup G_2 \cup G_3$

G_1 :



$$\Pi_2 = \{ (AC) (B) (D) (E) \}$$

$G_1 :$

$$A \xrightarrow{a} G_2$$

$$B \xrightarrow{a} G_2$$

$$A \xrightarrow{b} G_1$$

$$B \xrightarrow{b} G_1$$

$$\therefore \pi = \{ (AC), (B), (D), (E) \}$$

minimized DFA

$A \& C$ are indistinguishable



	a	b
A	B	ΔA
B	B	D
ΔA	B	ΔA ← redundant.
D	B	E { (char) (var) }
E	B	ΔA . { (3) (ABA) }



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

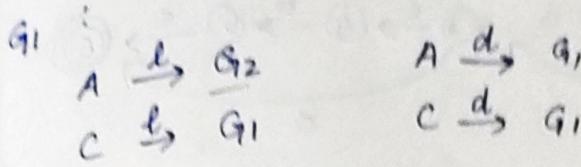
{ (3) (D) (ABA) }

minimize DFA of identifier. ($L(L/d)^*$)

	l	d
$\rightarrow A$	B	C
* B	D	E
C	C	C
* D	D	E
* E	D	E

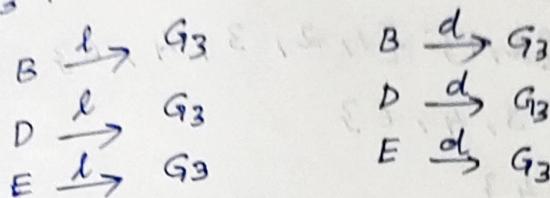
{ (S) (A) }

$$\Pi_0 = \{ (AC)_{G_1} (BDE)_{G_2} \}$$



$$\Pi_1 = \{ (A)_{G_1} (c)_{G_2} (BDE)_{G_3} \}$$

$G_3 :$

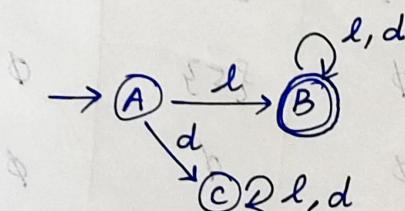


$$\therefore \Pi = \{ (A) (c) (BDE) \}$$

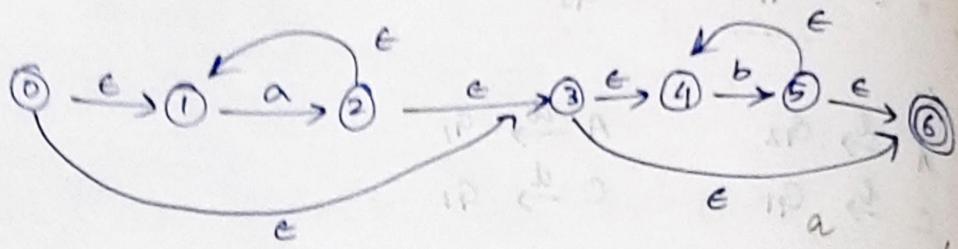
	l	d
$\rightarrow A$	B	C
$* B$	D \cancel{B}	E \cancel{B}
$\circ C$	C \cancel{C}	C
$* D$	D	E
$* E$	D \cancel{D}	E

	l	d
$\rightarrow A$	B	C
$* B$	B	B
c	c	c

minimized DFA.



Guess DFA for a^*b^*



s

$\bar{e}(S)$

0

$\{0, 1, 3, 4, 6\}$

1

$\{1\}$

2

$\{2\} \cup \{1, 2, 3, 4, 6\}$

3

$\{3, 4, 6\}$

4

$\{4\}$

5

$\{4, 5, 6\}$

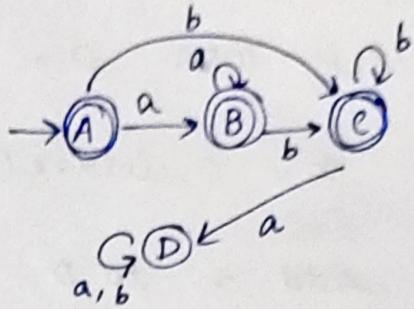
6

$\{6\}$

		A	$\bar{e}(T)$
		T	
A = $\{0, 1, 3, 4, 6\}$	$\{2\}$	$\{5\}$	$\{1, 2, 3, 4, 6\}$ $\Rightarrow B$
B = $\{1, 2, 3, 4, 6\}$	$\{2\}$	$\{5\}$	B
C = $\{4, 5, 6\}$	\emptyset	$\{5\}$	\emptyset
D = \emptyset	\emptyset	\emptyset	\emptyset

