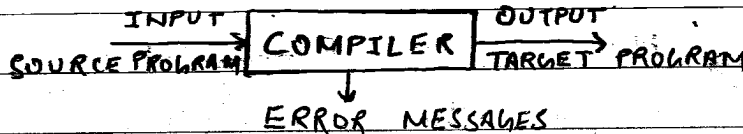


UNIT - 1 INTRODUCTION TO COMPILING & LEXICAL ANALYSIS

① Introduction to compiler -

Compiler is a program which takes one language (source program) as input and translates it into an equivalent another language (target program)



② Major Data Structures in a compiler -

Tokens - It represent basic program entities such as identifiers, literals, reserved words, operators, delimiters etc.

Syntax Tree - It is generated by the parser. It is usually constructed as a standard pointer-based structure that is dynamically allocated as parsing proceeds.

Symbol Table - It keeps information associated with all kind of identifiers. eg - constants, variables, functions, parameters, types, fields etc.

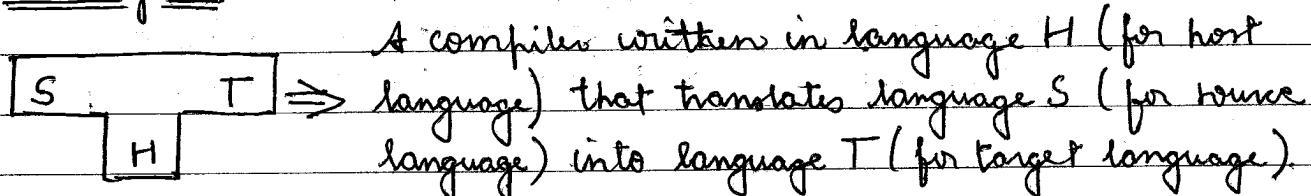
Literal Table - It stores constant and strings used in a program. Quick insertion and lookup are essential. Deletion is not necessary.

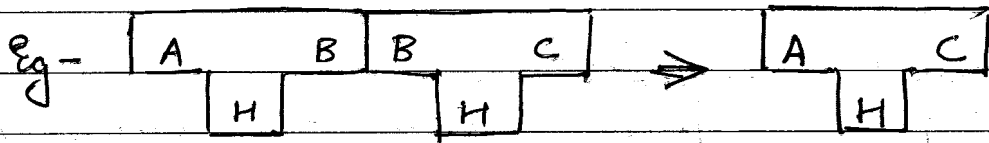
Temporary files - Used historically by old compilers due to memory constraints. Hold the data of various stages.

③ Bootstrapping and Porting -

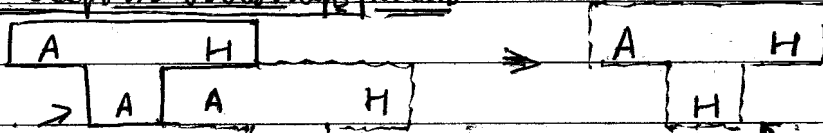
Bootstrapping is a process in which simple language is used to translate more complicated program which in turn may handle far more complicated program. This complicated can further handle even more complicated program and so on.

T-diagram -





1st Step in bootstrap process -

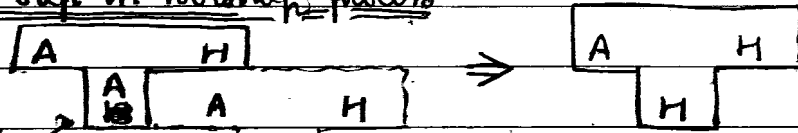


Compiler written in its own language A

'Quick & dirty' compiler written in machine language

Running but inefficient compiler

2nd Step in bootstrap process -

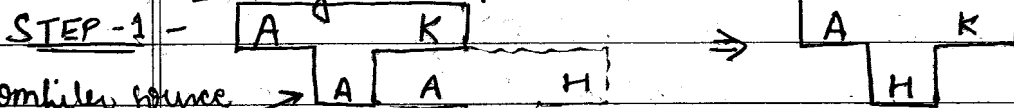


Compiler written in its own language A

Running but inefficient compiler (from the 1st step)

