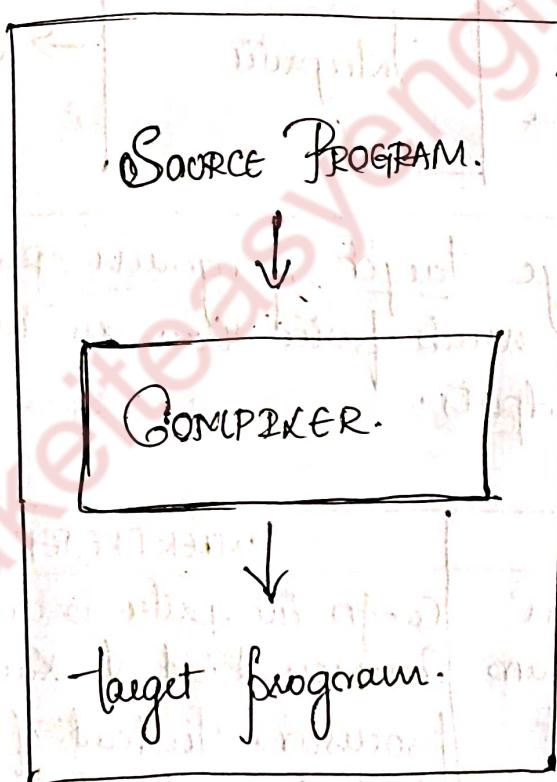
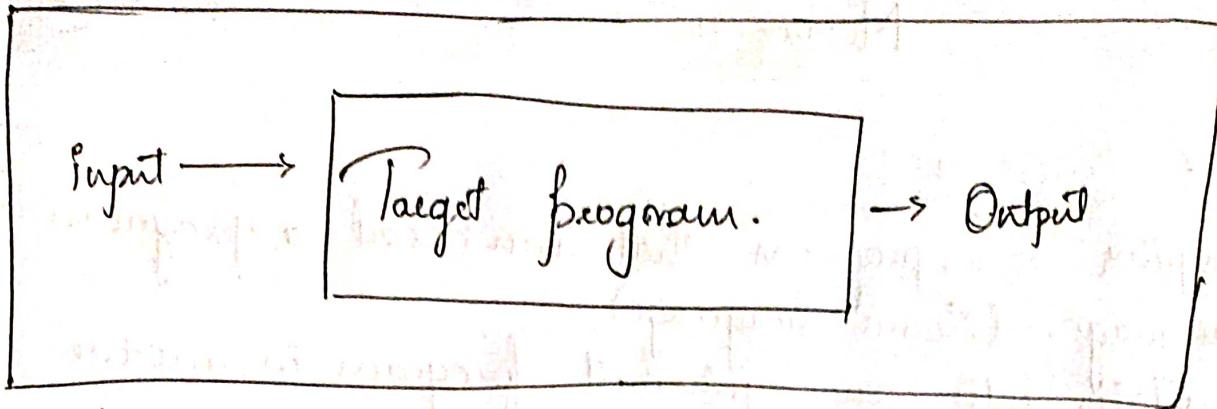


Module-IILANGUAGE PROCESSORS:

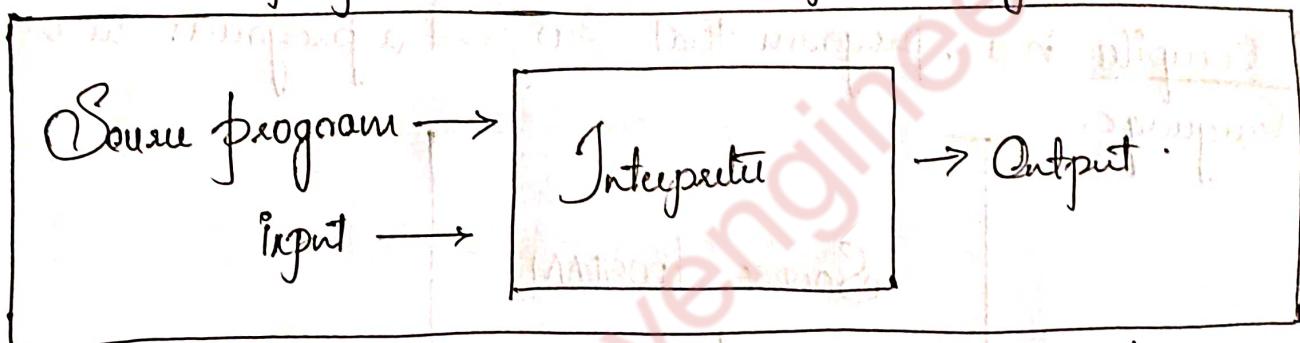
- Q A Compiler is a program that can read a program in one language. - (Source language).
- Q It translates into an equivalent program in another language. (Target language).
- Q An Important role of the Compiler is to report any errors in the Source program that it decides or detects during translation process.
- Q A Compiler is a program that can read a program in one language.



- Q If the target program is an executable machine-language program.
- Q It can then be called by the user to receive inputs and produce outputs.



- Q An Interpreter is another common kind of language processor
- Q Instead of producing a target program as a translation an Interpreter appears to directly execute the operations specified in the Source program, on inputs supplied by the user.



- Q The machine-language target program produced by a Compiler is usually much faster than an Interpreter at mapping inputs to outputs.

COMPIILER	INTERPRETER.
<ul style="list-style-type: none"> <li>Q A Compiler is a program that can read a program in a language and it is translate it into an equivalent program in another language.</li> </ul>	<ul style="list-style-type: none"> <li>Q An interpreter is another common kind of language processor. Instead of producing a target program as a translation, It appears directly to execute the operations.</li> </ul>
<ul style="list-style-type: none"> <li>Q The Machine - Language target program produced by a Compiler is usually much faster than Interpreter.</li> </ul>	<ul style="list-style-type: none"> <li>Q Interpreter is slower in Mapping inputs to outputs.</li> </ul>

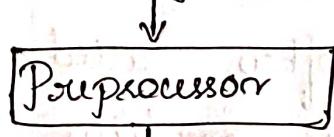
↳ Error diagnostics is less than Participant

↳ An Interpreter is slower however, can usually give better error diagnostics than a Compiler, because it executes the Source program statement by statement.

