

COMPILER DESIGN

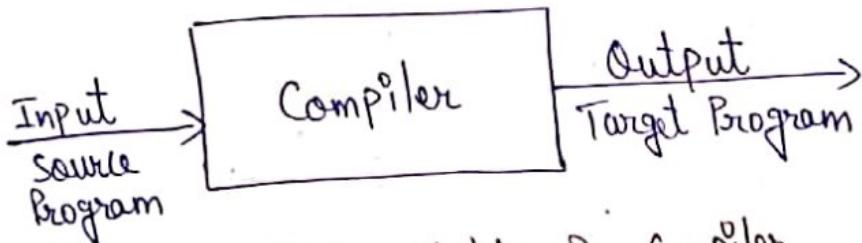
(1) to (14)

1

UNIT I : Overview of Compiler & Lexical Analysis

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

→ Compilers basically act as translators. The basic model of compiler can be represented as follows :



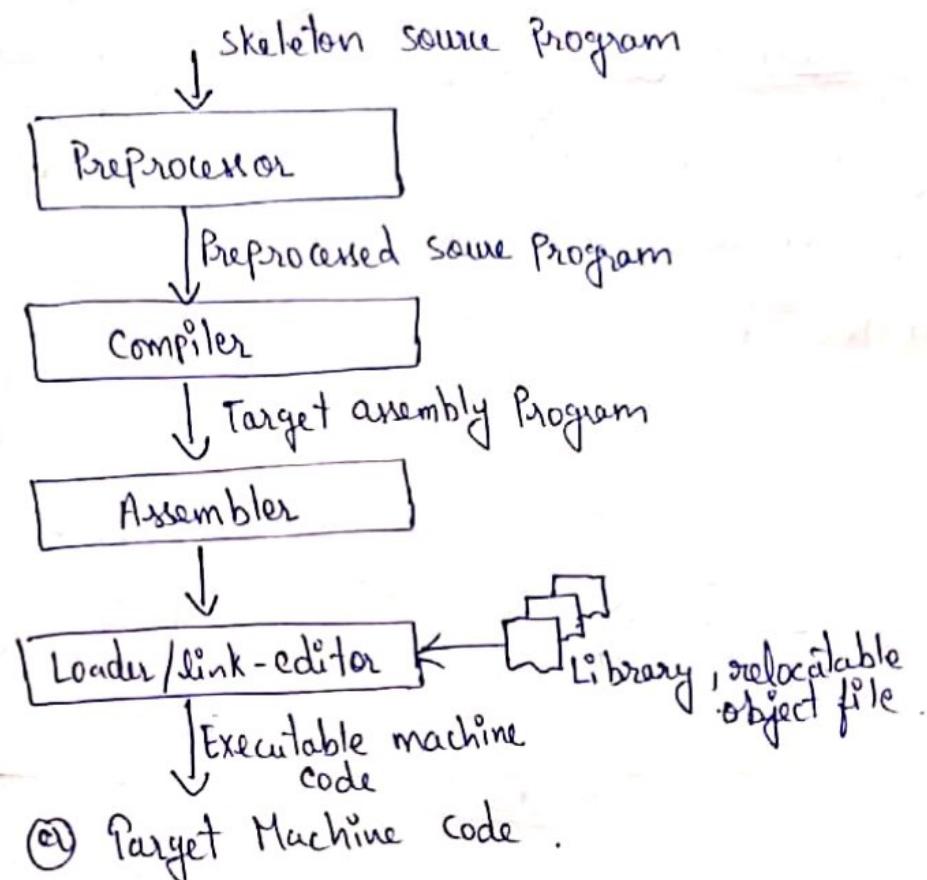
Basic Model of Compiler

→ The compiler takes a source program as higher level languages such as C, Pascal, Fortran & converts it into low level language or a machine level language such as assembly language.

* Properties of a Compiler:

- The compiler itself must be bug-free.
- It must generate correct machine-code
- The generated machine code must run fast.
- The compiler must be portable.
- It must give good diagnostics & error messages.
- It must have consistent optimization.

→ To create an executable form of your source program only a compiler program is not sufficient. we may require several other programs to create an executable target program

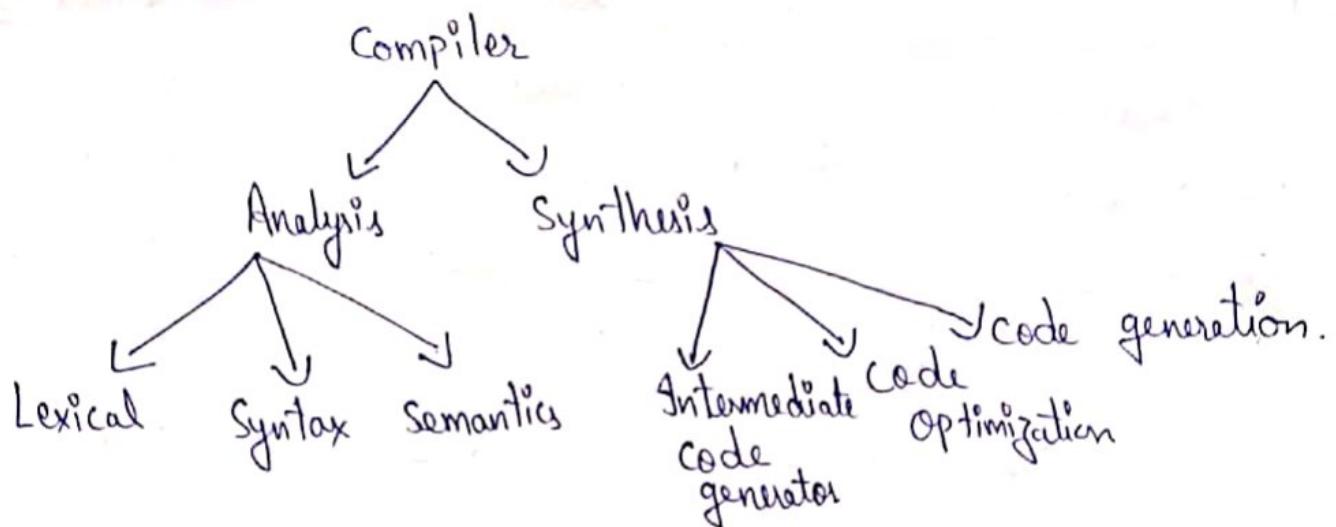


Process of Execution of Program.

* Structure of a Compiler:

The Process of Compilation is carried out in two parts namely:
"Analysis & Synthesis".

- Analysis Part is the front end of the compiler.
- Synthesis part is the back end of the compiler.
- Compilation process operates as a sequence of phases, each of which transforms one representation of the source program to another.



Phases of Compiler

(a) Lexical Analysis:

The lexical analysis is also called "Scanning". The complete source code is scanned & the source program is broken up into group of strings called "lexemes". For each lexeme, the lexical analyzer produces as output a "token" of the form : (token-name, attribute-value). Sequence of characters having a collective meaning

→ The blank characters which are used in the programming statement are eliminated during this phase.