Final version of the compiler

→ Porting a compiler -

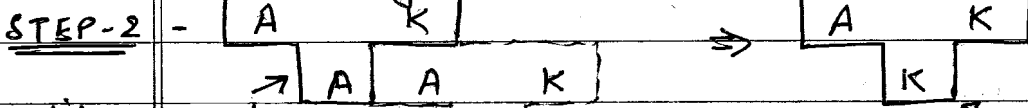


Compiler source

Code retargeted to K

Original compiler

Cross compiler



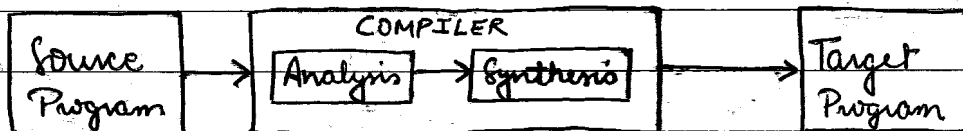
Compiler source code retargeted to K

Cross compiler

Retargeted compiler

COMPILER STRUCTURE -

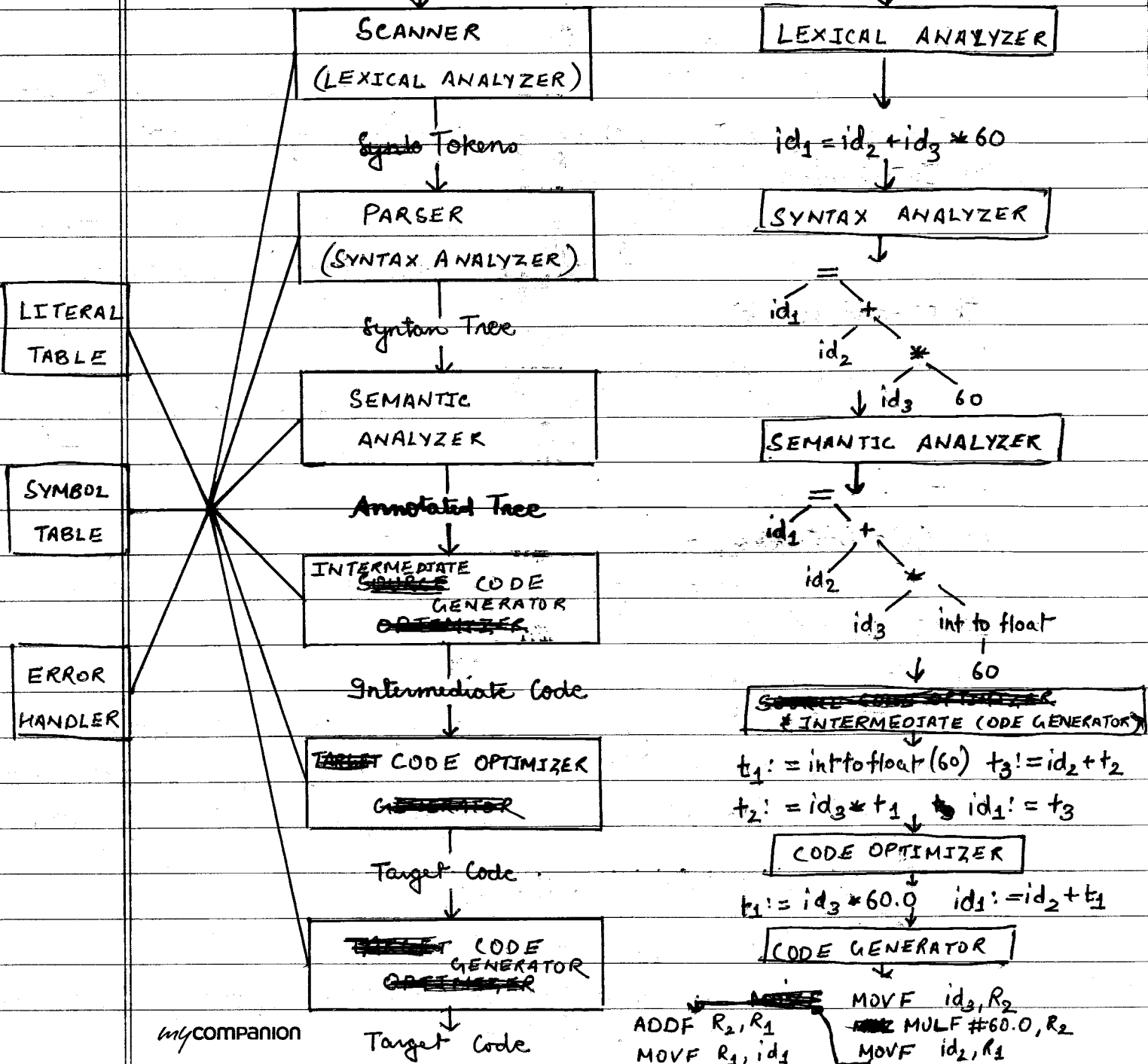
① Analysis-synthesis model of compilation



Analysis of the source program contains three steps that is lexical analysis, syntax analysis and semantic analysis. Intermediate code is generated from the input source program. In Synthesis part, three steps are involved that is intermediate code generation, code generation and code optimization.