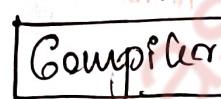
~~Language~~

## LANGUAGE PROCESSING

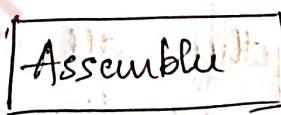
Source processing



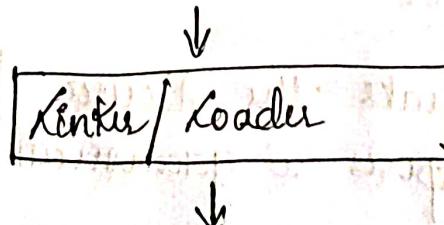
Modified Source program



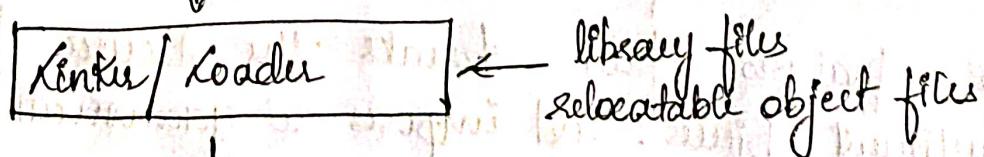
Target Assembly program



Relocatable machine code



Target machine code



- Q A source program may be divided into modules stored in a separate file.
- Q The task of collecting the source program is sometimes entrusted to a separate program. Called as Preprocessor.
- Q The preprocessor may also expand shorthands, called macros into source language statements.
- Q The modified source program is then fed to a compiler.
- Q The compiler tries to produce as op and easier to debug.
- Q The assembly language is thus processed by a program called an assembler. that produces relocatable machine code as its op.
- Q Large programs are often called as or compiled in pieces, so the relocatable machine code as its output.
- Q May have to be linked together with other relocated object files and library files into the code that actually runs on the Machine.
- Q The linker resolves categorical memory addresses, where the code in one file may refer to the location in another file.
- Q The loader then puts together all of the executable object files into memory for execution.

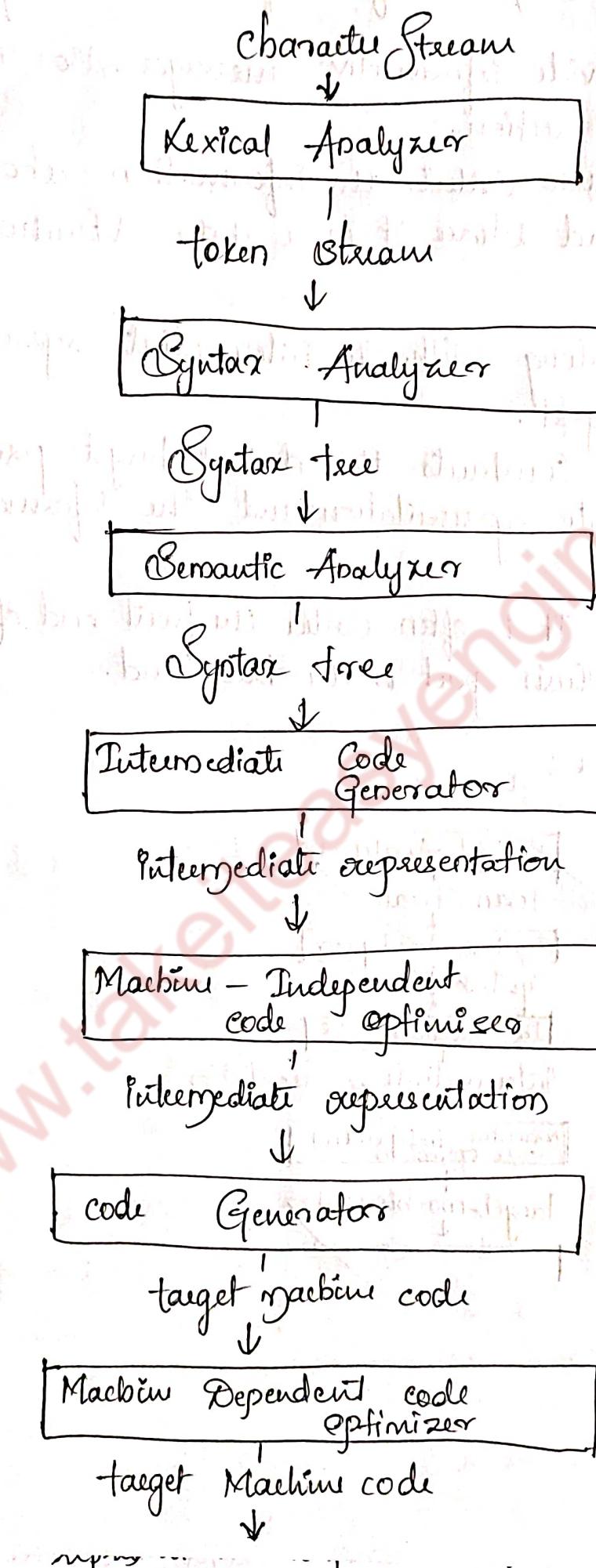
## THE STRUCTURE OF A COMPILER

- Q The analysis part breaks the source program into constituent pieces and imposes a grammatical structure on them.
- Q It then uses this structure to create an intermediate representation of the source program.

- ④ If the analysis part detects that the source program is either Semantically / Syntactically ill formed or Semantically Unsound.
- ⑤ Thus it must provide informative messages, so the user can take corrective action.
- ⑥ The analysis part also collects the information about the source program and stores it in a data structure called a Symbol table.
- ⑦ which is passed along with the intermediate representation to the Synthesis part.
- ⑧ The Synthesis part constructs the desired target program from the intermediate representation and the information in the symbol table.
- ⑨ The analysis part: It is often called the front end of the compiler the Synthesis part is the back end..

## PHASES OF COMPILER

(6)



## Lexical Analysis:

(F)

- The first phase of a Compiler is called lexical analysis or Scanning.
- The Lexical analyzer reads the stream of characters making up the source and groups the characters into meaningful sequences called lexemes.
- For each, lexeme the Lexical analyzer produces as output a token of the form that it passes on to the subsequent phase, Syntax analysis, and the Second. Component attribute value points to an entry in the symbol table for this token.
- Information from the symbol-table entry is needed for Semantic analysis and code generation.