Ex: c := a + 10 ;

c, a ⇒ identifiers
:= ⇒ assignment.
+ ⇒ operator
10 ⇒ constant.
Tokens

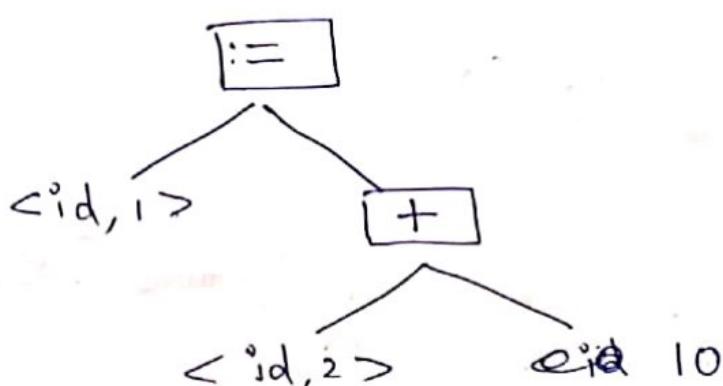
there are mapped to the tokens like this.

<id, 1> <:=> <id, 2> <+> <id, 3> 10 ;

→ we use regular expressions to recognize tokens.

(b) Syntax Analysis: The syntax analysis is also called "Parsing". In this phase the tokens generated by the lexical analyzer are grouped together to form a hierarchical tree-like structure, called as Parse tree or Syntax tree. we use context-free grammar (CFG) to recognize syntax or grammar.

Syntax tree for $\langle \text{id}, 1 \rangle \langle := \rangle \langle \text{id}, 2 \rangle \langle + \rangle \langle \text{id}, 3 \rangle \Rightarrow \text{id} \ 10$



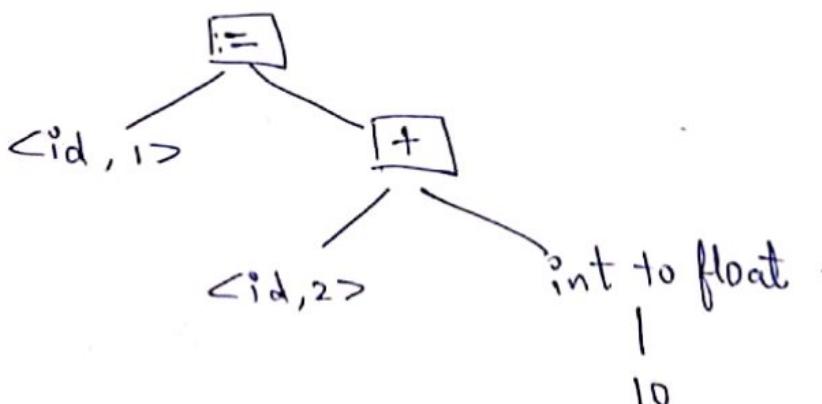
2 Parsing techniques in this phase

- Top-down
- Bottom-up

Parse tree or Syntax tree.

(c) Semantic Analysis:

Semantic Analysis determines the meaning of the source string. It also does "type checking" where the compiler checks that each operator has matching operands. For example, matching of Parenthesis in the expression, matching of the stmts, checking the scope of operation etc.



(d) Intermediate code generation:

The intermediate code is a kind of code which is easy to generate & this code can be easily converted to target code. Intermediate code is of many forms such as three address code, quadruple, triple, Posix etc. For the above example let us convert it into three address form (consists of instructions each of which has at most three operands).

$$\text{Ex: } t_1 = \text{int to float (10)}$$

$$t_2 = \text{id}_2 + t_1$$

$$t_3 = \text{id}_3 - t_2 \quad \text{id}_1 = t_2$$

$$id_1 \leq id_3$$

(e) Code optimization:

The code optimization phase attempts to improve the intermediate code. This is necessary to have a faster executing code or less consumption of memory. By optimizing the code the overall running time of the target program can be improved.

$$\text{Ex: } t_1 = \text{id}_2 + 10.0 \quad // \text{code is optimized here. unnecessary variables are removed.}$$

$$\text{id}_1 = t_1$$

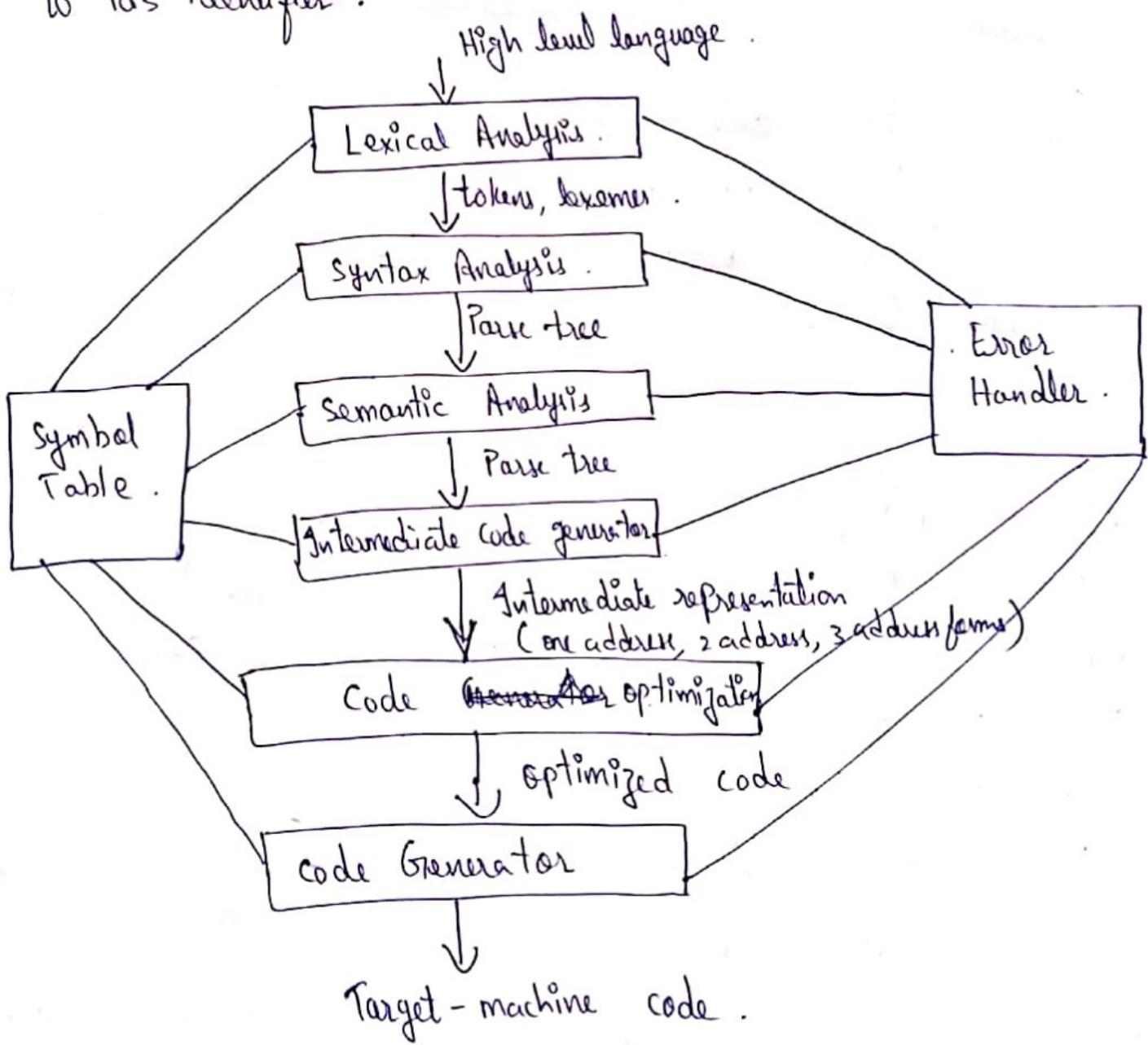
(f) Code Generation: In code generation phase the target code gets generated. The intermediate code instructions are translated into sequence of machine instructions.