DFA

	a	b
*A	B	c
*B	B	c
*C	D	c
D	D	D

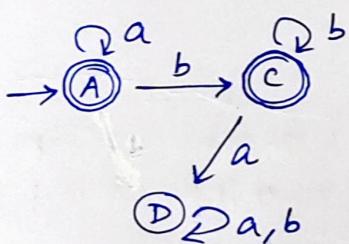


$$\pi_0 = \{ (A \ B \ C) \ (D) \}$$

$$\pi_1 = \{ (AB) \ (C) \ (D) \}$$

minimized DFA.

	a	b
*A	A	c
*C	D	c
D	D	D



Add minimization of DFA

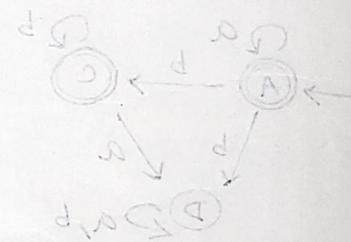
$$\{ (a \ b \ c) \ (a \ b \ c \ d) \} = \pi_1 \text{ (two terms)}$$

ab if ab is good then we can

test ab & if ab is redundant then

Algorithm for converting NFA to DFA.

- i. Let DFA $D = \{S_1, S_2, \dots, S_n\}$.
- ii. $x = \text{E-closure}(o)$ o - initial state
NFA.
Add x to D .
Unmark x .
- iii. While there is an unmarked state
 $'x' = \{S_1, S_2, \dots, S_n\}$ in D . do.
 { for each input symbol 'a' in Σ .
 {
 $T = \text{set of states reachable from } x \text{ through 'a'}$.
 $y = \text{e-closure}(T)$.
 if $y \notin D$ then {
 add y to D
 unmark y
 }
 }
 mark x ;
 }



Algorithm for minimizing DFA.

- i. Construct $\Pi = \{\{\text{Final states}\}, \{\text{Non Final states}\}\}$
- ii. For each group G of Π do.
 {
 a) partition the group G such that
 }

two states 's' and 't' are allowed to be in the same group if transition for every I/P symbol should go to the same group of II.

b) place all the groups in Π_{new}

c) repeat step 2 till $\Pi = \Pi_{\text{new}}$.

lex:

Auxiliary definitions	i
Translation rules	ii
Auxiliary function	iii

i. define regEx.

 % %
 pattern {action}
 % %

e.g.: ex. 1

letter [a-z]

digit [0-9]

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

(% %

if | then | goto | for

{printf("Keyword");}

{id}

{printf("Identifier");}

{digit}+

{printf("Const");}

< | <= | > | >=

{printf("relOp");}

% %

\$ lex ex1.l → lex.yy.c.

\$ cc lex.yy.c -lfl → a.out

\$./a.out.

Parser

i. Top Down.

ii. Bottom Up.

Grammar, $G = (N, T, P, S)$.

$N \rightarrow$ set of non terminals. (A, B, C, \dots)

$T \rightarrow$ set of terminals (a, b, c, d, \dots)

$P \rightarrow$ set of productions (of form $\alpha \rightarrow \beta$)

$S \rightarrow$ start symbol ($S \in N$)

Language of a grammar, G.

$L(G) = \{ w \mid S \xrightarrow{*} w \text{ and } w \in T^* \}$

left most derivation.

$S \xrightarrow{*} w$ (in every step always replace
lm the left most non terminal).

Right most derivation.

$S \xrightarrow{*} w$ (" " " " the right most non terminal).

left-sentential form.

$S \xrightarrow{*} \alpha$
lm
 $\alpha \in (NUT)^*$

Regular Grammar.

$$\alpha \rightarrow \beta .$$

$$\alpha \in N^*$$

β is of the form ' aA ' or ' Aa ' or ' a '.

$$A \rightarrow aA/a \quad (a^*)$$

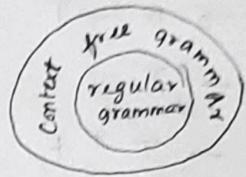
Context Free Grammar.

$$\alpha \rightarrow \beta$$

$$\alpha \in N^*$$

$$\beta \in (N \cup T)^*$$

expression grammar.



$$E \rightarrow E+E \quad T + bi$$

$$E \rightarrow E * E$$

$$E \rightarrow E - E$$

$$E \rightarrow E/E$$

$$E \rightarrow (E) / id.$$

$$N = \{E\}$$

$$T = \{+, -, *, /, (,), id\}$$

derive $id + id * id$.