② Various phases of a compiler -
Source Code

Eg - $a = b + c * 60$



(1) Lexical Analysis (Scanning) -

It is a phase of compilation in which the complete source code is scanned and source program is broken up into group of things called tokens. A token is a sequence of characters having a collective meaning.

(2) Syntax Analysis (Parsing) -

In this phase, tokens generated by the lexical analyzer are grouped together to form a hierarchical structure called parse tree or syntax tree.

(3) Semantic Analysis -

It determines the meaning of the source string. Eg like matching of parentheses, checking the scope of operation etc.

(4) Intermediate Code Generation -

It is a kind of code which is easy to generate and this code can be easily converted to target code. It can be represented as three-address code.

(5) Code Optimization -

It improves the intermediate code which reduces memory consumption and have a faster executing code, which improve the runtime of the target program.

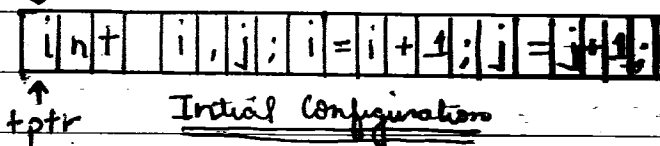
(6) Code Generation -

Target code is generated as assembly code or machine code

LEXICAL ANALYSIS -

① Input Buffering -

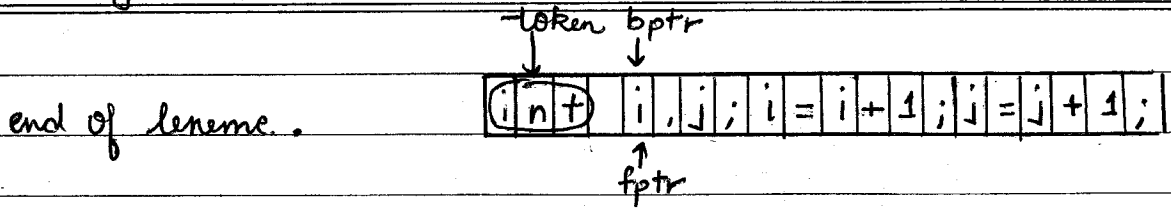
The lexical analyzer scans the input string from left to right one character at a time. It uses two pointers begin pointer (bptr) and forward pointer (fptr) to keep track of the position of the input stream.



As soon as the blank space is encountered by fptr, it indicates (or comma, or operator)

lexemes - Sequence of characters in the source program that are matched with the pattern of the token.

Eg - int, i, num, ans, choice;



→ Two methods are used -

→ One buffer scheme →

| | |
|-----|-----------|
| int | i = i + 1 |
|-----|-----------|

→ Two buffer scheme →

| | |
|-----|-----------|
| int | i = i + 1 |
|-----|-----------|

 buffer 1

| | |
|--------------|-----|
| ; j = j + 1; | eof |
|--------------|-----|

 buffer 2

If lexeme is short we ~~use~~ can use one buffer scheme but if lexeme is long we have to use two buffer scheme in which the 1st buffer and 2nd buffer are scanned alternatively and when the end of current buffer is reached the other buffer is filled.

If the length of lexeme is longer than the length of the buffer then the input can not be scanned completely.

eof is introduced at the end of both the buffers is called sentinel which is used to identify the end of buffer.

→ Code for input buffering -

```
if ( fptr == eof (buff1) ) /* encounter end of first buffer */
{
```

```
    /* Refill buffer2 */
```

```
    fptr ++;
```

```
}
```

```
else if ( fptr == eof (buff2) ) /* encounter end of second buffer */
```

```
{
```

```
    /* Refill buffer 1 */
```

```
    fptr ++;
```

```
}
```

```
else if ( fptr == eof (input) )
```

```
    return;
```

```
    /* terminate scanning */
```

```
else
```

```
    fptr ++ /* still remaining input has to be scanned */
```



② Specification of tokens -

To specify tokens regular expressions are used. When a pattern is matched by some regular expressions then token can be recognized.

→ Things and language -

String is a collection of finite number of alphabets or letters.

The things are synonymously called as words.

Collection of string is called the language.

$|S|$ → length of a string ϵ → empty string ϕ → empty set of strings

→ Operation on language -

Union ($L_1 \cup L_2$), Concatenation ($L_1.L_2$), Kleen Closure (L^*),

Positive Closure (L^+)

→ A language denoted by regular expressions is said to be a regular set or a regular language.

