

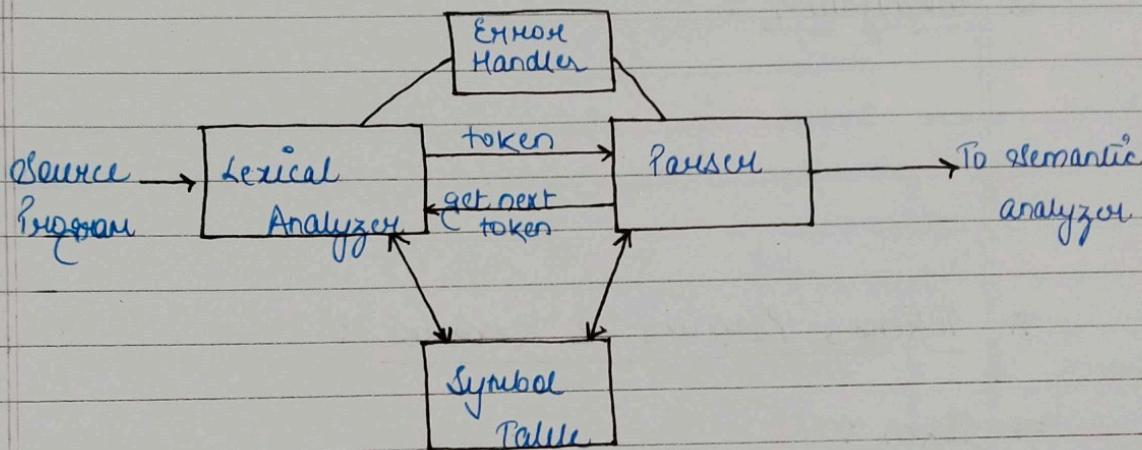
14th Jan '17.
Lecture-5

Unit-1

Date: / /
Page: / /

Lexical Analysis

Flow of lexical analyzer:



- removes white spaces
- if there are more than 1 white spaces then, LA does compaction; i.e; if there are 3 blank spaces then LA will compact it & make it 1 blank space.
- removes the comment lines.

* Generally lexical analyzers are divided into cascade of 2 processes :

1 Scanning :

consist of simple processes that do not require tokenization of inputs such as deletion of comments, compaction of white spaces in one

2 Tokenization :

portion

LA ~~process~~ is the more complex portion which produces the tokens from the output of scanner.

* Tokens, Patterns & Lexemes :

⇒ Token :

A token is a pair consisting of a token name & an optional attribute value.

The token name is an abstract symbol representing a kind of lexical unit.

Example : a particular keyword or a sequence of input characters denoting an identifier. ~~identifiers~~ identifiers, keywords, constants, literals, punctuations, tokens will be generated.

For operators tokens will be generated but ~~then~~ ~~any~~ every will not be mentioned in symbol table.

Date: _____
Page: _____

int a;

3 tokens for "int", "a", ";"

printf("y. d")
① ② ③ ④ ⑤ ⑥ ⑦ ;

8 tokens for "printf", "(", "'", "y. d", "'", ")", ":", ";"

2 & a ;

3 tokens.

17th Jan 17

Lecture 6

⇒ Pattern:

A pattern is a description of a form that the lexemes of a token may take.

In the case of a keyword as a token the pattern is just the sequence of characters that form the keyword.

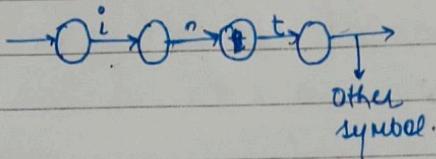
⇒ Lexeme:

A lexeme is a sequence of characters in the source program that match the pattern for a token & is identified by the lexical analyzer as an instance of that token.

Example: i > = j ; i : j ; 8 tokens

if i > j ; i > = j ; 3 tokens 3 tokens

int a = 100; 5 tokens



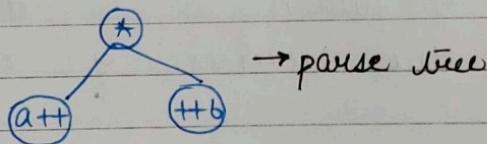
a++ + + + b; 6 tokens

no lexical error

but syntax error (parse tree can't be made)

a++ * + + b 5 tokens

no lexical & syntax error.

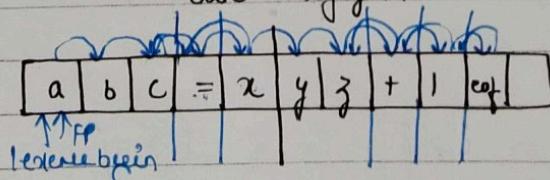


Input buffering : to speed up the slow reading of source prog. for tokenization.

We will introduce a buffer scheme that handles large look ahead safely. We then consider an implement involving sentinels that saves time checking for the end of buffer.

1. Buffer pairs :

$$abc = xyz + 1$$



18th Jan '17

Lecture 7

- Each buffer is of the same size N , and N is usually the size of a disk block. Using one system read command we can read N characters into a buffer.
 - If fewer than N characters remain in the input file then a special character represented by EOF marks the end of file & it is different from any possible character of the source program.
 - α pointers to the input are maintained:
- Pointer lexemeBegin: Marks the beginning of the current lexeme whose extent we are attempting to determine.
- Pointer forward: Scans ahead until a pattern match is found.
- Once the next lexeme is determined, forward is set to the character at its right end then after the lexeme is recorded as an attribute value of a token returned to the parser. lexemeBegin is set to the character immediately after the lexeme is just found.

- Advancing forward requires that we first test whether we have reached the end of one of the buffers & if so we must reload the other buffer from input & move forward to the beginning of the newly loaded buffer.
- * Code to advance the forward pointer :

{ if (FP is at end of the 1st half)

 reload the end half;

 FP = FP + 1;

}

use if [FP is at end of the end half]

{

 reload the 1st half;

 move the FP to the beginning of the 1st half;

 FP = FP + 1;

}

 FP + 1;

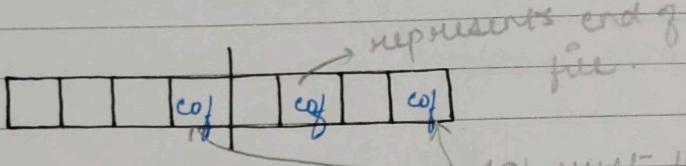
 → in this FP matches

 view of where as
 in previous FP match
 with size.

II. Buffer pair with sentinels :

- In the previous method, for each character read we make 2 tests, one for the end of the buffer and other to determine what character is read.

- We can combine the buffer end test with the test for the current character if we extend each buffer to hold a sentinel character at the end. Sentinel is a special character that can not be part of the source program & the natural choice is the character `eof`.