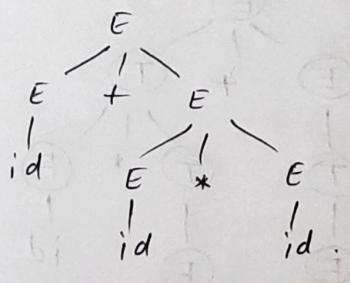
$$E \Rightarrow E + E$$

$$\Rightarrow E + E * E$$

$$\Rightarrow id + E * E$$

$$\Rightarrow id + id * E$$

$$\Rightarrow id + id * id$$



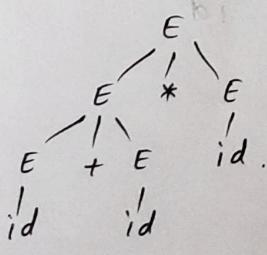
$$E \Rightarrow E * E$$

$$\Rightarrow E + E * E$$

$$\Rightarrow id + E * E$$

$$\Rightarrow id + id * E$$

$$\Rightarrow id + id * id$$



Two different derivation trees for same sentence
 \therefore ambiguous.
 introduce priority/
 precedence

F id ()

T * /

E + -

$E \rightarrow E + T$

$E \rightarrow E - T$

$E \rightarrow T$

$T \rightarrow T * F$

$T \rightarrow T / F$

$T \rightarrow F$

$F \rightarrow id$

$F \rightarrow (E)$

priority

$id + id * id$

$E \xrightarrow{lm} E + T$

$\xrightarrow{lm} T + T$

$\xrightarrow{lm} F + T$

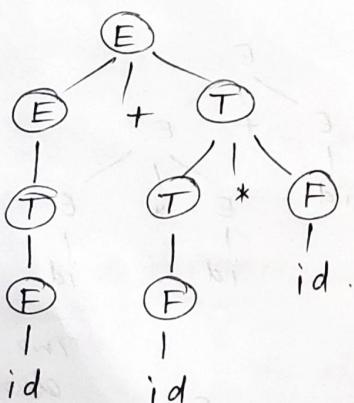
$\xrightarrow{lm} id + T$

$\xrightarrow{lm} id + T * F$

$\xrightarrow{lm} id + id * E$

$\xrightarrow{lm} id + id * id$

note: right most derivation also results in same tree.



If then-else grammar.

$$S \rightarrow i \underline{c} t s$$

$$S \rightarrow i c t s e s$$

$$C \rightarrow b$$

$$S \rightarrow a$$

$$N = \{ S, C \}$$

$$T = \{ i, t, e, a, b \}$$

ibtibtaea.

$$S \Rightarrow i \underline{c} t s m e s$$

①

$$\Rightarrow i b t \underline{s} e s$$

if then.

$$\Rightarrow i b t i \underline{c} t s e s$$

if then
else.

$$\Rightarrow i b t i b t \underline{s} e s$$

$$\Rightarrow i b t i b t a e s$$

$$\Rightarrow i b t i b t a e a$$

$$S \Rightarrow i \underline{c} t s$$

$$\Rightarrow i b t \underline{s}$$

②

if then ✓

$$\Rightarrow i b t i \underline{c} t s e s$$

if then
else.

$$\Rightarrow i b t i b t \underline{s} e s$$

(else must be
paired with the
closest unmatched if)

$$\Rightarrow i b t i b t a e s$$

$$\Rightarrow i b t i b t a e a$$

Two different Lm derivations
∴ ambiguous.

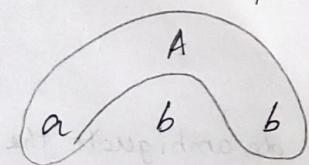
② is right Now how do we deambiguate the grammar.

$s \rightarrow U/M$ $M \rightarrow i c t M e M / a$ $U \rightarrow i c t M e U / i c t s.$ $C \rightarrow b$ ibtibtaea $s \Rightarrow U$ $U \Rightarrow i \underline{c} t M e U$ $\Rightarrow ibt \underline{M} e U$ $\Rightarrow ibt i c t M e M e U$

all possible parenthesis pairs.