For Example,

Suppose a C source program contains the assignment statement.

position = initial + rate \* 60

Shows the representation of the assignment statement after lexical analysis as the sequence of tokens.

(cid, 1) (=) (cid, 2) (+) (cid, 3) (\*) (60)

## Syntax Analysis:

- The second phase of the Compiler is syntax analysis or parsing
- The parser uses the first components of the tokens produced by the lexical analyzer to create a tree-like intermediate representation that depicts the grammatical structure of the token stream.
- A typical representation is a Syntax Tree in which each interior node represents an operation and the children of the node represent the arguments of the operations.

## SEMANTIC ANALYSIS

- ↳ The Semantic analyzer uses the Syntax tree and the information in the symbol table to check the source programs for Semantic Consistency with the language definition.
- ↳ It also gathers type information and saves it in either the Syntax tree or the Symbol table, for subsequent use during Intermediate-code generation.
- ↳ An important part of Semantic analysis is type checking, where the Compiler checks that each Operator has matching Operands.

### For Example:

- ↳ Many programming language definitions require an array index to be an integer.
- ↳ the Compiler must report an error if a floating point number is used to index an array.
- ↳ The Language Specification may permit some type conversions called coercions.

## INTERMEDIATE CODE GENERATION:

- ↳ After Syntax and Semantic analysis of the source program,
- ↳ Many Compilers generates an explicit low-level or machine like Intermediate representation. ~~should have two~~
- ↳ which we can think of as a program for an abstract machine.
- ↳ This Intermediate representation should have two important properties :
  - (i) It should be easy to produce and
  - (ii) It should be easy to translate into the target machine

- ↳ An intermediate form called three-address code.
- ↳ which consist of a sequence of assembly-like instructions with three operands per instruction.
- ↳ Each operand can act like a register.
- ↳ The output of the intermediate code generator consists of a sequence of assembly-like instructions with three operands per instruction.
- ↳ Each operand can act like a register.
- ↳ The output of the intermediate code generator consists of the three-address code sequence.
- ↳ There are several points worth noting about three-address instructions.
  - ↳ First, each three address assignment instruction has at most one operator on the right side.
  - ↳ Thus, these instructions fix the order in which operations are to be done.
  - ↳ The multiplication precedes the addition in the source program.
  - ↳ Second, the compiler must generate a temporary name to hold the value computed by a three-address instruction.
  - ↳ Third, some "three-address instructions" like the first and last in the sequence

### Code OPTIMIZATION:

- ↳ The Machine ~~dependent~~ independent code-optimization phase attempts to improve the intermediate code so that better target code will result.
- ↳ Usually better means faster, but other objectives may be desired, such as shorted code, or target code that consumes less power.

## For Example:

Q A straight forward algorithm generates the intermediate code, using an instruction for each operator in the tree representation that comes from the semantic analysis.

$$Q. t_1 = id_3 * 60.0$$

$$id_1 = id_2 + t_1$$

## Code Generation:

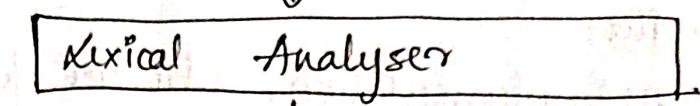
- Q The code generator takes as input an intermediate representation of the source program and maps it into the target language.
- Q If the target language is machine code, registers or memory locations are selected for each of the variables used by the program.
- Q Then, the intermediate instructions are translated into sequences of machine instructions that performs the same task.
- Q A crucial aspect of code generation is the judicious assignment of registers to hold variables.

## SYMBOL - TABLE MANAGEMENT

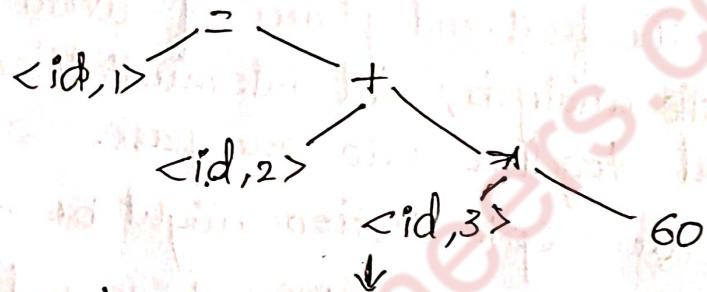
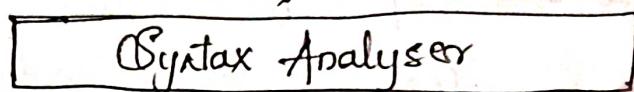
- Q The symbol table is a data structure containing a record for each variable name, with fields for the attributes of the name.
- Q The data structure should be designed to allow the compiler to find the record for each name quickly.

Position = initial + rate \* 60

u-

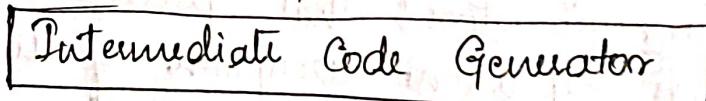
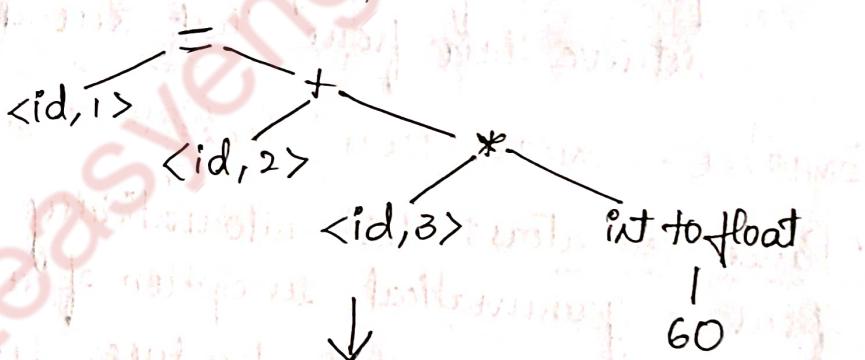


`<id ,1><=><id ,2><+><id .3><*><60>`



1	position	• • •
2	initial	• • •
3	rate	• • •

## SYMBOL TABLE



$t_1 = \text{int} \downarrow \text{to float}(60)$

$$t_2 = id_3 * t_1$$

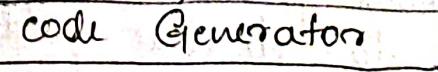
$$t_3 = \text{id}_2 + b_2$$

$$id_1 = b_3$$



$$h = id_3 * g_{0,0}$$

$$id_1 = id_2 + t$$



LDF R<sub>2</sub>, id<sub>3</sub>  
MULF R<sub>2</sub>, R<sub>2</sub>, #60.0  
LDF R<sub>1</sub>, id<sub>2</sub>  
ADD R<sub>1</sub>, R<sub>1</sub>, R<sub>2</sub>  
STF id<sub>1</sub>, R<sub>1</sub>.