Eg : Mov / id₂ / R₁
ADD R₁, #10.0 , R₁

Mov id₂, R₁
ADD #10.0 , R₁

Mov R₁, id₃

we are moving value of id₂ to 'R₁' register. Then we add id₂ & 11
& result is stored in 'R₁'. Then value of R₁ is assigned
to id₃ identifier.



Phases of Compiler

* Symbol Table

Symbol table is a data structure containing a record for each variable name, with fields for the attributes of the name. Symbol table also stores information about the subroutines used in the program.

→ It allows us to find the record for each identifier quickly & to store or retrieve data from that record efficiently.

* Error detection & Handling:

In compilation, each phase detects errors. These errors must be reported to error handler, whose task is to handle the errors so that the compilation can proceed. Errors are reported in the form of 'messages'.

→ Large no. of errors can be detected in syntax analysis phase. Such errors are called as: "syntax errors".

* Pass: when several phases are grouped together, we call it as Pass.

* The Science of Building a Computer:

A compiler must accept all source programs that conform to the specification of the language. Every transformation performed by the compiler while translating a source program must preserve the meaning of the program being compiled.

Compiler writers thus have influence over not just the

compilers they create, but all the programs that their compilers compile.

a) Modeling in Compiler Design & Implementation

The study of compilers is mainly a study of how we design the right mathematical models & choose the right algorithms, while balancing the need for generality & power against simplicity & efficiency.

→ Some of most fundamental models are finite-state machines & regular expressions (for describing the lexical ^{units of} programs), context-free grammars (to describe syntactic structure).

b) The science of code optimization:

The term 'optimization' in compiler design refers to the attempts that a compiler makes to produce code that is more efficient than the obvious code.

→ Compiler optimizations must have the following design objectives:

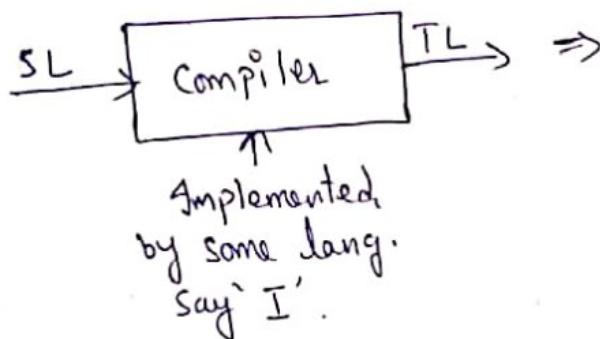
- The optimization must be correct, i.e. preserve the meaning of the compiled program.
- The optimization must improve the performance of many programs.
- The compilation time must be kept reasonable.
- The engineering effort required must be manageable.

* Bootstrapping :

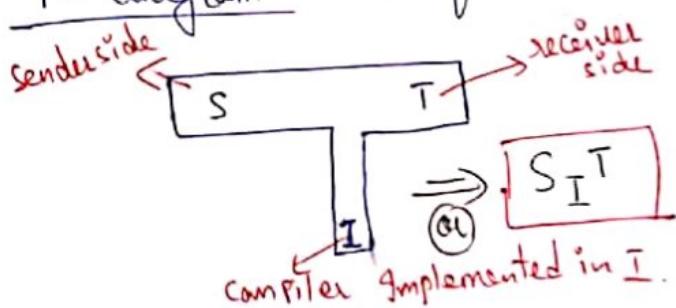
A process by which simple language is used to translate more complicated program, which in turn may handle an even more complicated program.

→ There are three types of languages.

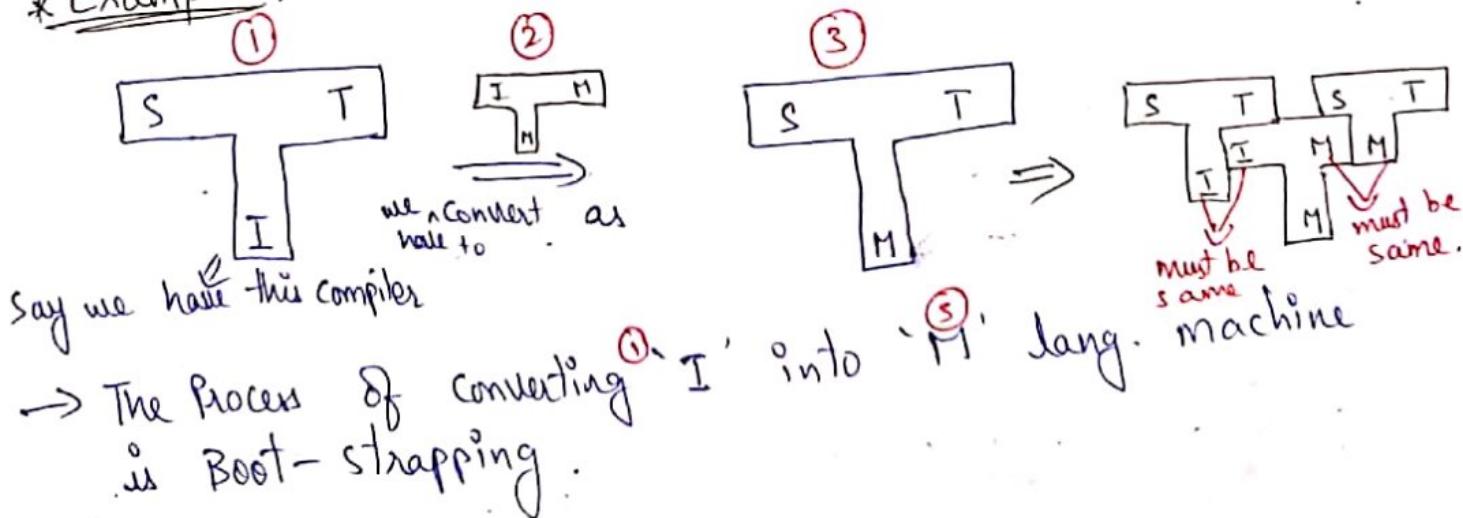
- Source language (SL); i.e. the application program
- Target language (TL) in which machine code is written.
- Implementation language in which a compiler is written.