- * Code to advance ~~the~~ FP:

if ($FP == eof$)

if (FP is at the end of 1st half)

reload the 2nd half;

$FP = FP + 1;$

?

else if (FP is at the end of 2nd half)

?

reload 1st half;

move to the beginning of the 1st half;

?

else

terminate lexical analysis;

?

Q. Find number of tokens in : `int a[4][5];`

for 1D arrays:

`a[i]` can be written as $*(a+i)$

for 2D arrays:

`a[i][j]` can be written as $*(*(a+i)+j)$

in `int a[4][5];`

Replace with above form

`int | * | | * | (a | + | 4) | + | 5 | ;` \Rightarrow 13 tokens.

19th Jan '17

Lecture - 8

Specification of tokens:

* Formal notations for RE:

1. Strings & Languages

2. Operations on languages:

a.) Union

b.) Concatenation

c.) Closure (Kleen's & Positive)

Example: Let the set of letters :

$$L = \{ A, B, C, \dots, Z, a, b, c, \dots, z \}$$

Let the set of digits :

$$D = \{ 0, 1, \dots, 9 \}$$

LUD = language with 62 strings of length 1.

LD = set of 520 strings of length 2 each consisting of 1 letter followed by 1 digit.

L^4 = set of all 4 letter strings.

L^* = set of all strings of letters including null.

$L(LUD)^*$ = set of all strings of letters & digits but beginning with letter.
It is an identifier.

D^+ = set of all strings of 1 or more digit, null is not present.

3. Regular Expressions :

Rules that define the regular expressions over some alphabet Σ and language that those expressions denote.

a) Basis:

2 rules that form the basis:

- "Null" is a regular expression and

$$L(\epsilon) = \{ \epsilon \}.$$

- if "a" is a symbol in Σ then, a is regular expression and $L(a) = \{ a \}$.

It means language of one string of length one with a in its position.

b) Induction rule:

Suppose r and s are regular expression denoting languages $L(r) \subseteq L(s)$.

r, s is regular expression denoting language $L(r) \cup L(s)$.

r, s is a regular expression denoting language $L(r) L(s)$

r^* is a regular expression denoting $L(r)^*$.

r is a regular expression denoting $L(r)$.

4. Regular definitions :

We may wish give names to certain RE and use those names in subsequent expressions as if the names were themselves symbols.

If Σ is an alphabet of basic symbols then a regular definition is a sequence of definitions of the form :

$$\begin{aligned} d_1 &\rightarrow r_1 \Rightarrow \Sigma \\ d_2 &\rightarrow r_2 \Rightarrow \Sigma \cup d_1 \\ d_3 &\rightarrow r_3 \\ &\vdots \\ &\vdots \\ d_n &\rightarrow r_n \end{aligned}$$

where,

each r_i is a RE over $\Sigma \cup \{d_1 \cup d_2 \cup \dots \cup d_{i-1}\}$

Example : Give a regular definition for identifiers.

$$\text{identifier} \rightarrow \text{letter} (\text{letter/digit})^*$$

$$\text{letter} \rightarrow a/b/c/\dots/\beta/A/B/C/\dots/\zeta$$

$$\text{digit} \rightarrow 0/1/\dots/9/$$

20th Jan '17

Lecture-9 Example: Give the regular definition for unsigned numbers (int & float).

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

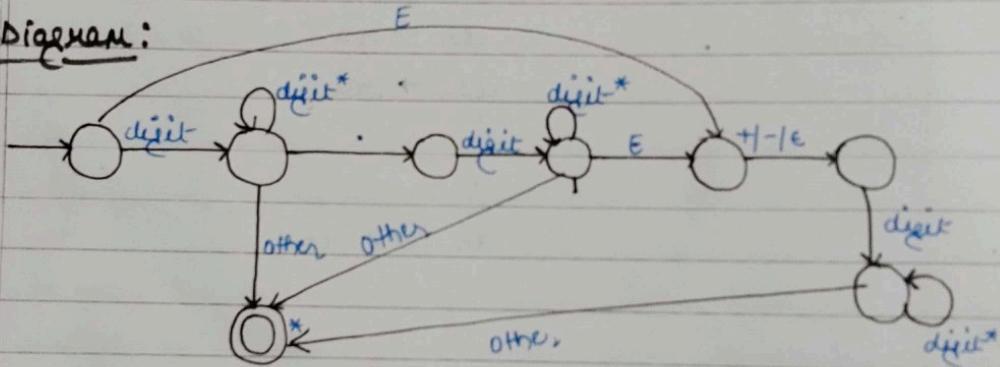
digits \rightarrow digit (digit)*

optional fraction $\rightarrow .\ digits | \epsilon$

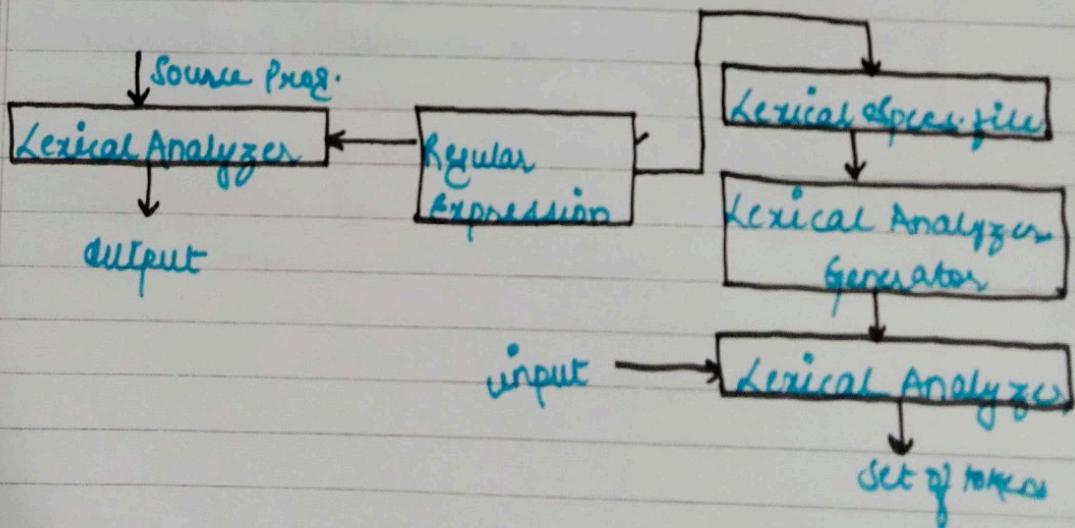
optional exponent $\rightarrow (E (+|-|E) digits) | \epsilon$

number \rightarrow digits optional fraction optional exponent

Transition Diagram:



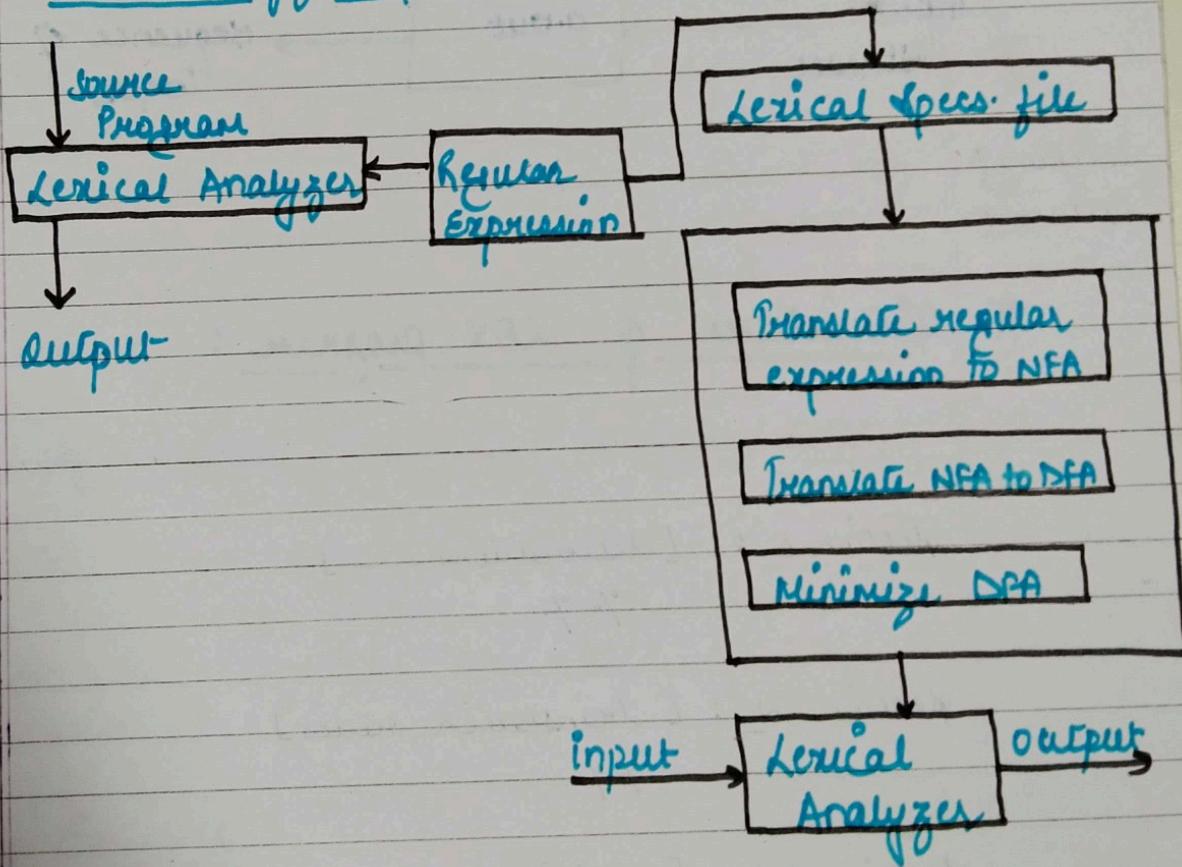
Elements of Lexical Analyzer 8



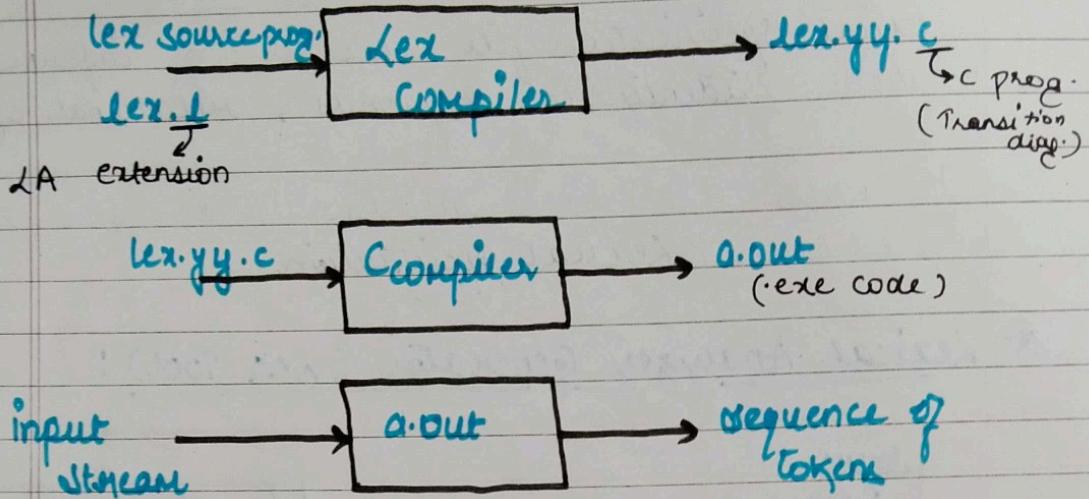
Lexical Specification file consists of information about identifiers, keywords, relational operators.

The rules of basic building blocks for a language are called Lexical specification.

* Lexical Analyzer Generator (Lex Tool):



* Use of Lex :



21st Jan '17
Lecture - 10

#. Structure of LEX program :

A LEX program is in following form:

Section 1 : [declarations ...]
% %

Section 2 : [translation rules]
% %

Section 3 : [auxiliary functions/ subroutines]

Section 1 : [declaration ...]

This section includes variables, constants & regular definition.

Section 2 : [Translation rules]

The translation rules each have the form pattern and corresponding action.

Each pattern is a RE which may use the regular definitions of the declaration section.

The actions are the fragments of code written in C.

Section 3 : [auxiliary functions/ subroutines]

It holds whatever additional functions are used in the actions. Alternatively, these functions can be compiled separately & can be loaded with LA.

Example :

```
% {
    #include <stdio.h>
    int count;
} // section 1

%?
letter [A-Z a-z]
digit [0-9]
%.
```

```
IF { printf("It is a KEYWORD %s", yytext); }
FOR { printf("It is a KEYWORD %s", yytext); }
while { printf("It is a KEYWORD %s", yytext); }

{ letter? (letter? | digit?)* printf("Identifier %s", yytext); }

%.
```

```
main()
{
    yytext();
} // section 2
```

Q. Write a lex program to find the longest word in the input file & print the length of word.

Date: _____
Page: _____

%.%

#include <stdio.h>

int k=0;

%.%

%.%

+ → occurrences

a-z, A-Z

can be more

than 1

a is not inc.