 $s \rightarrow (s)s$ $s \rightarrow \epsilon$ Bottom up parser: (There is a reduction.
rhs is replaced with lhs). $s \rightarrow aACBe$ $A \rightarrow Ab/b$ $B \rightarrow d$

abbe de

The parser thinks the sentence
does not belong to the grammar,
but there is a derivation actually.

c d e.

Handle:
 substring of a given sentence and should be
 right hand side of a production and if it is
 reduced to LHS of a production, it should lead
 to start symbol.

The problem with bottom up parser is the
 difficulty to identify the handle.

Handle pruning Technique:

Right Most Derivation in reverse.

rmd of $id + id * id$.

$$E \xrightarrow{rm} E + \underline{E}$$

$$\xrightarrow{rm} E + E^* \underline{E}$$

$$\xrightarrow{rm} E + \underline{E}^* id$$

$$\xrightarrow{rm} \underline{E} + id^* id$$

$$\xrightarrow{rm} id + id^* id$$

right sentential form | Handle

Production used.

$id_1 + id_2 * id_3$

id_1

$E \rightarrow id$.

$E + id_2 * id_3$

id_2

$E \rightarrow id$.

$E + E^* id_3$

id_3

$E \rightarrow id$.

$E + E^* E$

$E^* E$

$E \rightarrow E^* E$

$E + E$

$E + E$

$E \rightarrow E + E$

E

Operator Precedence Parser.

Method - 2.

Three relationships defined b/w terrains

$\langle \cdot, \cdot \rangle, \doteq$

Terminals are operators.

- i. if θ_1 and θ_2 are operators and
 θ_1 has higher precedence than θ_2
then set
 $\theta_1 \rightarrow \theta_2$ and $\theta_2 \leftarrow \theta_1$.

- ii. if θ_1 and θ_2 are equal in precedence,
left association. - set $\theta_1 > \theta_2$
and
 $\theta_2 > \theta_1$.

- right association - set $\theta_1 < \theta_2$
and
 $\theta_2 < \theta_1$.

- iii. $\text{id} \rightarrow \theta$ $\theta < \text{id}$ (\doteq)

delimiter \$ < . @

	+	-	*	/	id	\$
+	→	→	←	←	←	→
-	→	→	←	←	←	→
*	→	→	→	→	←	→
/	→	→	→	→	←	→
id	→	→	→	→		→
\$	←	←	←	←	←	←

Stack	relation	input buffer
\$ \rightsquigarrow top most terminal	<	<u>id + id * id \$</u>
P < . id	>	<u>+ id * id \$</u>
\$ E	<	<u>+ id * id \$</u>
\$ < . E +	<	<u>id * id \$</u>
\$ < . E + < . id + handle	>	<u>* id \$</u>
\$ < . E + E	<	<u>* id \$</u>
\$ < . E + < . E *	<	<u>id \$</u>
\$ < . E + < . E * < . id	>	<u>\$</u>
\$ < . E + < . E * E	>	<u>\$</u>

< or = : shift

> : handle at top of stack.

string between < and > is the handle

\$ < . E + E	> >	\$
\$ E	accept.	\$

This is called shift reduce parser.

Reduce the foll.

i. id - id + id

ii. id * id * id - id

stack	relation	input buffer
\$	<	id - id + id \$
\$ < . id	<	- id + id \$
\$ E	<	id + id \$
\$ < . E -	>	+ id \$
\$ < . E - < . id	>	+ id \$
\$ < . E - E	<	+ id \$
\$ E	<	id \$
\$ < . E +	>	\$
\$ < . E + < . id	>	\$
\$ < . E + E	>	\$
\$ E	accept	

ii. stack	relation	input buffer
\$	<	id * id * id - id \$
\$ < . id	>	* id * id - id \$
\$ E	<	* id * id - id \$
\$ < . E *	<	id * id - id \$
\$ < . E * < . id	>	* id - id \$
\$ < . E * E	>	* id - id \$
\$ E	<	id - id \$
\$ < . E *	<	- id \$
\$ < . E * < . id	>	- id \$
\$ < . E * E	>	- id \$
\$ E	<	- id \$

\$ < . E -	<	id
\$ < . E - < . id	>	\$
\$ < . E - E	>	\$
\$ E	accept.	\$

Method - 2

using operator precedence grammar.