This can be represented by "T-diagram" as follows:



* Example :

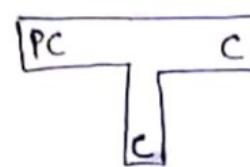


Example :-

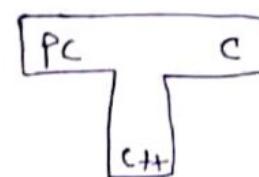
We have Pascal Translator written in 'c' lang.

i/p is Pascal code(Pc), o/p - 'c'. Create Pascal translator in C++.

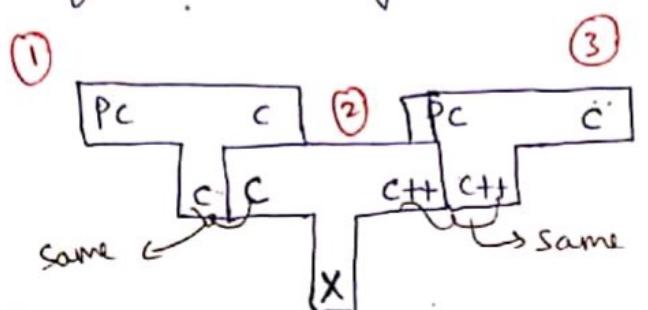
Sol:- Given :-



Convert it into



With help of Bootstrapping :-



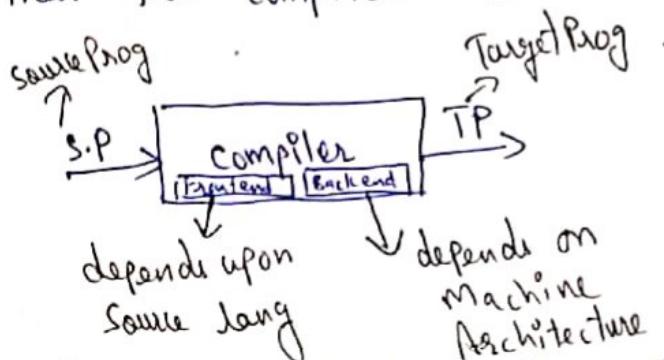
$x \Rightarrow$ can be anything
like java or
any other lang

(or) $P_c C + C_x C++ \Rightarrow P_{C++} C$

* Cross Compiler :-

Cross compiler is a compiler which is capable of creating executable code for a platform other than ^{the one} on which compiler is running.

Eg:- A compiler which runs on windows 7 but this " generates a code which can run on Android smartphone. Then this compiler is known as Cross-compiler.



\Rightarrow we need to change only the backend to implement cross compiler.

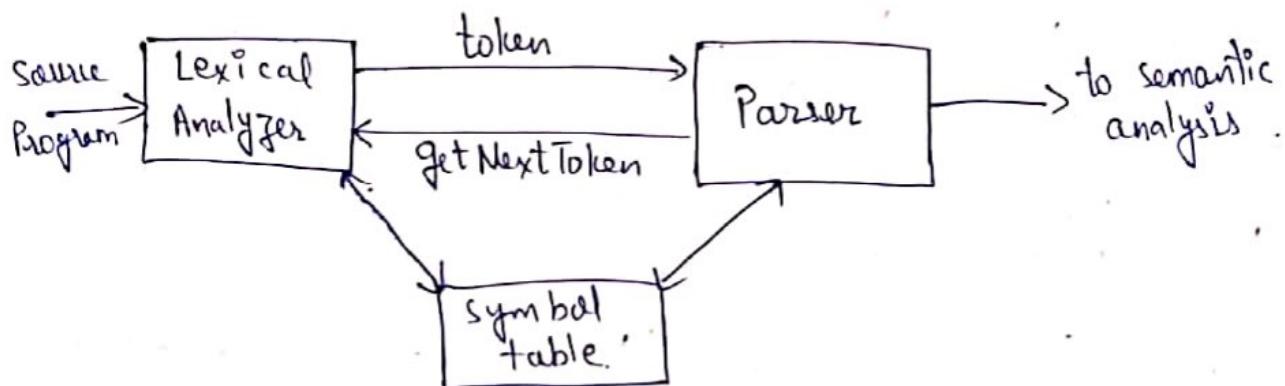
\rightarrow using cross compilation, platform independency can be achieved.

LEXICAL ANALYSIS

* Role of Lexical Analyzer:

Lexical Analyzer reads the input source program from left to right one character at a time & generates the sequence of tokens. It is the first phase of a compiler.

→ Each token is a single logical cohesive unit such as Identifier, keywords, operators & punctuation marks. Then the Parser to determine the syntax of the source program can use these tokens.



Interactions between the Lexical analyzer & the Parser

→ As the lexical analyzer scans the source program to recognize the tokens, it is also called an Scanner.

* Functions of Lexical Analyzer:

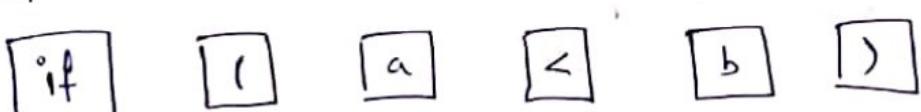
- It produces stream of tokens
- It eliminates whitespaces (blank, newline, tab) & comments.
- It generates symbol table which stores the information about identifiers, constants encountered in the i/p.
- It keeps track of line no's.
- It reports the error encountered while generating the tokens.

* Token : A token is a pair consisting of a token name & an optional attribute value. It describes the class or category of input string. For example identifiers, keywords, constants are called tokens.

* Pattern : A pattern is a description of the form that the lexemes of a token may take. It is a set of rules that describe the token.

* Lexemes : Sequence of characters in the source program that are matched with the pattern of the token.
For example int, i, num, choice;

Example : if (a < b)

 tokens which are generated

here "if", "(", "a", "<", "b", ")" \Rightarrow are all lexemes.

\rightarrow "if" is a keyword

\rightarrow "(" is opening parenthesis

\rightarrow "a" is identifier \rightarrow collection of letters

\rightarrow "<" is an operator

\rightarrow "b" is identifier

\rightarrow ")" is closing parenthesis.

* Attributes for tokens:

If Token: lexeme
 number — 0, 1, 2, ... 0, 6.2, 1, ...

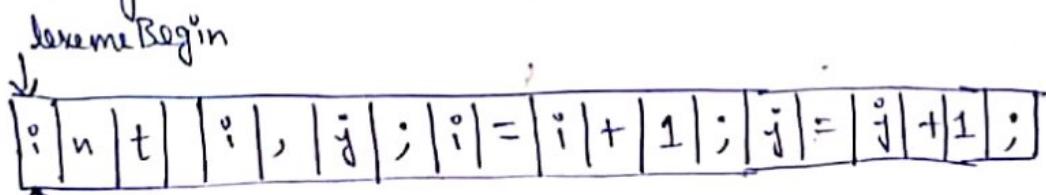
When more than one lexeme can match a pattern, the lexical analyzer must provide the subsequent compiler phases additional information about the particular lexeme that matched.

Eg: Pattern for "token: Number" matches both '0' and '1'.

→ Lexical analyzer returns to the parser not only a token name but an attribute value that describes the lexeme represented by the token. Token name influences parsing decisions, while the attribute value influences translation of tokens after the phx parse.