## THE GROUPING OF PHASES INTO PASSES

(12)

- Q The discussion of phases deals with the logical organisation of a compiler.
- C In an implementation, activities from several phases may be grouped together into a pass that reads an input file and writes an output file.

### For Example:

- C The front-end phases of lexical analysis, syntax analysis, semantic analysis, and intermediate code generation might be grouped together into one pass.
- C Code optimization might be an optional pass.
- C Thus there could be a back-end pass consisting of code generation for a particular target machine to retrieve data from that record quickly.

## COMPILER CONSTRUCTION TOOLS

1. Parser Generators: that automatically produce syntax analyzers from a grammatical description of a programming language.
2. Scanner Generators: that produce lexical analyzers from a regular-expressions description of the tokens of a language.
3. Syntax-directed translation: engines that produce collection of routines for walking a parse tree and generating intermediate code.
4. Code-generator generators: that produce a code generator from a collection of rules for translating each operation of the intermediate language into the machine language for a target machine.
5. Data-flow analysis engines: that facilitates the gathering of information about how values are transmitted from one part of a program to another part.

## The Evolution Of Programming Languages:

### C The move to higher-level language:

- The first step towards more people-friendly programming languages was the development of mnemonic assembly languages in early 1950's.
- Initially, the instructions in an assembly language were just mnemonic representation of machine instructions.
- Later, macro instructions were added to assembly language so that a programmer could define parameterized shorthands for frequently used sequences of machine instructions.

### C Impacts On Compilers:

- Compilers can help promote the use of high-level languages by minimizing the execution overhead of the programs written in these languages.
- Compilers are also critical in making high-performance computer architectures effective on user applications.

G

## The Science Of Building A Computer

- A compiler must accept all source program that conform to the specification of the language.
- The set of source programs is infinite and any program can be very large, consisting of possibly millions of lines of code.
- Any transformation performed by the compiler while translating a source program must preserve the meaning of the program being compiled.

# APPLICATIONS OF COMPILER TECHNOLOGY

(14)

## Implementations of high-level programming languages:

- A high-level programming language defines a programming abstraction.
- The programmer expresses an algorithm using the language, and the compiler tool must translate that program to the target language.
- Optimizing Compilers include techniques to improve the performance of generated code, thus offsetting the in-efficiency introduced by high-level

## Optimizations for Computer Architectures

- The rapid evolution of Computer Architectures has also led to an insatiable demand for new Compiler technology.
- Almost all high-performance systems take advantage of the same two techniques: parallelism and memory hierarchies.
- Parallelism can be found at several levels:
  - At instruction level
  - The processor level.

### Parallelism:

- All modern microprocessors exploit instruction-level parallelism.
- Programs are written as if all instructions were executed in sequence.
- The hardware dynamically checks for dependencies in the sequential instructions stream and issues them in parallel when possible.
- In some cases, the machine includes a hardware scheduler that can change the instruction ordering to increase the parallelism in the program.

## Memory Hierarchies:

- A memory hierarchy consists of several levels of storage with different speeds and sizes, with the level closest to the processor being the fastest but smallest.
- The average memory-access time of a program is reduced if most of its access are satisfied by the faster levels of the hierarchy.
- Both parallelism and the existence of a memory hierarchy improve the potential performance of a machine, but they must be harnessed effectively by the compiler to deliver real performance on an application.

## Design of new Computer architectures:

- In the early days of computer architectures design, compilers were developed after the machines were built.
- Since programming in high-level languages is the norm, the performance of a computer system is determined not by its raw speed but also by how well compilers can exploit its features.

## Program translations:

- While we normally think of compiling as a translation from a high-level language to the machine level.
- The same technology can be applied to translate between different kinds of languages.

## Software productivity tools

- Programs are arguably the most complicated engineering artifacts ever produced.
- They consist of many details, every one of which must be correct before the program will work correctly.

## LEXICAL ANALYSIS:

- Q Lexical analysis read characters from left to right and groups into tokens.
- Q A simple way to build lexical analyzer is to construct a diagram to illustrate the structure of tokens of the source program.
- Q This approach makes it easier to modify a lexical analyzer since we have only to rewrite the affected patterns, not the entire program.

Three important approaches for implementing lexical analyzer are:

- (i) Use lexical analyzer generator (lex) from a regular expression based specification that provides routines for reading and buffering the input.
- (ii) Write lexical analyzer in conventional language using I/O facilities to read input.
- (iii) Write lexical analyzer in assembly language and explicitly manage the reading of input.

## THE ROLE OF THE LEXICAL ANALYSER.



- Q Since the lexical analyzer is the part of the compiler that reads the source text.

- Q It may perform certain other tasks besides identification of lexemes.

- Q One such task is stripping out comments and whitespace (blank, newline, tab and perhaps other characters that are used to separate tokens in the input).

- (T-1)
- Q Another task is correlating error messages generated by the Compiler with the Source program.
  - Q For instance, the lexical analyzer may keep track of the number of newline characters seen, so it can associate a line number with each error message.
  - Q In some Compilers, the lexical analyzer makes a copy of the source program with the error messages inserted at the top appropriate positions.
  - Q If the source program uses a macro-processor, the expansion of macros may also be performed by the lexical analyzer.
  - Q It is the first phase of a Compiler. It reads source code as input and sequence of tokens as output.
  - Q This will be used as input by the parser in Syntax analysis.
  - Q Upon receiving 'getNextToken' from parser, lexical analyzer searches for the next token.