- operator grammar. (No two consecutive non terminals allowed , No ϵ productions).
- relations $<$ \rightarrow $=$ are disjoint.

Leading (A)

A - Non Terminal.

1. If there is a production of the form
 $A \rightarrow r a \delta$, where

r is a single non terminal or ϵ
 then 'a' is in Leading (A) .

" LEADING (A) is also sometimes referred to
 as FIRST OF (A) "

2. If there is a production of the form

$$A \rightarrow B \alpha ,$$

all LEADING (B) are in LEADING (A) .

$$E \rightarrow E + T \mid T$$

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

$$F \rightarrow (E) \mid id.$$

	LEADING
F	(, id
T	, *, (, id
E	+ , *, (, id

TRAILING (A)

A - non

terminal.

1. If there is a production of the form
 $A \rightarrow aB$, where

B is a single non terminal or
 then ' a ' is in TRAILING (A).

2. If there is a production of the form
 $A \rightarrow \alpha B$ then

All trailing (B) are in TRAILING (A).

TRAILING	
F), id.
T	*,), id
E	+ , *,), id.

$$S \rightarrow aAcBe$$

$$A \rightarrow Ab/b$$

$$B \rightarrow d.$$

	LEADING	TRAILING
S	a	e
A	b	b.
B	d.	d.

$$S \rightarrow (L) / a *$$

$$L \rightarrow L, S / S$$

	LEADING	TRAILING
S	'(, 'a')	')', 'a'
L	'', 'c', 'a'	', ', ')!', 'a'

"The method"

1. $a \equiv b$, if there is an RHS of a production of the form $\alpha a \beta b \gamma$, where β is a single non terminal or ϵ .

2. $a < b$, if there is an RHS of a production of the form $\alpha a A \beta$, $\forall b$ in $\text{LEADING}(A)$, $a < \text{LEADING}(A)$.

3. $a > b$, if there is an RHS of a production of the form $\alpha A a \beta$, $\forall b$ in $\text{trailing}(A)$, $\text{TRAILING}(A) \rightarrow a$

4. $\$ \leftarrow \text{LEADING}(S)$.

$\text{TRAILING}(S) \rightarrow \$$

S - start symbol.

$$E \rightarrow E^+ T \mid T$$

$$T \rightarrow T^* F \mid F$$

$$F \rightarrow (E) \mid id$$

	LEADING	TRAILING
E	$+ * (id$	$+ *) id$
T	$* (id$	$*) id$
F	$(id$	$) id$

	+	*	c	>	id	\$
+	$\cdot >$	$<.$	$<.$	\rightarrow	$<.$	\rightarrow
*	$\cdot >$	$\cdot >$	$<.$	$\cdot >$	$<.$	\rightarrow
c	$\cdot <$	$\cdot <$	$<.$	\doteq	$<.$	
>	$\cdot >$	$\cdot >$		$\cdot >$		$\cdot >$
id	$\cdot >$	$\cdot >$		$\cdot >$		$\cdot >$
\$	$<.$	$<.$	$<.$		$<.$	

	a	b	c	d	e	f
a	<	\doteq				
b	\gt	\gt				
c			\lt	\doteq		
d					\gt	
e						\gt
f	\lt					

	()	a	,	\$
(<	=	<	<	
)		>		>	>
a		>		>	>
,	<	>	<	>	
\$	<		<		

Top-down parser:

Elimination of left recursion:

$$A \rightarrow A\alpha \mid \beta$$



$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' \mid \epsilon$$

In general,

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid A\alpha_3 \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \beta_3 \dots \mid \beta_n$$



$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \beta_3 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \alpha_3 A' \mid \dots \mid \alpha_m A' \mid \epsilon$$

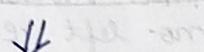
i. $A \rightarrow Aa \mid Ad \mid e \mid f \dots$

$$A \rightarrow eA' \mid fA'$$

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

ii. $s \rightarrow Aa \mid b$

$$A \rightarrow Aa \mid Sd \mid e$$



$$A \rightarrow Aa \mid Aad \mid bd \mid e$$

$$A \rightarrow eA' \mid bdA'$$

$$A' \rightarrow cA' \mid adA' \mid \epsilon$$

iii. $E \rightarrow E + T \mid E - T \mid T$

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

$F \rightarrow (E) \mid id.$

$E \rightarrow TE'$

$E' \rightarrow +TE' \mid -TE' \mid \epsilon$

$T \rightarrow FT'$

$T' \rightarrow *FT' \mid /FT' \mid \epsilon$

$F \rightarrow (E) \mid id$

Left factoring.