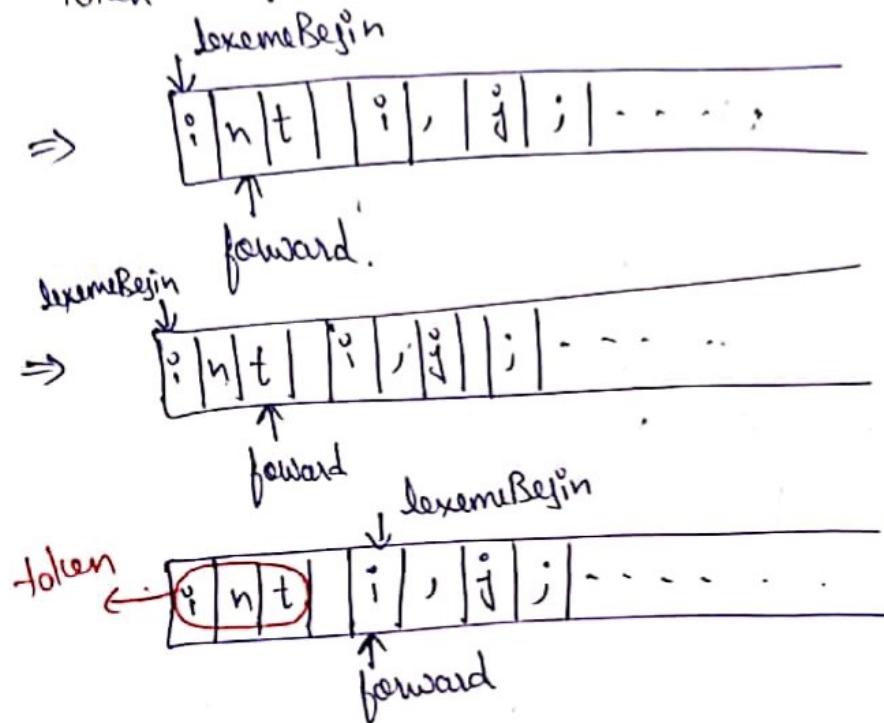
* Input Buffering:

The lexical analyzer scans the input string from left to right one character at a time. It uses two pointers: begin pointer "lexemeBegin" and forward pointer "forward" to begin pointer "lexemeBegin" and forward pointer "forward" to keep track of the portion of the input scanned. Initially both the pointers point to the first character of the input string as:



Initial configuration

The 'lexemeBegin' pointer remains at the beginning of the string to be read & the 'forward' pointer moves ahead to search for end of lexeme. As soon as the blank space is encountered, it indicates end of lexeme. & then both 'lexemeBegin' & 'forward' is set at next token 'i'.



→ Reading I/P's from secondary storage is costly. So, a block of data is first read into a buffer & then scanned by lexical analyzer.

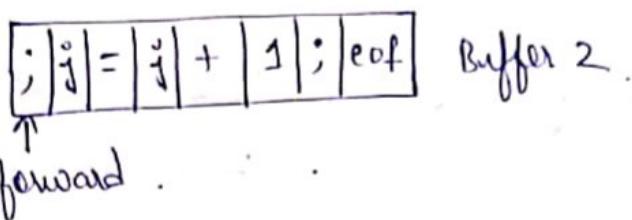
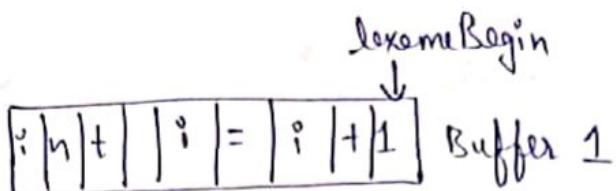
Input Buffering

→ Before one buffer scheme was used but now we are using two buffer scheme.

* Two Buffer Scheme

Two buffers each of size 'N' (usually size of disk block) are used. Each buffer can read 'N' characters. In this method, first buffer & second buffer are scanned alternately. When the end of current buffer is reached the other buffer is filled.

to store I/P
string.



Two Buffer scheme

* Sentinel :

Both buffers have a special character "eof" at the end. When 'lexemeBegin' pointer encounters 'eof' at Buffer 1, then it starts filling up second buffer. Similarly when 'eof' is encountered at Buffer 2, it starts filling Buffer 1 & so on.

→ This 'eof' character introduced at the end is called "Sentinel" which is used to identify the end of buffer.

* Specifications of tokens :

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

* Alphabet : Alphabet is a finite, non-empty set of symbols. It is represented by the symbol ' Σ '.

Iys: $\Sigma = \{0, 1\}$, $\Sigma = \{a, b\}$.

$\Sigma = \{a, b, c, \dots\}$, $\Sigma = \{\text{e}, \text{f}, \dots, \text{g}\}$

* String : String is a finite sequence of symbols chosen from some alphabet.

Eg: $\Sigma = \{0, 1\}$.

0011 is a string of Σ , 102 is not a string of Σ .
1010 " " " " " . , 234 " " " " " . . .

→ Length of string is denoted by $|s|$.

→ Empty string can be denoted by ' ϵ '.

* Language : Set of strings which belong to some alphabet ' Σ ' is called as a Language (L).

Eg: $\Sigma = \{0, 1\}$.

$L = \{00, 01, 10, 11\}$.

$L = \{0^n 1^n | n \geq 1\}$.

* Operations on Language:

There are various operations which can be performed on a language as follows: below. Let ' L ', ' M ' be two languages.

<u>Operation</u>	<u>definition & Notation</u>
• Union of L & M	$L \cup M = \{s s \text{ is in } L \text{ or } s \text{ is in } M\}$
• Concatenation of L & M	$LM = \{st s \text{ is in } L \text{ and } t \text{ is in } M\}$
• Kleen closure of ' L '	$L^* = \bigcup_{i=0}^{\infty} L^i$
• Positive closure of ' L '	$L^+ = \bigcup_{i=1}^{\infty} L^i$

Definitions of operations on languages

Eg: $L = \{A, B, c, \dots Z, a, b, c, \dots z\}$, $D = \{0, 1, 2, \dots 9\}$
are two languages.

- LUD is a set of letters & digits
- LD is a set of strings consisting of letters followed by digits.
- L^5 is a set of strings having length of 5 each.
- L^* is " " " all strings including ' ϵ '.
- L^+ " " " except ' ϵ '.

* Regular Expressions: A Regular Expression is a notation to represent certain sets of strings in an algebraic fashion. This notation involves a combination of strings of symbols from some alphabet ' Σ ', the symbol of null string ' ϵ ', star operator (*) & plus operator (+).

* Example:

(1) Write a Regular Expression (RE) for a language containing the strings of length two over $\Sigma = \{0, 1\}$.

$$\text{sol: } RE = (0+1)(0+1)$$

(2) Write a RE for lang. containing strings which end with 'abb' over $\Sigma = \{a, b\}$.

$$\text{sol: } (a+b)^*abb$$

(3) Write a RE for a recognizing identifier.

Sol: For denoting identifier, we will consider a set of letters

& digits because identifier is a combination of letters or letter & digits but having first character as letter always.
Hence RE can be denoted as:

$$RE = \text{letter} (\text{letter} + \text{digit})^*$$

where letter = (A, B, ... Z, a, b, ... z) & digit = (0, 1, 2, ... 9).