[a-z A-Z] +

{

if (yylen > k)

k=yylen;

{

%.%

int main (int argc[], char **argv[])

{

yyin = fopen("abc.txt", "r");

yyerr();

printf("Largest : %d", k);

printf("\n");

return 0;

{

24th Jan '17
Lecture 11

Date: ___/___/___
Page: ___

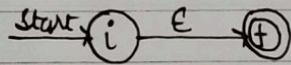
Thompson's Algorithm:

Input: A RE Σ over alphabet Σ .

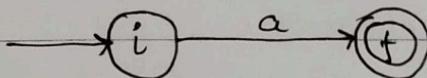
Output: An NFA 'N' accepting $L(\Sigma)$.

Method/Procedure:

a) Basis: For expression ϵ construct the NFA

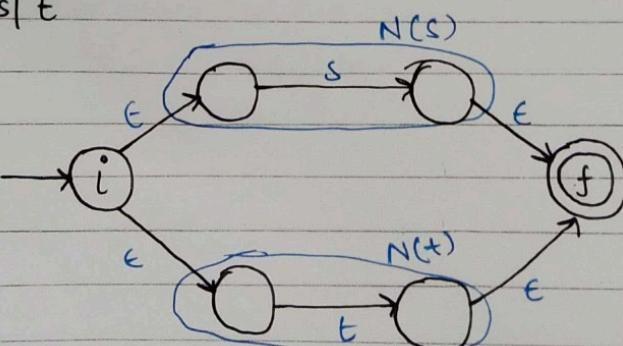


For an subexpression ~~is~~ 'a' in Σ .

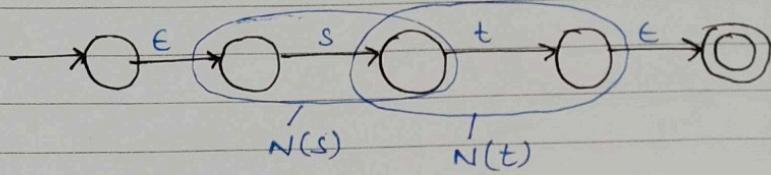


b) Inductive: Suppose $N(s)$ & $N(t)$ are NFAs for REs s & t resp.

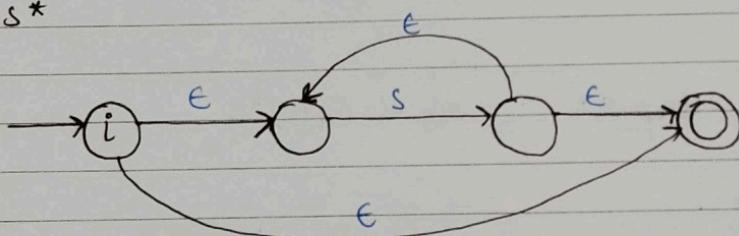
i) $g = s \mid t$



ii) $M = st$



iii) $M = s^*$

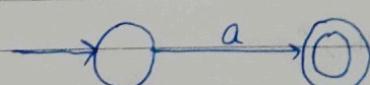


iv) Suppose $M = (s)$ then $\delta(M) = \delta(s)$ and then we can use the NFA $N(s)$ as $N(M)$.

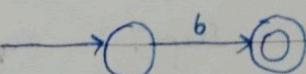
Ques. Use Thompson's construction algo. to construct NFA for given RE :