### TOKENS PATTERNS AND LEXEMES :

- Q Token is a terminal symbol in the grammar for the source language.
- Q When the character sequence 'pi' appears in the source program, a token representing identifier is returned to the parser.
- Q A token name is an abstract symbol representing a kind of lexical unit.  
Eg: a particular keyword, or a sequence of input characters denoting an identifier.
- Q The token names are the input symbols that the parser processes.
- Q We often refer to a token by its token name.

- C) Pattern is a rule describing the set of lexemes that can represent a particular token in source programs.
- C) A pattern is a description of the form that the lexemes of a token may take.
- C) In the case of a keyword as a token, the pattern is just the sequence of characters that form the keyword.
- C) For identifiers and some other tokens, the pattern is a more complex structure that is matched by many strings.

### LEXEME:

- C) Lexeme is a sequence of characters in the source program that is matched by the pattern for a token.
- C) A lexeme is a sequence of characters and is identified by the lexical analyzer as an instance of that token.

In programming languages, the following classes cover most or all of the tokens:

1. One token for each keyword. The pattern for a keyword is the same as the keyword itself.
2. Tokens for the  $\neq$  operators, either individually or in classes such as the token Comparison.
3. One token representing all identifiers.
4. One or more tokens representing constants, such as numbers and literal.
5. Tokens for each punctuation symbol, such as left and right parenthesis, comma and semicolon.

Ex: Identify the lexemes and write the pattern and token for  $a=b*10$ .

Lexemes are  $a=$ ,  $b$ ,  $*$ ,  $10$

pattern for identifiers = letter followed by letter or digit (or)  
letter (letter | digit)\*

Pattern for numbers = any digits 0-9 (or) digit (+)

Token for identifiers is  $\langle id \rangle$ , points to symbol table entry

Tokens for numbers =  $\langle \text{Num}, \text{Value} \rangle$

### ATTRIBUTES Of TOKENS:

- Each token has only a single attribute - a pointer to the symbol-table entry in which the information about the token is kept.
- The token names and associated attribute values for the statement.

$$E = M * C ** 2.$$

are written below as a sequence of pairs

$\langle id, \text{pointer to symbol-table entry for } E \rangle$

$\langle \text{assign\_op} \rangle$

$\langle id, \text{pointer to symbol-table entry for } M \rangle$

$\langle \text{mult\_op} \rangle$

$\langle id, \text{pointer to symbol-table entry for } C \rangle$

$\langle \text{exp\_op} \rangle$

$\langle \text{number, integer value } 27 \rangle$

### LEXICAL ERRORS:

- It is hard for a lexical analyser to tell, without the aid of other components, that there is a source-code error.

For instance, if the string  $f_i$  is encountered for the first time in a program in the context:

$f_i (a == f(x))$

- A lexical analyser cannot tell whether  $f_i$  is a misspelling of the keyword  $if$  or an undeclared function identifier.

Since  $f_i$  is a valid lexeme for the token  $id$ , the parser and

let some other phase of the compiler.

- Probably the parser is this case - handle an error due to transposition of the letters.

However, suppose a situation arises in which the lexical analyser is unable to proceed because none of the patterns for tokens matches any prefix of the remaining input.

- The simplest recovery strategy is "panic mode" recovery, i.e. we delete successive characters from the remaining input, until the lexical analyser can find a well-formed token at the beginning of what input is left.
- This recovery technique may confuse the parser, but in an interactive computing environment it may be quite adequate.

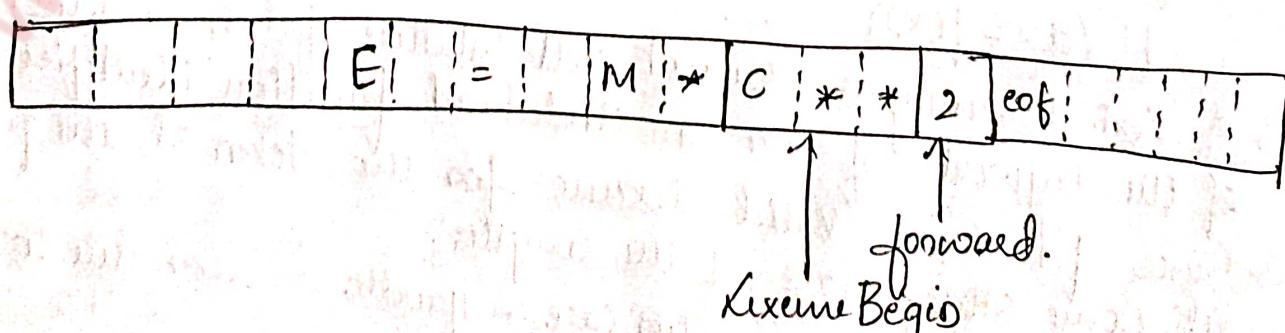
Other possible error-recovery actions are:

1. Delete one character from the remaining input.
2. Insert a missing character into the remaining input.
3. Replace a character by another character.
4. Transpose two adjacent characters.

### INPUT BUFFERING:

- Before discussing the problem of recognizing lexemes in the input, let us examine some ways that the simple but important task of reading the source program can be speeded.
- This task is made difficult by the fact that we often have to look one or more characters beyond the next lexeme before we can be sure we have the right lexeme.

Buffer pairs:



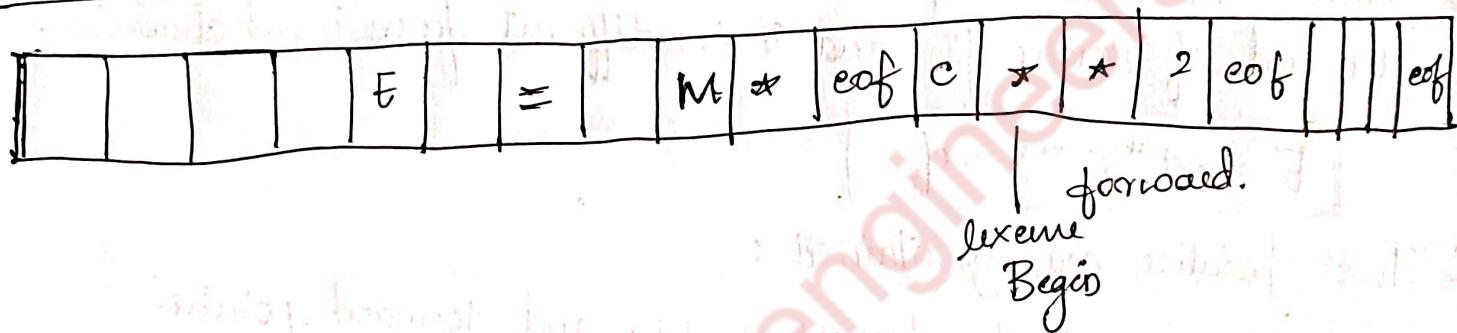
- Because of the amount of time taken to process characters and the large number of characters that must be processed during the compilation of a large source program. 21
- Specialised buffering techniques have been developed to reduce the amount of overhead required to process a single input character.
- Look ahead of characters in the input is necessary to identify tokens.
- Specified buffering techniques have been developed to reduce the large amount of time consumed in moving characters.
- A buffer (array) divided into two N-character halves.
- where  $N = \text{number of characters on one disk block}$  'eof' marks the end of source file and it is different from input character.

$$E = M * C ** 2 \text{eof}.$$

- Two pointers are maintained :
  - beginning of the lexeme pointer and forward pointer.
- both pointers point to the first characters of the next lexeme to be found.
- Forward pointer scans ahead until a match for a pattern is found.
- Once the next lexeme is determined, processed and both pointers are set to the character immediately past the lexeme.
- If the forward pointer moves halfway back, the right N half is filled with new characters.
- If the forward pointer moves right end of the buffer then left N half is filled with new characters.
- The disadvantage is look ahead is limited and it is impossible to recognise tokens, when distance between the two pointers is more than the length of the buffer.

- 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 the input, and move forward to the beginning of the newly loaded buffer.
- As long as we never need to look so far ahead of the actual lexeme that the sum of the lexeme's length plus the distance we look ahead is greater than  $N$ , we shall never overwrite the lexeme in its buffer before determining it.

### SENTINELS:



- It is an extra key inserted at the end of the array.
- It is special, dummy character that can't be part of your program.
- With respect to buffer pairs, the code for advancing forward pointer is

code, with sentinels

Switch (\* forward ++)

{

case eof :

if (forward is at end of first buffer)

{

reload second buffer;

forward = beginning of second buffer;

{

else if (forward is at end of second buffer)

{

reload first buffer;

forward = beginning of first buffer;

{

else /\* eof within a buffer marks the end of input \*/  
terminate lexical analysis;

break

Cases for the other characters.

3.

### SPECIFICATION OF TOKENS

- Q An Alphabet is any finite set of symbols.
- Q Typical examples of symbols are letters, digits, and punctuation
- Q The set  $\{0,1\}$  is the binary Alphabet
- Q A string over an Alphabet is a finite sequence of symbols drawn from that Alphabet
- Q In language theory, the terms "Sentence" and "word" are often used as synonyms for "string".

The length of string  $S$ , usually written  $|S|$ , is the number of occurrences of symbols in  $S$ .

For Example:

Banana is a string of length six, The empty string denoted  $e$ , is the string length of zero.

The following string-related terms are commonly used:

1. A prefix of string  $S$  is any string obtained by removing zero or more symbols from the end of  $S$ .

For Ex:

= ban, banana, and  $e$  are prefixes of banana.

2. A suffix of string  $S$  is any string obtained by removing zero or more symbols from the beginning of  $S$ .

ForEx :

= nana, banana and  $e$  are suffixes of banana.

3. A substring of  $S$  is obtained by deleting any prefix and any suffix from  $S$ .

For instance,

banana, nan, and  $e$  are substrings of banana.

4. The proper prefixes, suffixes, and substrings of a string  $S$  are those, prefixes, suffixes and substrings respectively, of  $S$  that are not  $e$  or not equal to  $S$  itself.

5. A subsequence of  $S$  is any string formed by deleting zero or more, not necessarily consecutive positions of  $S$ .

for ex : baan is a subsequence of banana.

String ( $S$ ) : Sentence/word : finite set/sequence of symbols.

Q  $|S|$  = number of symbols in string  $S$

eg:  $S = \text{banana}$ , then  $|S|=6$

Prefix ( $S$ ):  $S = \text{banana}, \text{ba}, \text{ban}, \text{banan}, \text{banan}$

Suffix ( $S$ ):  $S = \text{banana}, \text{ana}, \text{nana}, \text{anana}, \text{nana}$

Substring ( $S$ ):  $S = \text{banana}, \text{ba}, \text{na}, \text{na}$

$\_\_$ : empty string,  $S = \_\_$  then  $|S|=0$ ,  $\emptyset$  = empty set.

Language: set of strings.

If  $X$  and  $Y$  are strings then  $XY$  is concatenation.

$$SE = ES = S$$

$$S\emptyset = G$$

$$S_1 = S$$

$$S_2 = SS$$

$$S_3 = SSS$$

$$S^i = S^{i-1}S \quad (\text{if } i > 0)$$

Q A language is any countable set of strings over some fixed alphabet.

Q This definition very broad.

Q Abstract languages like  $\emptyset$ , the empty set, or  $\{\epsilon\}$ , the set containing only the empty string, all languages under this definition.

Q C programs and the set of all grammatically correct English sentences, although the latter two languages are difficult to specify exactly.

LUD: Union operation, where  $\lambda$  = set of alphabet  $\{A, Z, a..z\}$  and  $\emptyset$  = set of digits  $\{0..9\}$

LD: Concatenation

LA: exponentiation: set of strings with  $A$  letters.

$$LO = E$$

$$L^i = L^{i-1}L$$

$L^*$  = all strings with  $E$ : Kleen closure of  $L$ .

: Set of all strings of digits of one or more.

## Operations On Languages

1.  $\Sigma^0$  is the set of letters or digits - strictly speaking the language with 62 strings of length one, each of which strings is either one letter or one digit
2.  $\Sigma^1$  is the set of strings of length two, each consisting of one letter followed by one digit
3.  $\Sigma^4$  is the set of all 4-letter strings
4.  $\Sigma^*$  is the set of all strings of letters, including the empty string
5.  $L(\Sigma^0)^*$  is the set of all strings of one or more digits.
6.  $\Sigma^1$  is the set of all strings of one or more digits.