* Regular Definitions:

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

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

$$\vdots \quad \vdots$$

$$d_n \rightarrow r_n$$

where,

- Each d_i is a new symbol, not in Σ & not the same as any other of the d 's &
- Each r_i is a regular expression over the alphabet $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$.

→ Some of the notations used for writing the regular expressions are as follows:

(1) one or more instances : To represent one or more instances " $+$ " sign is used. $y \vdash r^+$.

(2) zero or more instances : To represent zero or more instances " $*$ " sign is used. $y \vdash r^*$.

(3) character classes : A class of symbols can be denoted by []. $s \vdash [012]$ means 0 or 1 or 2.

* Examples :-

(1) write Regular definition for 'c' language identifiers.

Sol: 'c' identifiers are strings of letters, digits & underscores.

Regular definition :-

$$\text{letter} \rightarrow [A-Z a-z]$$

$$\text{digit} \rightarrow [0-9]$$

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

(2) write Regular definition for Unsigned numbers.

Sol: unsigned nos (integer or floating point) are strings like:
5080, 0.01234, 6.33458, 0.89E-4 etc.

Regular definition :-

$$\text{digit} \rightarrow [0-9]$$

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

$$\text{number} \rightarrow \text{digits} (\cdot \text{digit})? (E [+-]? \text{digit})?$$

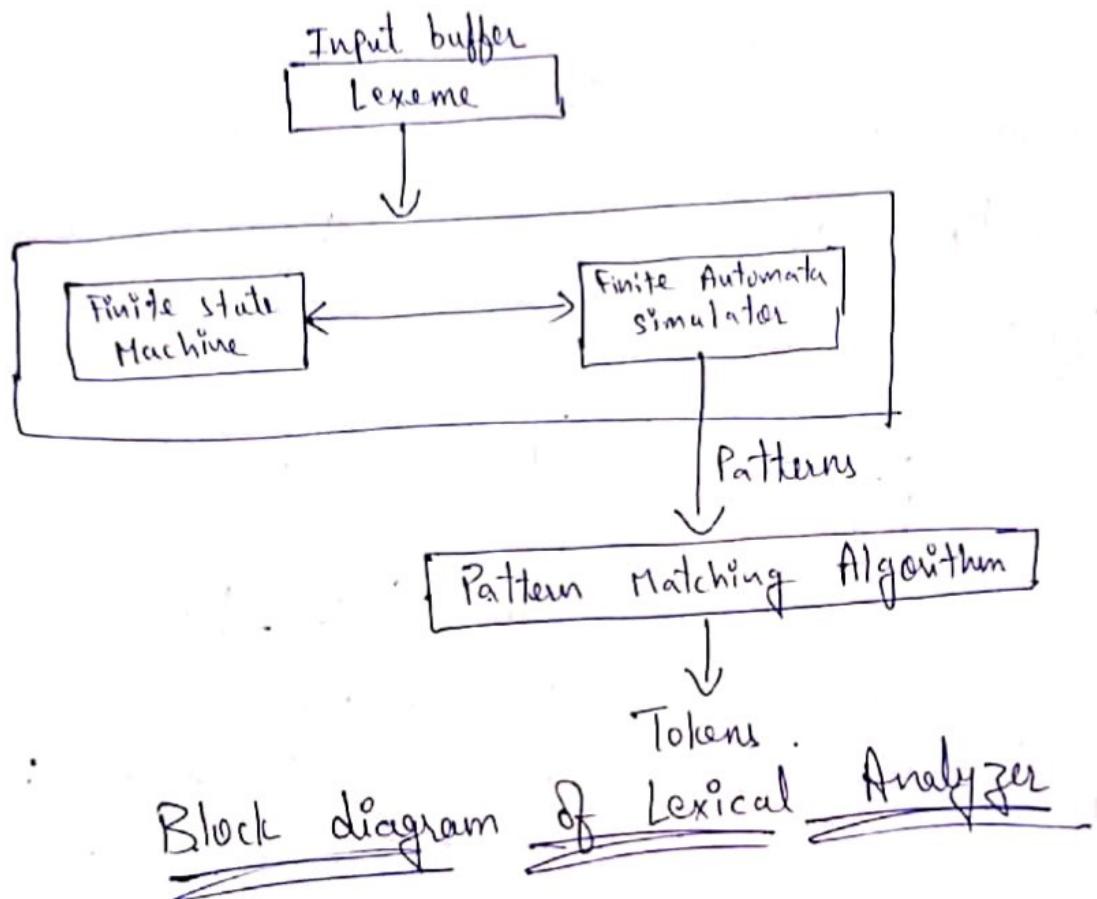
↑ zero or one instance

* Recognition of Tokens :-

The token is usually represented as:

Token type | Token value/attribute

→ The token type tells us the category of token & token value gives us the information regarding analysis, process, the symbol table is maintained. The token value can be a pointer to symbol table in case of identifiers & constants.



* Example :

digit	$\rightarrow [0-9]$
digits	$\rightarrow \text{digit}^+$
number	$\rightarrow \text{digits} (\cdot \text{digit})? (E [\text{+}-]?) \text{digit}^?$
letter	$\rightarrow [A-Za-z]$
:id	$\rightarrow \text{letter} (\text{letter} \text{digit})^*$
if	$\rightarrow \text{if}$
then	$\rightarrow \text{then}$
else	$\rightarrow \text{else}$
relational operator	$\rightarrow < > <= >= = < >$

Sample Regular definition

Lexemes	Token Name	Attribute Value
Any ws	-	-
if	if	-
then	then	-
else	else	-
Any id	id	Pointer to table entry
Any number	number	Pointer to table entry
<	relOp	LT
<=	relOp	LE
=	relOp	EQ
<>	relOp	NE
>	relOp	GT
>=	relOp	GE

Tokens, their patterns & attribute values

* Transition Diagrams:

In the middle of lexical analyzer process we convert patterns into transition diagrams. Transition diagrams have a collection of nodes or circles called states. Each state represents a condition that could occur during the process of scanning the input looking for a lexeme that matches one of several patterns.

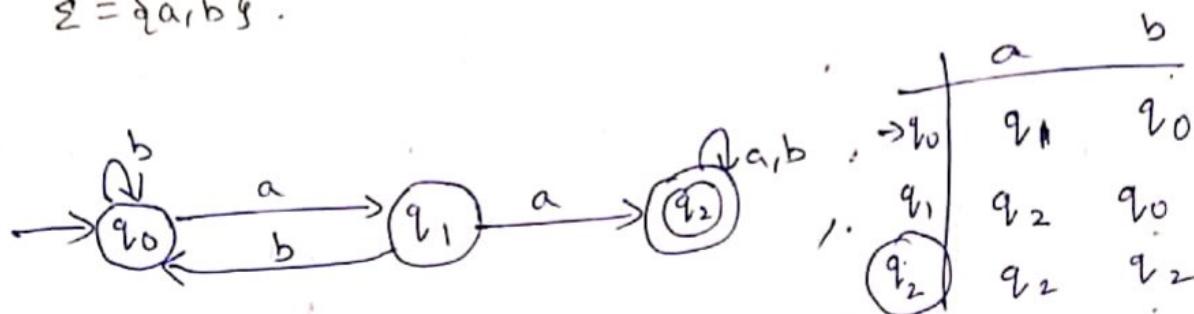
→ Edges are directed from one state to another. Each edge is labeled by a symbol or set of symbols.