$$M = (a|b)^*abb$$

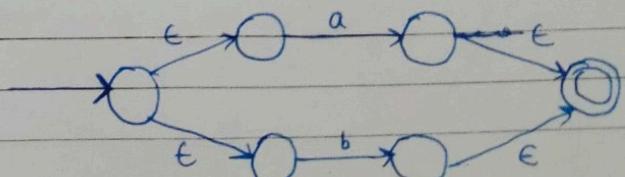
Step-1 $M_1 = a$



Step-2: $M_2 = b$

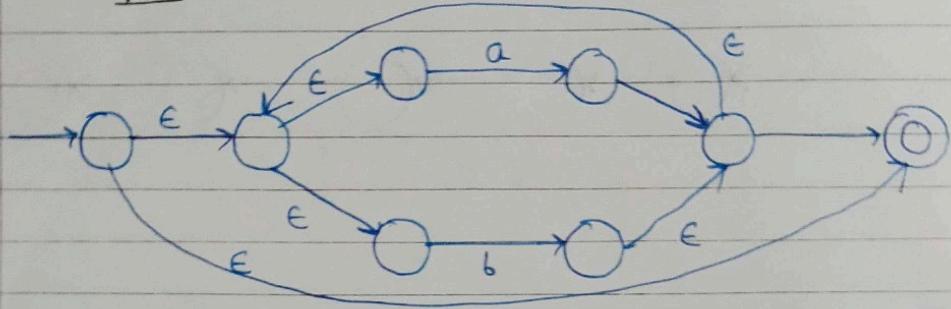


Step-3: $M_3 = a/b$

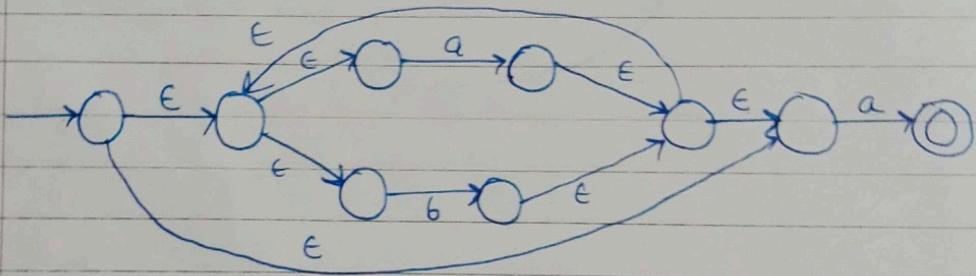


NFA

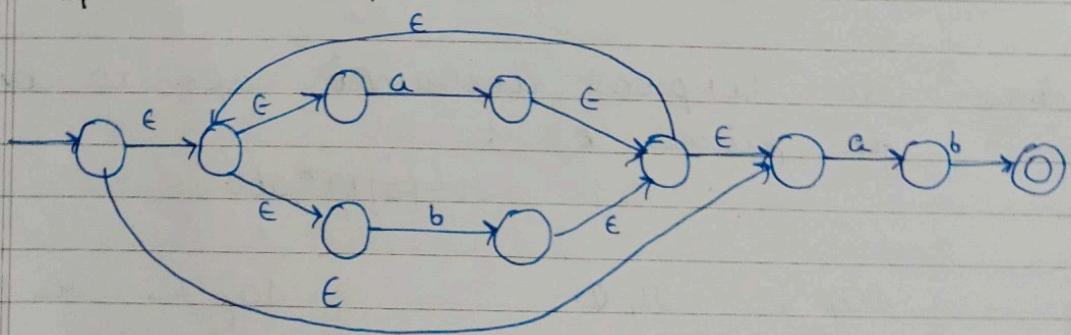
Step-4 : $g_4 = (a/b)^*$



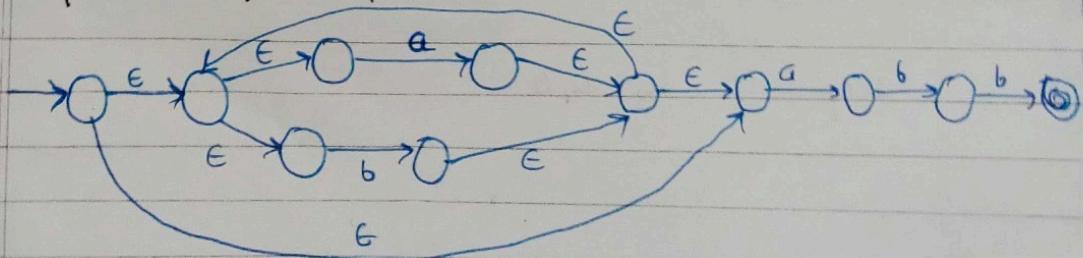
Step-5 : $g_5 = (a/b)^* a$



Step-6 : $g_6 = (a/b)^* ab$



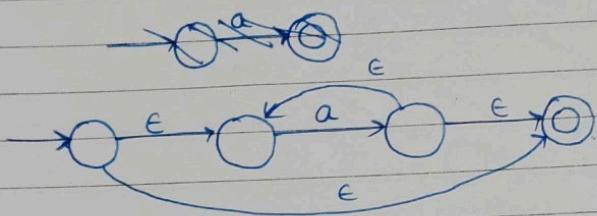
Step-7 : $g_7 = (a/b)^* abb$



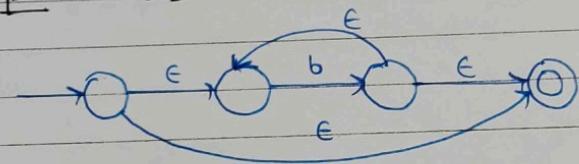
$$M = (a^* \mid b^*)^*$$

$$\underline{\text{Step 1}} : M_1 = a^*$$

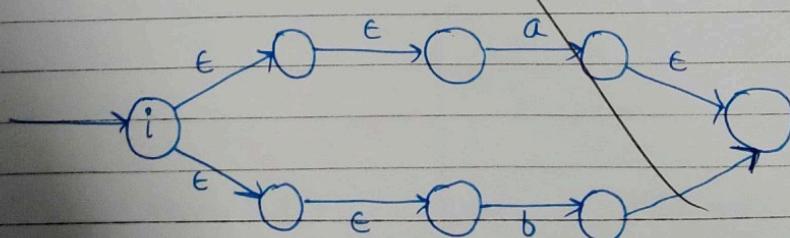
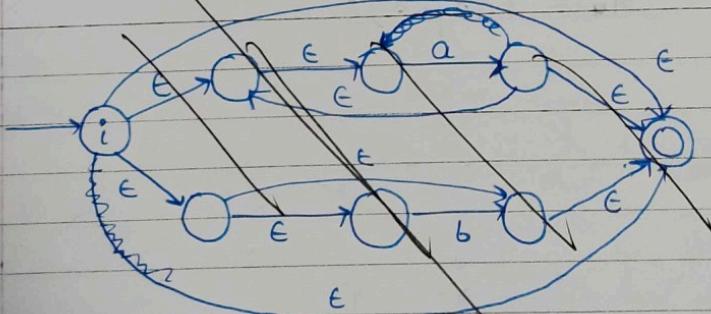
Step 2



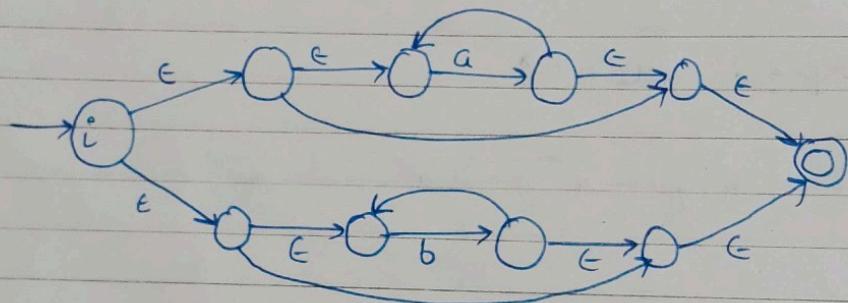
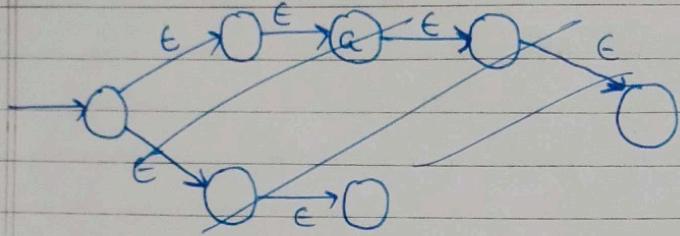
$$\underline{\text{Step 2}} : M_2 = b^*$$



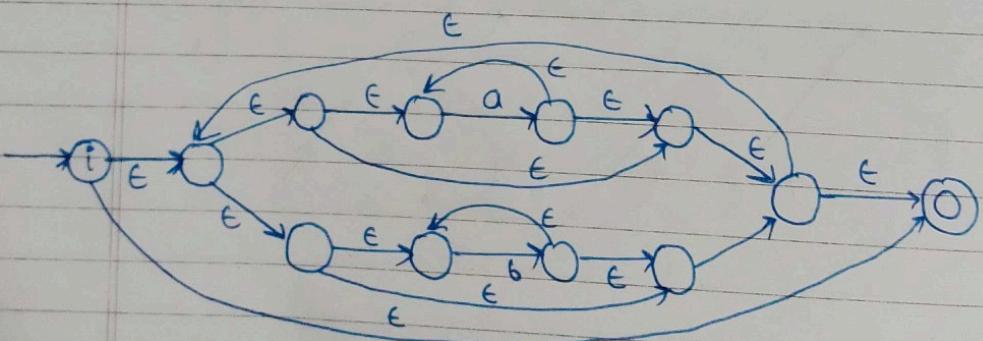
$$\underline{\text{Step 3}} : M_3 = (a^* \mid b^*)^*$$



Step-3 $R_3 = (a^* \mid b^*)$



Step-4 : $R_4 = (a^* \mid b^*)^*$



Conversion of NFA to DFA :

(Subset construction algorithm)

Input : A NFA N .