## Regular Expressions

It is a notation which allows defining the sets precisely.

Eg:  $L(\Sigma^0)^*$  if regular expression is:  
letter (letter/digit)\*

Regular expression over alphabet has following data or rules.

If  $a$  is a regular expression, the set containing empty string  $a$  is a regular expression.

If  $a$  belongs to  $\Sigma$  this is  $\{a\}$  set containing the string  $a$ .

Suppose  $r$  and  $s$  are regular expression.

denoting the languages  $L(r)$  and  $L(s)$  then

$(r \cup s)$  is a regular expression denoting  $L(r) \cup L(s)$

$rs$  is a regular expression denoting  $L(r) L(s)$

$(r)^*$  is a regular expression is said to be a regular set, concatenation and  $*$  have highest to lowest precedence with left associative

If two regular expressions ' $r$ ' and ' $s$ ' denote the same language we say ' $r$ ' and ' $s$ ' are equivalent and write  $r \equiv s$ .

## RECOGNITION OF TOKENS

(N-7)

Q So, let's see how to take the patterns for all the needed tokens and built a piece of code that examines that input string and finds a prefix that is a lexeme matching one of the patterns.

Q Consider the following grammar, fragment/regular definitions:

stmt → if expr then stmt

expr → num relop terms

terms → id | num

if → if

then → then

else → else

relop → < | <= | = | > | >= | = | <>

letter → [a -> A - Z]

digit → [0 - 9]

id → letter (letter|digit)\*

digits → digit+

num → digits (. digits)? (E (+|-) ? digits)?

Q Keywords Cannot be used as identifiers

Num represents Unsigned int and real numbers.

Q The regular definition, ws,

delim → blank | tab | newline

ws → delim+

Q If 'ws' was found, the lexical analyzer does not return a token to the parser.

Q Our Goal is to Construct a lexical analyzer to produce a pair consisting of token and attribute-value as output using the Transition table

## Lexemes

## TOKEN NAME

## ATTRIBUTE VALUE

Any ws	-	-
if	f	-
then	then	-
else	else	-
Any id	id	points to table entry
Any Number	number	Points to table entry
<	relOp	LT
<=	relOp	LE
=	relOp	EQ
<>	relOp	NE
>	relOp	GT
>=	relOp	

## TRANSITION DIAGRAMS

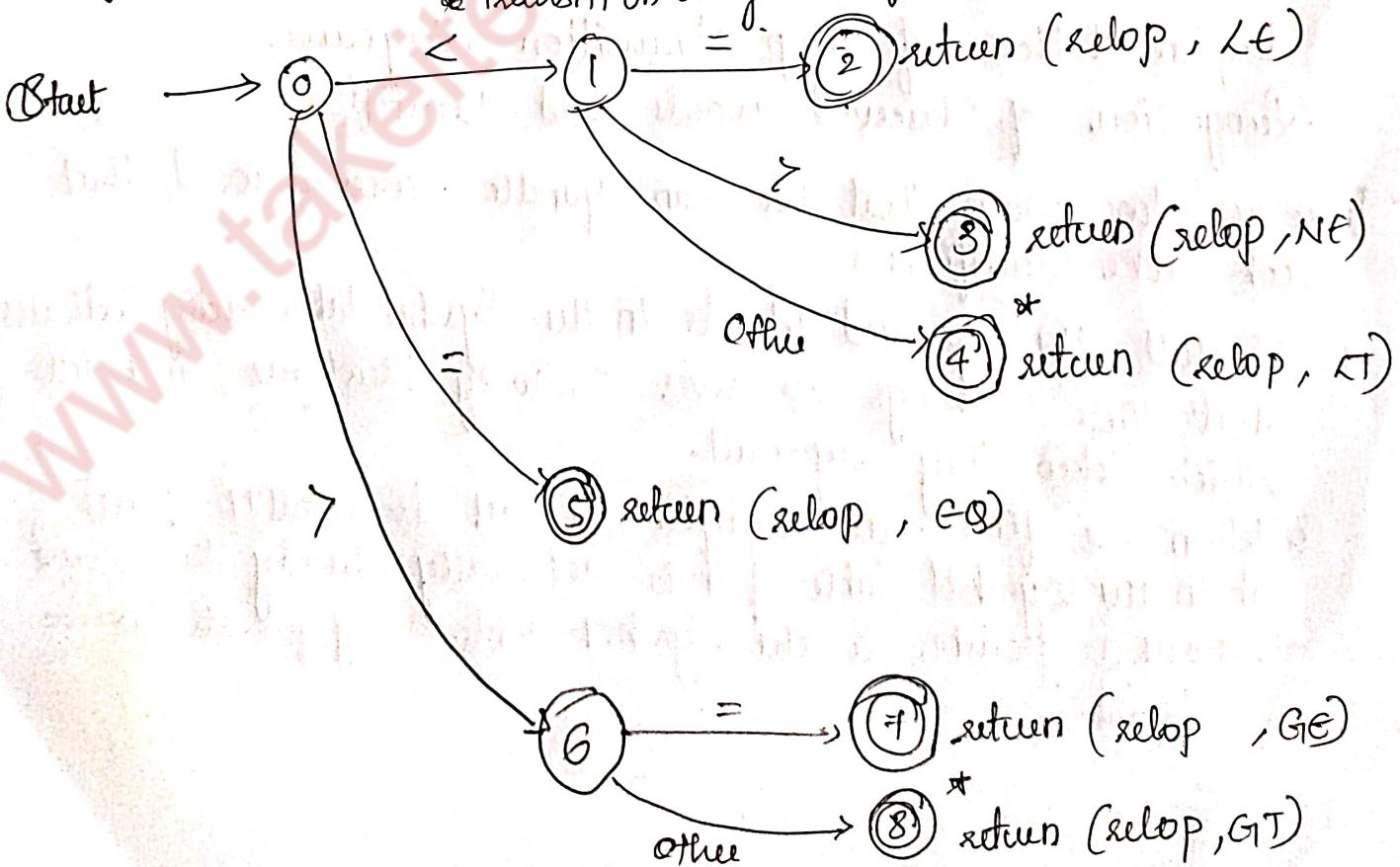
- These are the flowcharts, as an intermediate step in the construction of a lexical analyzer.
- This takes actions like a lexical analyzer is called by the parser to get the next token.
- We use transition diagram to keep track of information about characters that are seen as and where the forward pointer scans the input.
- Lexeme beginning point points to the character following the last lexeme found.

$$T = M^* C^* \star 2 \text{eof}$$

- Transition diagrams have a collection of nodes or circles called states.