Eg - Regular Expression for identifier = letter (letter + digit)*

*** Unspecific information cannot be represented by regular expressions

③ Recognition of Tokens -

Token Representation →

| | |
|------------|-------------|
| Token Type | Token value |
|------------|-------------|

Token Type → category of token Token value → Information regarding token.

The lexical analyzer reads the input program and generates a

symbol table for tokens.

| eg - | TOKEN | CODE | VALUE | TOKEN | CODE | VALUE |
|------|------------|------|---------------------|-------------|------|-------|
| | if | 1 | - | != | 7 | 5 |
| | else | 2 | - | (| 8 | 1 |
| | while | 3 | - |) | 8 | 2 |
| | for | 4 | - | + | 9 | 1 |
| | identifier | 5 | Ptr to symbol table | - | 9 | 2 |
| | constant | 6 | Ptr to symbol table | = | 10 | - |
| | < | 7 | 1 | if (a < 10) | | |
| | <= | 7 | 2 | i = i + 2; | | |
| | > | 7 | 3 | else | | |
| | >= | 7 | 4 | i = i - 2; | | |

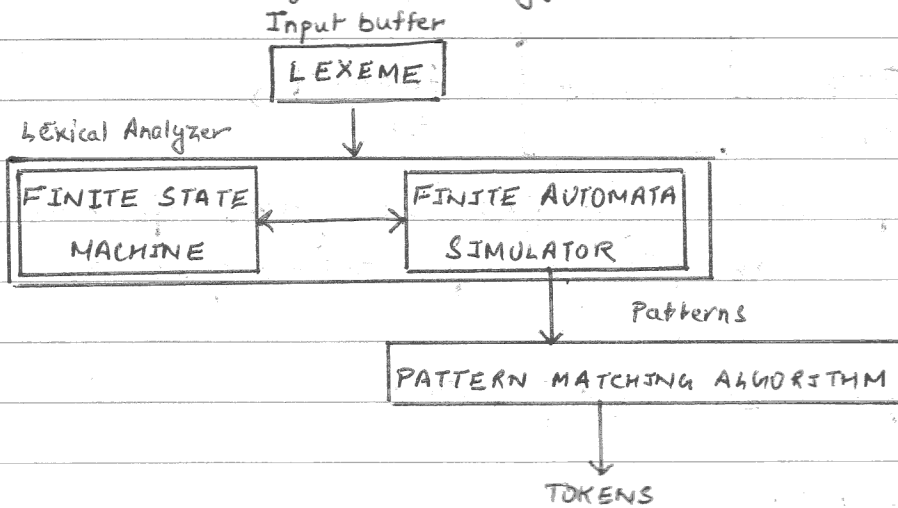
Our lexical analyzer will generate following token stream.

1, (8,1), (5,100), (7,1), (6,105), (8,2), (5,107), 10, (5,107), (9,1), (6,110), 2, (5,107), 10, (5,107), (9,2), (6,110).

Corresponding symbol table for identifiers and constants will be -

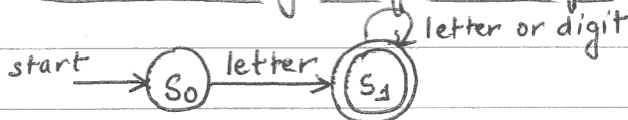
| LOCATION COUNTER | TYPE | VALUE |
|------------------|------------|-------|
| 100 | IDENTIFIER | a |
| 105 | CONSTANT | 10 |
| 107 | IDENTIFIER | L |
| 110 | CONSTANT | 2 |

④ Block schematic of lexical analyzer -



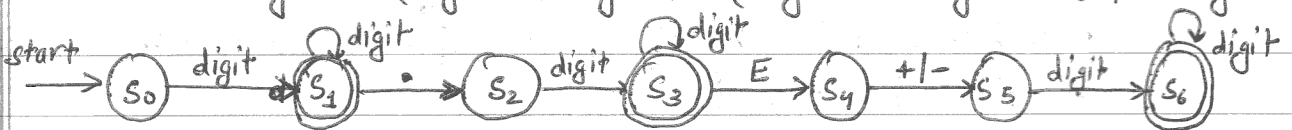
Transition

⑤ Transition diagram for Identifier - R.E. = letter (letter + digit)*



Transition diagram for constant -

R.E = $digit^+ + (digit^+ \cdot digit^+) + (digit^+ \cdot digit^+ E(+/-) digit^+)$