Output : A DFA D accepting the same language as N .

Method : our algorithm constructs a transition table D_{tran} for D . Each state q of D is a set of NFA states and we construct D_{tran} .

If E is a single state of N while 'T' is the set of states of N .

25th Jan '17
Lecture 12

Operation

Description

- | <u>Operation</u> | <u>Description</u> |
|----------------------------|---|
| ϵ -closure(s) | Set of NFA states reachable from NFA state s on ϵ transitions alone. |
| ϵ -closure(T) | Set of NFA states reachable from some NFA state s in set T on ϵ transitions alone. |
| $move(T, a)$ | Set of NFA states to which there is a transition on i/p symbol a from some state s in T . |

Operation

ϵ closure (S_0)

Description

The start state of DFA is one ϵ closure so. The accepting state of DFA are all these set of NFA's state that include atleast one accepting state of N.

Initially ϵ closure S_0 is the only state in Dstates & it is unmarked.

* Algorithm:

while (there is an unmarked state T in Dstates)

{

mark T;

for (each input symbol 'a')

{

$U = \epsilon$ -closure (move (T, a));

if (U is not in Dstates)

add U as an unmarked state of Dstates;

$D_{trans}[T, a] = U$;

}

}

Q

$SL = (a/b)^* a b b$

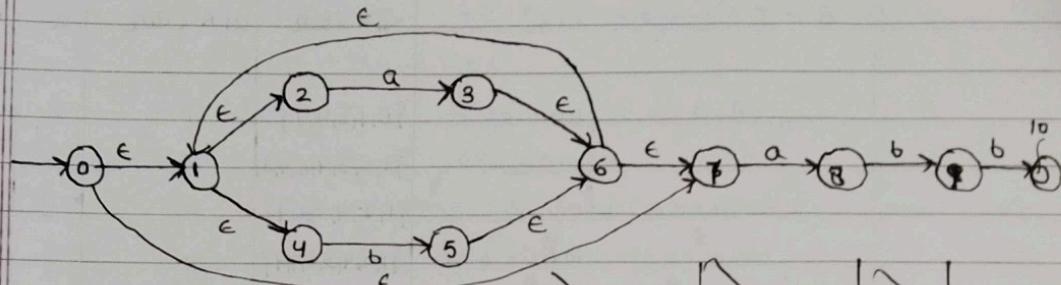
FA is
if state
g N
at least

the
e

states

to:

Date _____
Page _____



$$A = \text{E-closure}(0)$$

$$D\text{States} = A = \{0, 1, 2, 4, 7\}$$

NFA state	DFA state	a	b
{0, 1, 2, 4, 7}	A		
{1, 2, 3, 4, 6, 7, 8}	B		
{1, 2, 4, 5, 6, 7}	C		

$$\begin{aligned}D\text{trans}[A, a] &= \text{E-closure}(\text{move}(A, a)) \\&= \text{E-closure}(\text{move}(0, a), \text{move}(1, a), \text{move}(2, a), \text{move}(3, a), \text{move}(4, a), \text{move}(5, a), \text{move}(6, a), \text{move}(7, a)) \\&= \text{E-closure}(3, 8) \\&= \text{E-closure}(3) \cup \text{E-closure}(8) \\&= \{3, 6, 7, 1, 2, 4\} \cup \{8\} \\D\text{trans}[A, a] &= \{1, 2, 3, 4, 6, 7, 8\} = B\end{aligned}$$

$$\begin{aligned}D\text{trans}[A, b] &= \text{E-closure}(\text{move}(A, b)) \\&= \text{E-closure}(\text{move}(0, b), \text{move}(1, b), \text{move}(2, b), \text{move}(3, b), \text{move}(4, b), \text{move}(5, b)) \\&= \text{E-closure}(5) \\&= \{5, 6, 7, 1, 2, 4\} \\D\text{trans}[A, b] &= \{1, 2, 4, 5, 6, 7\} = C\end{aligned}$$

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

$$\begin{aligned}
 D_{\text{tran}}[B, a] &= \text{E-closure}(\text{move}(B, a)) \\
 &= \text{E-closure}(\text{move}(1, a), \\
 &\quad \text{move}(2, a), \text{move}(3, a), \\
 &\quad \text{move}(4, a), \text{move}(6, a), \\
 &\quad \text{move}(7, a), \text{move}(8, a)) \\
 &= \text{E-closure}(3) \cup \text{E-closure}(8) \\
 &= \{1, 2, 3, 4, 6, 7, 8\}
 \end{aligned}$$

NFA state	DFA state	a	b
$\{0, 1, 2, 4, 7\}$	A	B	C
$\{1, 2, 3, 4, 6, 7, 8\}$	B	B	D
$\{1, 2, 4, 5, 6, 7\}$	C	B	C
$\{1, 2, 4, 5, 6, 7, 9\}$	D	B	E
$\{1, 2, 4, 5, 6, 7, 10\}$	(E)	B	C

$$\begin{aligned}
 D_{\text{tran}}[B, b] &= \text{E-closure}(\text{move}(B, b)) \\
 &= \text{E-closure}(\text{move}(1, b), \text{move}(2, b), \\
 &\quad \text{move}(3, b), \text{move}(4, b), \text{move}(6, b), \\
 &\quad \text{move}(7, b), \text{move}(8, b)) \\
 &= \text{E-closure}(5) \\
 &= \text{E-closure}(5) \cup \text{E-closure}(9) \\
 &= \{5, 6, 7, 1, 2, 4\} \cup \{9\}
 \end{aligned}$$

$$D_{\text{tran}}[B, b] = \{1, 2, 4, 5, 6, 7, 9\} = D$$

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

$$\begin{aligned}
 D_{\text{tran}}[C, a] &= \text{E-closure}(\text{move}(C, a)) \\
 &= \text{E-closure}(\text{move}(1, a), \text{move}(2, a), \\
 &\quad \text{move}(3, a), \text{move}(4, a), \text{move}(5, a), \\
 &\quad \text{move}(6, a), \text{move}(7, a)) \\
 &= \text{E-closure}(3, 8)
 \end{aligned}$$