$A \rightarrow \alpha\beta \mid \alpha\gamma$

α : The common prefix.

eliminate it using.

$A \rightarrow \alpha A'$

$A' \rightarrow \beta \mid \gamma$

Recursive Descent Parser:

Procedure for every non terminal.

sample code.

```
char token[50];  
int i;
```

```
int main() {  
    cin >> token;
```

Remember:

grammar should have
no left recursion
and no common prefix

```
i = 0 ;  
E() ;  
if ( i == strlen(token) ) cout << "ACCEPTED" ;  
else cout << "NOT ACCEPTED" ;  
}  
void E()  
{  
    T();  
    EPRIME();  
}  
void EPRIME()  
{  
    if ( token[i] == '+' )  
    {  
        i++ ;  
        T();  
        EPRIME();  
    }  
    else if ( token[i] == '-' )  
    {  
        i++ ;  
        T();  
        EPRIME();  
    }  
    else  
    {  
        //  $E' \rightarrow E$ .  
    }  
}  
void T()  
{  
    F();  
    TPRIME();  
}
```

```
void TPRIME()
{
    if (token[i] == '*')
    {
        i++;
        F();
        TPRIME();
    }
    else if (token[i] == '/')
    {
        i++;
        F();
        TPRIME();
    }
    else
    {
        // T' → ε
    }
}

void F()
{
    if (token[i] == '(')
    {
        i++;
        E();
        if (token[i] == ')')
        {
            i++;
        }
        else
        {
            cout << "ERROR" << i;
            exit(0);
        }
    }
}
```

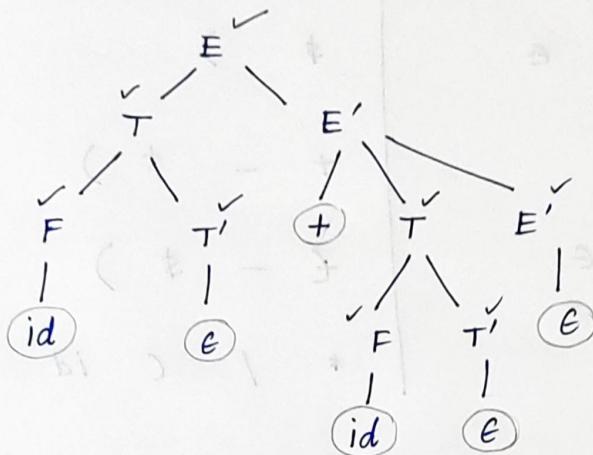
```

else if ( token[id] == 'i' )
{
    i++;
}
else
{
    cout << "ERROR" << i;
    exit(0);
}

```

3 ~~WILL NOT~~ ~~YES~~

id + id .



Predictive Parser :

FIRST(X)

1. if X is a terminal then $\text{FIRST}(X) = \{X\}$
2. if X is a non-terminal and there is a production $X \rightarrow a\alpha$, then 'a' is added to $\text{FIRST}(X)$.
3. if X is a non-terminal and there is a production $X \rightarrow Y_1 Y_2 Y_3 \dots Y_k$ Y_i - Non terminal $\text{FIRST}(Y_i)$ added to $\text{FIRST}(X)$. and

if $\text{FIRST}(Y_1)$ contains ϵ

$\text{FIRST}(Y_2)$ is added to $\text{FIRST}(x)$.

if $\text{FIRST}(Y_2)$ contains ϵ

$\text{FIRST}(Y_3)$

	FIRST	FOLLOW
E	(id	\$)
E'	+ - ϵ	\$)
T	(id	+ - \$)
T'	* / ϵ	+ - \$)
F	(id	* / (id

FOLLOW(A).

A - non terminal

1. \$ is added to FOLLOW(s), s - start symbol.
2. if there is a production of the form:
 $A \rightarrow \alpha B \beta$.
then $\text{FIRST}(\beta)$ ^{except ϵ} is added to FOLLOW(B).
- if $\text{FIRST}(\beta)$ contains ϵ , then add all FOLLOW(A) to FOLLOW(B).

3. If there is a production $A \rightarrow \alpha B$, add
all FOLLOW(A) to FOLLOW(B).