- (27)
- Some important conventions about transition diagrams are:
- Certain states are said to be accepting, or final. These states indicate that a lexeme has been found, although the actual lexeme may not consist of all positions from the lexemeBegin and forward pointers.
  - We always indicate an accepting state by a double circle, and if there is an action to be taken.
  - Typically returning a token and an attribute value to parser.
  - Shall attach that action to the accepting state.
  - If it is necessary to retract the forward pointer one position.
  - If it is necessary to retract the forward pointer one position.
  - Then we shall additionally place a \* near that accepting state.
  - One state is designated the start state, or initial state; it is indicated by an edge, labelled "start" entering from nowhere.
  - The transition state always begins in the start state before any input symbols have been read.

transition diagram of RElop



## Transition Diagram for Relop

TOKEN getRelop ()

{

TOKEN retToken = new(REFOP);

while (1)

{

On failure occurs

Switch(state)

Case 0 : c = nextChar();

If (c == '<') & state = 1;

else if (c == '=') & state = 5;

else if (c == '>') & state = 6;

else fail();

break;

Case 1 : ...

Case 8 : retreat();

retToken.setAttribute = GT;

return (retToken);

Implementation of relop transition diagram.

Recognition of Reserved words and Identifiers

There are two ways that we can handle reserved word that look like identifiers :

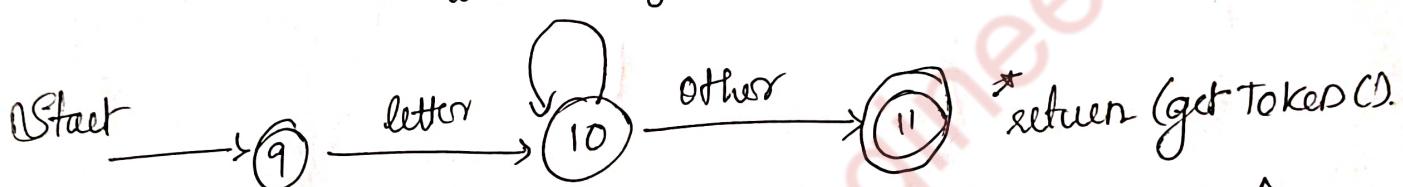
1. Install the reserved words in the symbol table entry indicates that these strings are never ordinary identifiers, and tells which tokens they represent.

When we find an identifier, a call to installID places it in the symbol table if it is not ~~early~~ already there and returns a pointer to the symbol-table entry for the lexeme found.

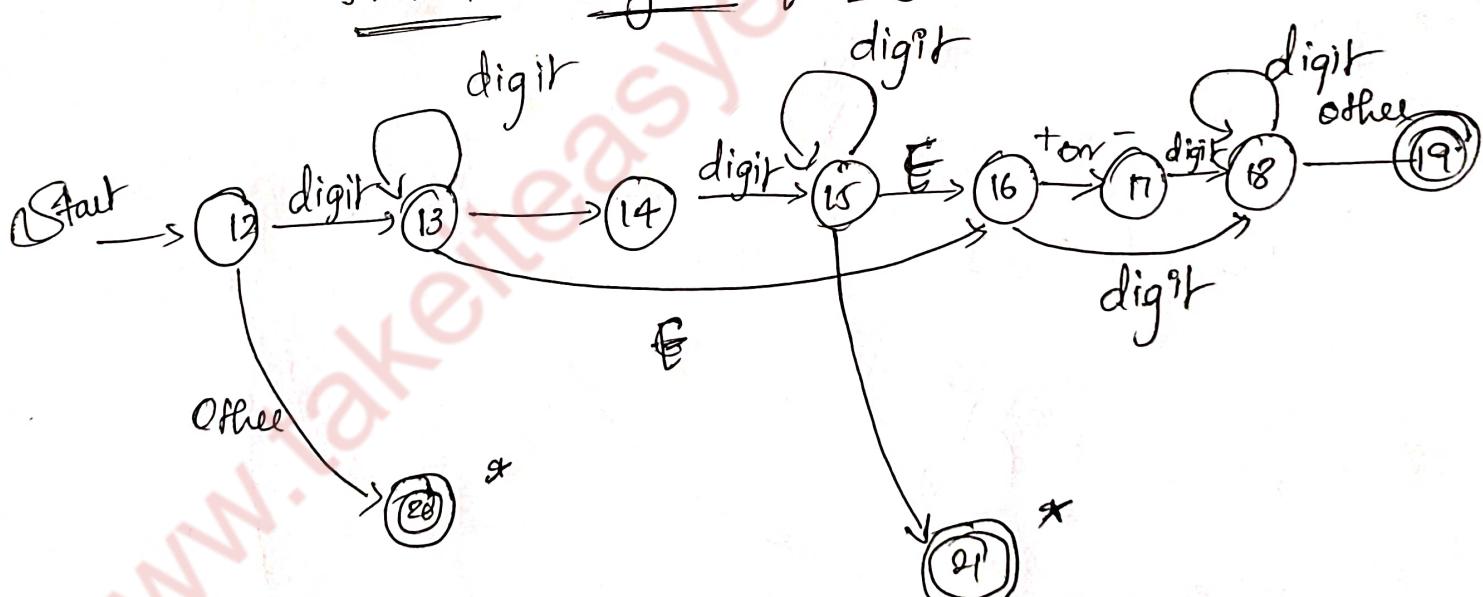
Q) Any identifier not in the symbol-table entry for the lexeme found  
 Q) either id or one of the keyword tokens that was initially installed in the table.

Q. Create separate transition diagrams for each keyword; as example for the keyword then.

Q) It is necessary to check the identifier has ended, or else we would return token then to situations where the correct token was id, with a lexeme like the next value that has then as a proper prefix.  
 letter or digit



Transition diagram of signed number



Transition diagram of white space