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

$$\begin{aligned}
 D_{\text{tran}}[C, b] &= \text{E-closure}(\text{move}(C, b)) \\
 &= \text{E-closure}(\text{move}(1, b), \text{move}(2, b), \\
 &\quad \text{move}(4, b), \text{move}(5, b), \text{move}(6, b), \\
 &\quad \text{move}(7, b))
 \end{aligned}$$

$$= \text{E-closure}(5)$$

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

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

$$\begin{aligned}
 D_{\text{tran}}[D, a] &= \text{E-closure}(\text{move}(D, a)) \\
 &= \text{E-closure}(\text{move}(1, a), \text{move}(2, a), \\
 &\quad \text{move}(4, a), \text{move}(5, a), \text{move}(6, a), \\
 &\quad \text{move}(7, a), \text{move}(9, a))
 \end{aligned}$$

$$D_{tran}[D, a] = \text{Eclosure}(\{3, a\})$$

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

$$D_{tran}[D, b] = \text{Eclosure}(\{a, b\})$$

$$= \text{Eclosure}(\{a, b\}, \{a, b\})$$

$$= \text{Eclosure}(\cancel{a, b}, 5, 10)$$

$$= \text{Eclosure}(5) \cup \text{Eclosure}(10).$$

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

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

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

$$D_{tran}[E, a] = \text{Eclosure}(\{a, E\})$$

$$= \text{Eclosure}(\cancel{a, a}, \{a, a\}, \{a, a\})$$

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

$$= B$$

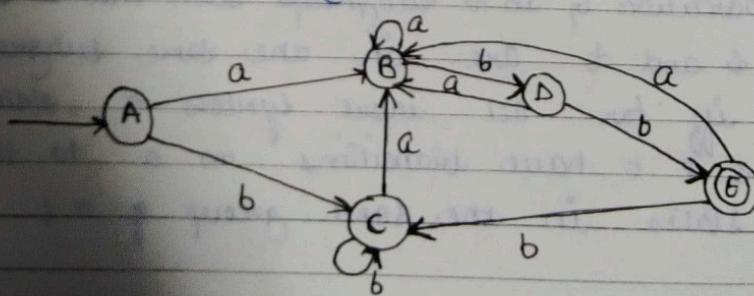
$$D_{tran}[E, b] = \text{Eclosure}(\{E, b\})$$

$$= \text{Eclosure}(\{b, b\}, \{b, b\})$$

$$= \text{Eclosure}(5)$$

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

$$= C$$



20th Jan '17

Lecture-13

Minimization of DFA (States):

Input : A DFA D with set of states S ,
input alphabet Σ , start state s_0 &
set of accepting states F .

Output : A DFA D' accepting the same language as
 D and having as few states as
possible.

Method :

Step-1 : Start with an initial partition π with
2 groups F & $S-F$, the accepting &
non-accepting states of D .

Step-2 : Apply the following procedure to construct
a new partition π_{new} .

Initially let, $\pi_{\text{new}} = \pi$
for (each group G of π)

partition G into subgroups such that 2 states
 s and t are in the same subgroup
iff for all input symbols a , states
 s & t have transitions on a to the
states in the same group of π ;

Replace G in π_{new} by the set of

all subgroups formed;

Step-3:

if ($\pi_{\text{new}} = \pi$), let $\pi_{\text{final}} = \pi$ and continue with step-4. otherwise repeat step-2 with π_{new} in place of π .

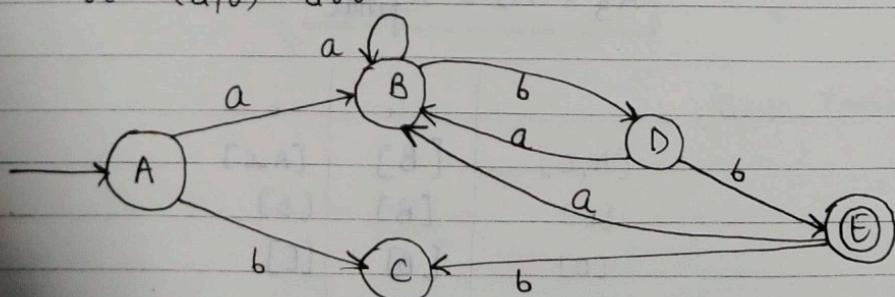
Step-4:

choose one state in each group of π_{final} as the representative for that group. The representatives will be the states of the minimum state DFA D' . The other components of D' are constructed as follows:

- the start state of D' is the representative of the group containing the start state of D .
- the accepting states of D' are the representatives of those groups that contain an accepting state of D .

Q

$$M = (a/b)^* abb$$

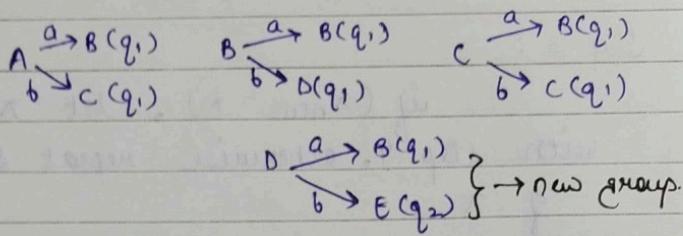


Date: ___/___/___
Page: 93

$$F = \{E\}$$

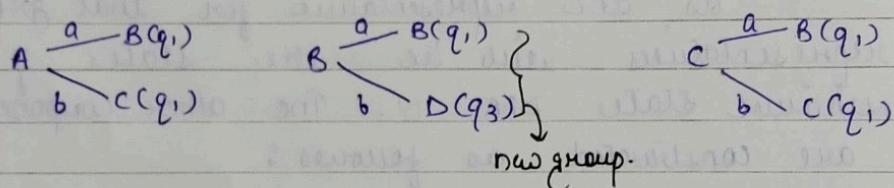
$$S - F = \{A, B, C, D\}$$

$$\pi_0 = \{A, B, C, D\}, \{E\}$$



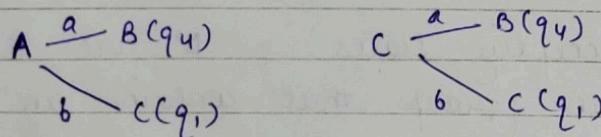
$$\pi_1 = \{A, B, C\}, \{D\}, \{E\}$$

$$q_1 \quad q_3 \quad q_2$$



$$\pi_2 = \{A, C\}, \{B\}, \{D\}, \{E\}$$

$$q_1 \quad q_4 \quad q_3 \quad q_2$$

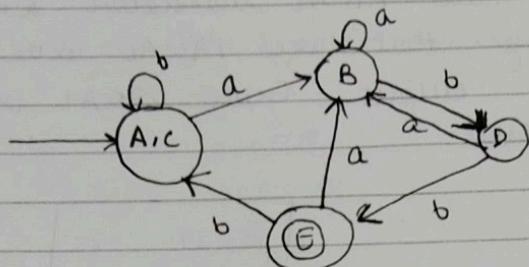


$$\pi_3 = \pi_2 = \pi_{\text{final}}$$

Now;

	a	b
[A, C]	[B]	[A, C]
[B]	[B]	[D]
[D]	[B]	[E]
* [E]	[B]	[A, C]

1st Feb '14
Lecture 1

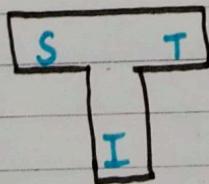


1st Feb '17

Lecture-14

Bootstrapping:

- The process by which simple language is used to translate more complicated program which in turn may handle an even more complicated program is known as Bootstrapping.
- The compiler is characterized by three languages:
- Source language (S)
- Target language (T)
- Implementation language (I).