- $\rightarrow q_0 \Rightarrow$ initial state
- $\textcircled{0} \Rightarrow$ final state

* Examples :-

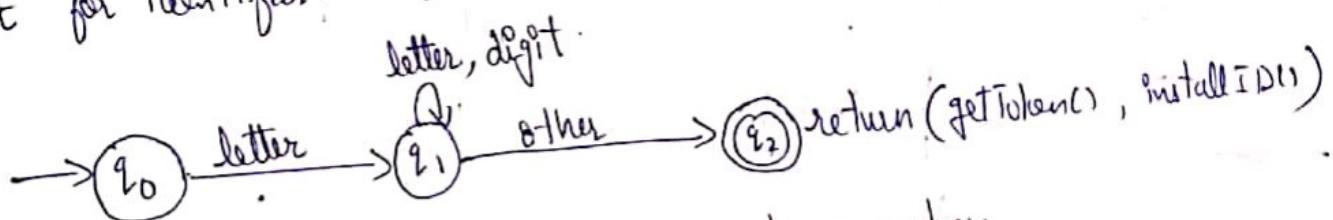
- ① Design a transition diagram for the language consisting of all the strings containing at least one pair of consecutive a's over $\Sigma = \{a, b\}$.

Sol:-

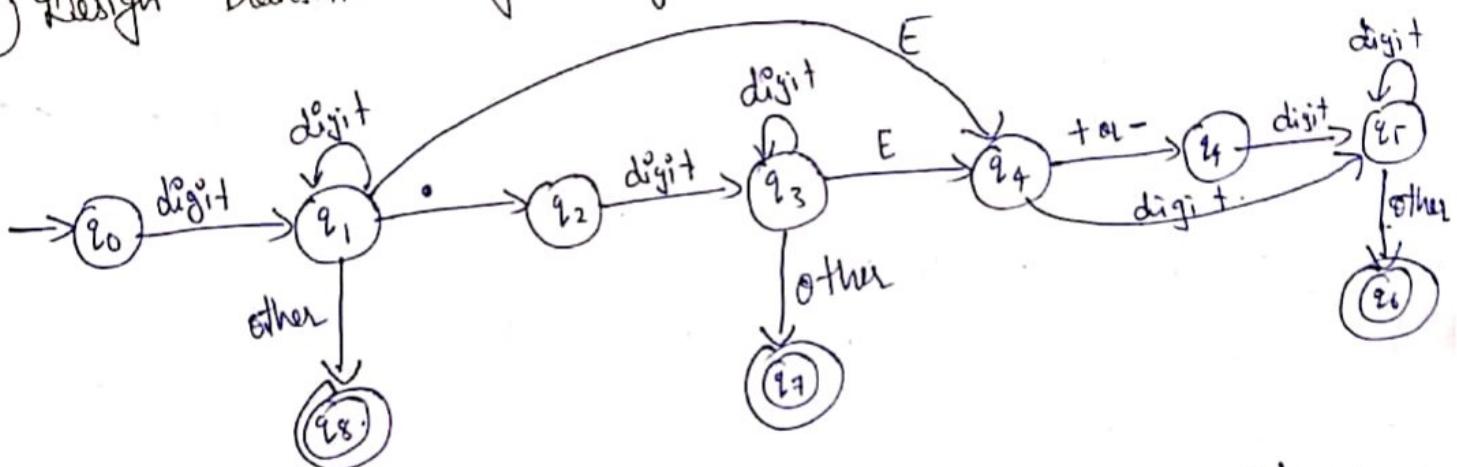


- ② Design transition diagram for identifiers.

Sol:- RE for identifier: e.g. letter (letter | digit)*

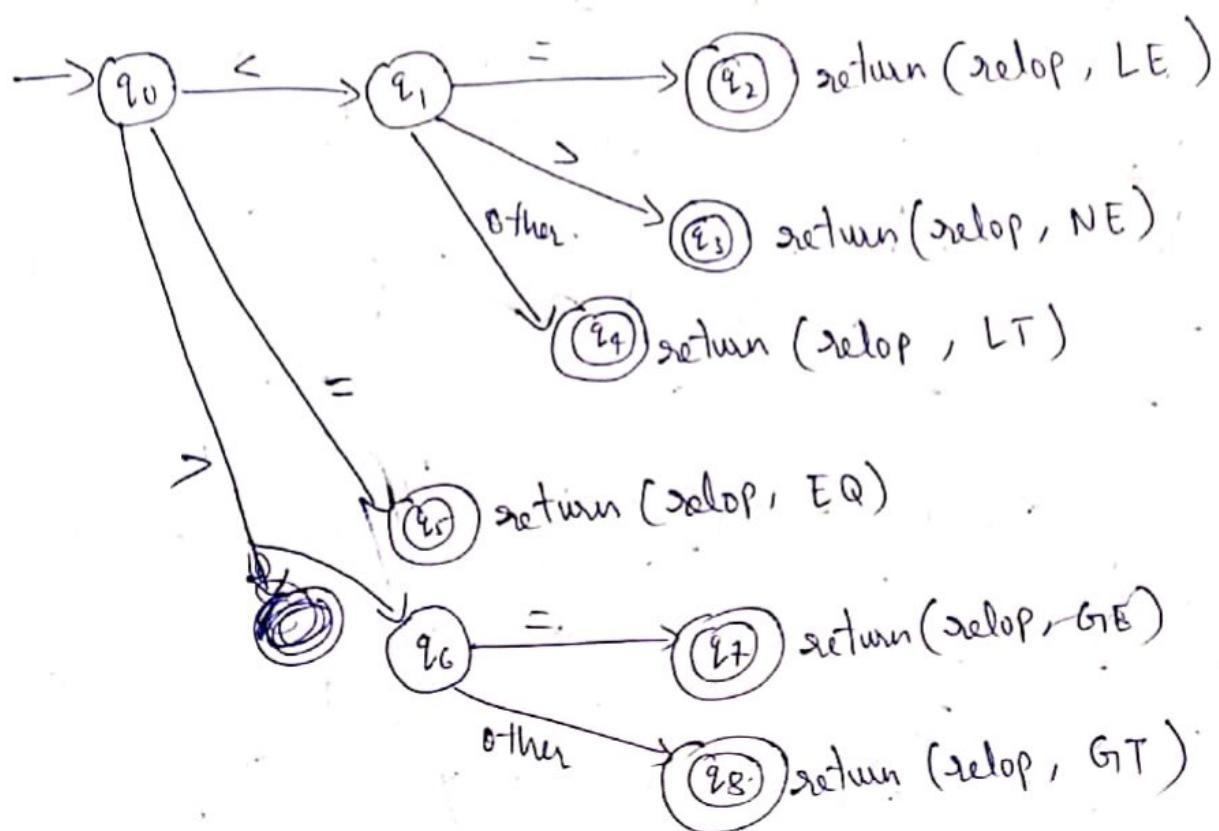


- ③ Design transition diagram for unsigned numbers.