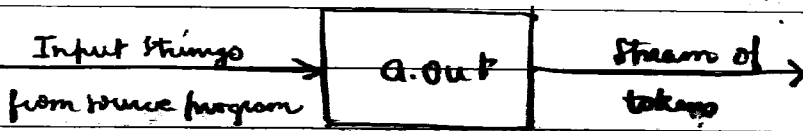
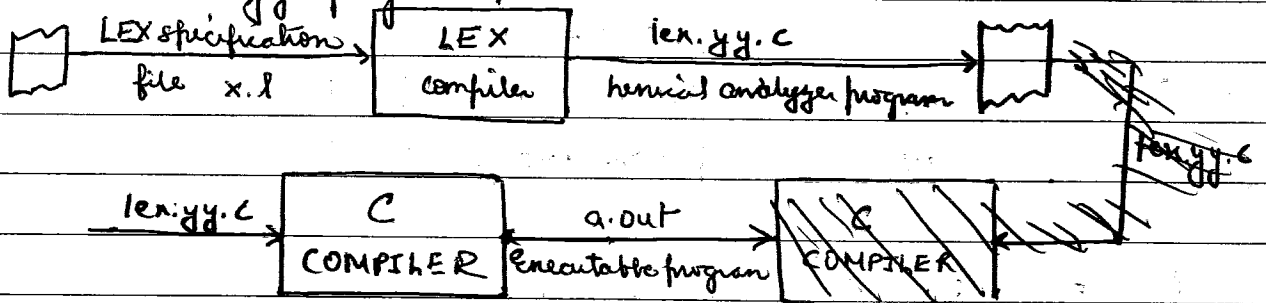


⑥ LEX (Lexical Analyzer Generator) -

It is a UNIX utility. In LEX tool, designing the regular expressions for corresponding tokens is a logical task.

LEX is used to write specification file with .l (dot l) extension.

→ LEX converted .l to .c where .c is a C program which is actually a lexical analyzer program.



→ (x.l) Specification file stores the regular expressions for the tokens

→ lex.yy.c consists of tabular representations of the transition diagrams

→ LEX program consists of three parts -

- { % ?
- DECLARATION SECTION
- { % }
- { % %
- RULE SECTION
- { % %
- AUXILIARY PROCEDURE SECTION

- Declaration section → declaration of variable constants

- Rule section → consists of regular expression with associated action.

eg + R₁ { action₁ }

- Auxiliary procedure section → Required procedures are defined

Eg:-

%i

%f

%%

"Rama" |

"Seeta" |

"Geeta" |

"Neeta" } printf ("\n Noun");

"Sings" |

"dances" |

"eats" } printf ("\n Verb");

%%

main ()

{

yywrap(); → Routine, defined in yywrap.c program.

}

Scanner starts scanning the source program

int yywrap() → called when scanner encounters end of file.

{

return 1;

}

Following commands are used to run the lex program n.s. in UNIX.

\$ lex x.l

\$ cc lex.yy.c

\$./a.out

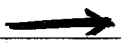
After entering this commands a blank space for entering input gets available. Thus we can give some valid input.

Rama eats

Noun

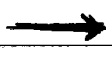
verb

Then press either control + c or control + d to come out of the output



LEX actions -

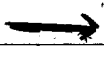
- (1) BEGIN - It indicates the start state.
- (2) ECHO - It emits the input as it is.
- (3) yytext - When lexer matches or recognizes the tokens from input tokens then the lexeme is stored in a null terminated string called yytext.
- (4) yyin - Standard input file that stores input source program.
- (5) yyleng - stores the length of the input string.
- (6) yylen()
- (7) yywrap()



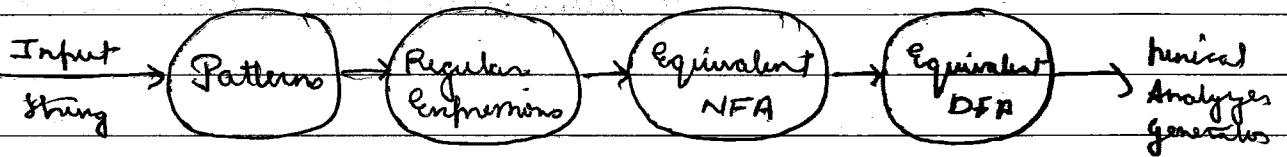
Design of lexical analyzer generator -

Two approaches -

- Pattern matching using NFA
- Using DFA for lexical analyzer.



Pattern matching using NFA -



⑦

Advantages

Issues of lexical analyzer -

- (1) lexical analysis and syntax analysis are separated out which reduces the burden on parsing phase.
- (2) Compiler efficiency get increased due to separation.

⑧

Issues of lexical analyzer -

- (1) Ambiguity of words (same word can have different meaning)
- (2) lookahead (identifying limits for sentence boundaries. Eg - Dr., Mr.)