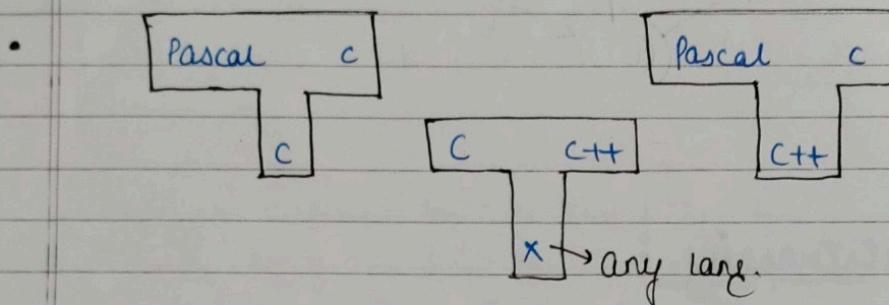


S T
I

fig: Textual form

fig: Power T-Diagram

Example: we have PASCAL translator written in C language that takes PASCAL code & produces C as output. Create PASCAL translator in C++ for same.



Pascal c + ~~Pascal~~ C++ = Pascal c++
 C J C++

• compilers are of 2 types :

I. Native compiler :

They are written in same language as the target language.

Example: SMM (Source Implementation's ^S Target) is a compiler for the language of ; i.e; in the language that runs on m/c M. and generates o/p code that runs on m/c M.

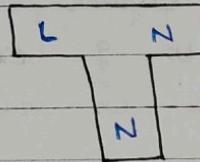
II. Cross compiler : (Retargeting compiler)

They are written in different languages as the target language.

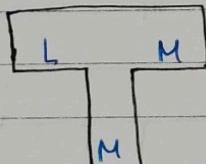
Example: SNM is a compiler for the language S that is in a language that runs on m/c N & generate o/p code that runs on m/c M.

Q. Write a compiler of language L for machine N with the help of an existing compiler of language L for machine M.

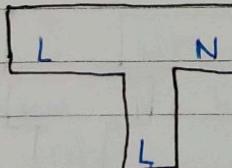
Required compiler :



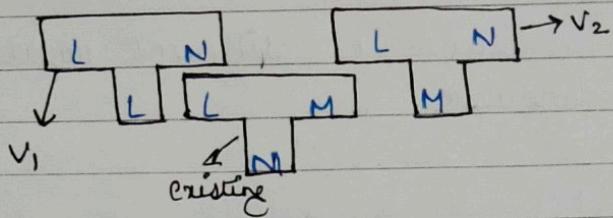
Existing compiler :



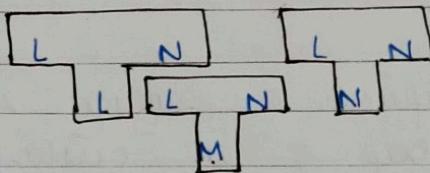
Step-1 : Create a compiler V_1



Step-2 : compile the compiler V_1 with the existing compiler.

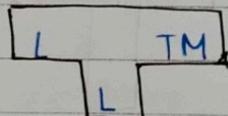


Step-3 : compile V_1 on V_2

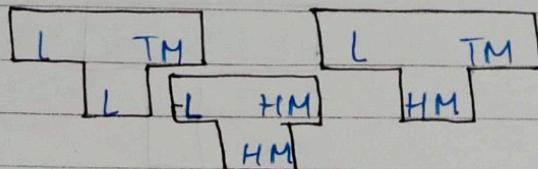


Q. Write a compiler L for machine TM with the help of existing compiler for the same language on machine HM .

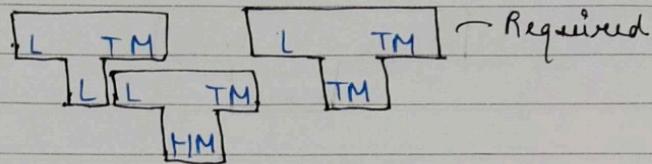
Step-1 : V_1



Step-2 : V_1 with existing.

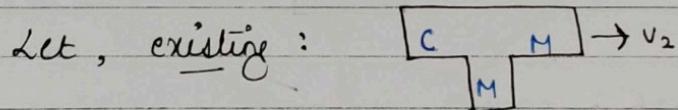
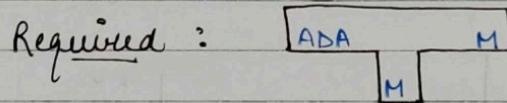


Step-3 : $v_1 \& v_2$

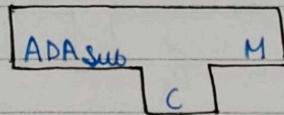


2nd Feb '17
Lecture 15

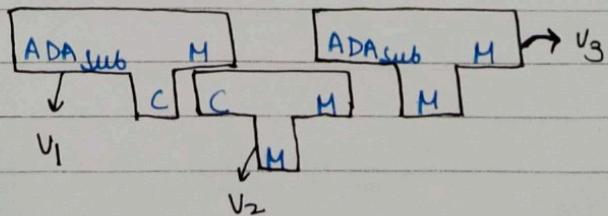
- Q. Implement an ADA compiler for machine M and no existing compiler is given.



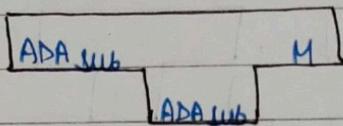
Step-1 : Create a compiler v_1 .



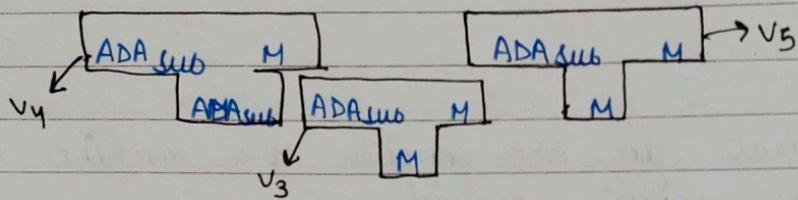
Step-2 : compile compiler v_1 on v_2



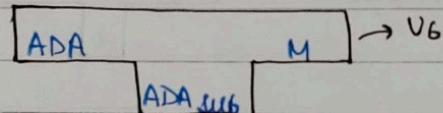
Step-3 : Create a compiler v_4 .



Step-4 : compile v_4 on v_3 .



Step-5 : Create a compiler v_6 .



Step-6 : compile v_6 on v_5