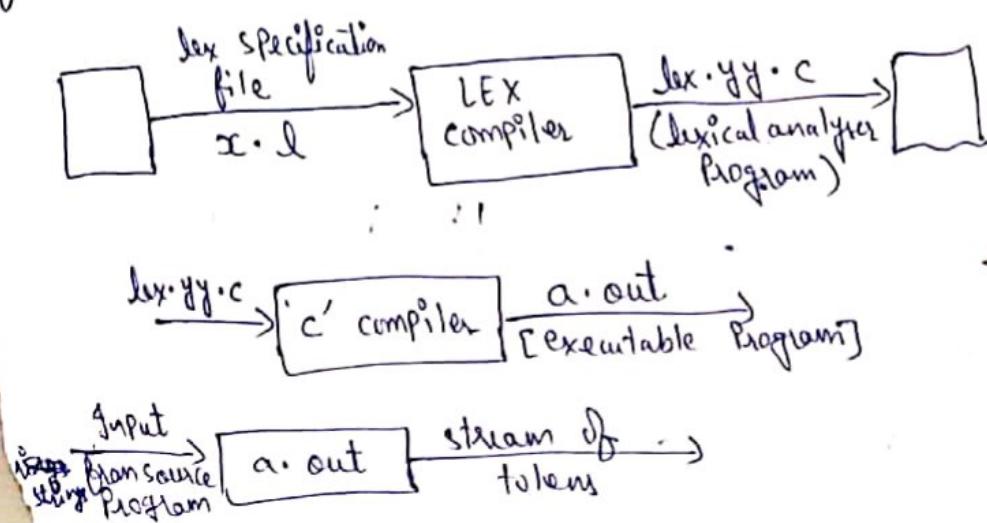
RE = digit⁺ | (digit⁻)⁺ (.) (digit⁺)⁺ | (digit⁺)⁺ (.) (digit⁺)⁺ E (+|-) (digit⁺).

④ Design transition diagram for relational operators (relOp).



* LEX ^{(@) FLEX} — Lexical Analyzer Generator:

Regular expressions are used in recognizing the tokens. A tool called LEX or FLEX (recent implementation) allows one to specify a lexical analyzer by specifying regular expressions to describe patterns for tokens.

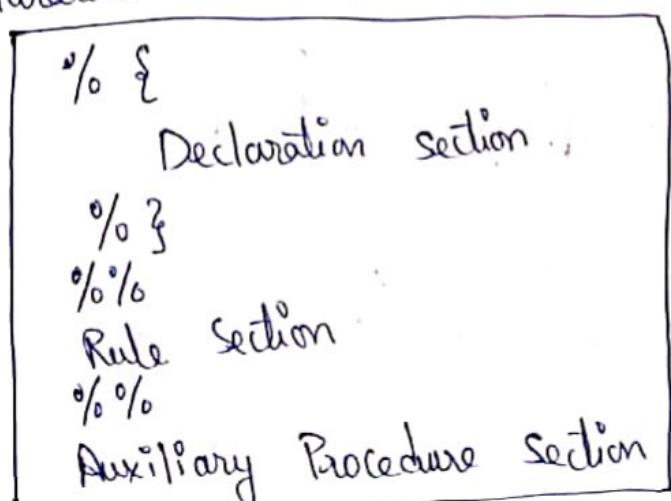


Generation of lexical analyzer using LEX

→ Specification file can be created using a file with extension ".l" (dot 'l'). say for example "xc.l". This xc.l is then given to LEX compiler to produce lex.yy.c. This lex.yy.c is a 'c' program which is actually a lexical analyzer program. We know that the specification file stores the regular expressions for the tokens, the lex.yy.c file consists of the tabular representation of the transition diagrams constructed for the regular expressions of specification file say xc.l.

→ The lexemes can be recognized with help of tabular transitions diagrams & some standard routines. In the specification file of LEX, actions are associated with each regular expression. These actions are simply pieces of 'c' code. These pieces of 'c' code are directly carried over to the lex.yy.c. Finally the 'c' compiler compiles this generated lex.yy.c & produces an object program a.out. When some input stream is given to a.out then sequence of tokens get generated.

→ LEX Program consists of three parts:
a) Declaration section, b) Rule section & c) Procedure section.



LEX Program format

(a) Declaration section: In this section, declaration of variables, constants can be done. Some regular definitions can also be written in this section. The regular definitions are basically components of regular expressions appearing in the rule section.

(b) Rule section: It consists of regular expressions with associated actions. These translation rules can be given as:

R_1	{action ₁ }
R_2	{action ₂ }
:	:
R_n	{action _n }

where each R_i is a regular expression & each action is a program fragment describing what action is to be taken for corresponding regular expression. These actions can be specified by piece of 'C' code.

(c) Auxiliary Procedure section: In this section, all the required procedures are defined. Sometimes these procedures are required by the actions in the rule section.

→ The lexical analyzer or scanner works in co-ordination with Parser. When activated by the Parser, the lexical analyzer begins reading its remaining input, character by character at a time. When the string is matched with one of the regular expressions R_i then corresponding action will

get executed on this action; returns the control to the Parser.

* Example LEX Programming:

// Program name: xc.l

% { } // Declaration section

% % // Rule section

" Rama " |

" seeta " |

" Geeta " |

" Neeta " | printf (" In Noun ");

" sings " |

" dances " |

" eats " printf (" In Verb ");

% % main () // Auxiliary Procedure section

{ yyflex (); }

int yywrap ()

{ return 1; }

}

→ consists of RE's & actions

→ This is a simple program that recognizes noun & verb from the string.

* Execution :

\$ lex x.l ↳

\$ cc lex.yy.c ↳

\$.\a.out ↳

Op : ip1 : Rama eats ↳

op : Noun
Verb.

ip2 : Seeta sings ↳

op : Noun
Verb

* LEX Actions :

① BEGIN : It indicates the start state. The lexical analyzer starts at state '0'.

② ECHO : It emits the input as it is.

③ yyflex() : As soon as call to yyflex() is encountered, scanner starts scanning the source program.

④ yywrap() : The function yywrap() is called when scanner encounters end of file. If yywrap() returns '0'(no) then scanner continues scanning. When " " '1'(one) that means end of file is encountered.

⑤ yyin : It is the standard file that stores ipx.

*FLEX :

Flex is a tool for generating scanners. It scanner, sometimes called a 'tokenizer' is a program which recognizes lexical patterns in text. The flex program reads user-specified input files or its standard input if no file names are given for a description of a scanner to generate. The description is in the form of pairs of regular expressions & 'c' code called rules. Flex generates a 'c' source file named "lex.yy.c" which defines the function `yylex()`. The file `lex.yy.c` can be compiled & linked to produce an executable.

→ when the executable is run, it analyzes its I/P for occurrences of text matching ^{with} the regular expression for each rule. whenever it finds a match, it executes the corresponding 'c' code.