# CS-402 Compiler Construction
# Lexical Analysis
# (Documenting the Design of Lexer)

## 1. Language Tokens and Lexemes

| Token | Description | Lexemes |
|-------|-------------|---------|
| INT | The data type int i.e., letters/characters i, n, t | int |
| CHAR | The data type char i.e., letters/character c, h, a, r | char |
| IF | It is a keyword made of letters i, f | if |
| ELIF | It is a keyword made of letters e, l, i, f | elif |
| ELSE | It is a keyword made of letters e, l, s, e | else |
| WHILE | It is a keyword made of letters w, h, i, l, e | while |
| INPUT | It is a keyword made of letters i, n, p, u, t | input |
| PRINT | It is a keyword made of letters p, r, i, n, t | print |

| | | |
|---|---|---|
| PRINTLN | It is a keyword made of letters p, r, i, n, t, l, n | println |
| REL_OP | Relational operators that compare two operands | < |
| | | <= |
| | | > |
| | | >= |
| | | == |
| | | ~= |
| '+' | The arithmetic operator for addition | + |
| '-' | The arithmetic operator for subtraction | - |
| '*' | The arithmetic operator for multiplication | * |
| '/' | The arithmetic operator for division | / |
| ID | Letter followed by letters, digits or underscore | p2, max_ etc. |
| NUM | Digits that form only integers | 0, 123, 99 etc. |
| LIT | A single letter/character enclosed in single quotes | 'x', 'a' etc. |
| STR | Sequence of letters/characters and while spaces enclosed in double quotes | "hello there    you" etc. |
| S_COMMENT | It's a single line comment | // comment |

| M_COMMENT | It's a multi-line/ block comment | /* comment */ |
|---|---|---|
| '=' | It is the assignment operator | = |
| INPUT_OP | It is an operator used for taking input | -> |
| ':' | It is the punctuation mark ":" | : |
| ';' | It is the punctuation mark ";" | ; |
| ',' | It is the punctuation mark "," | , |
| '(' | It is the punctuation mark "(" | ( |
| ')' | It is the punctuation mark ")" | ) |
| '{' | It is the punctuation mark "{" | { |
| '}' | It is the punctuation mark "}" | } |
| '[' | It is the punctuation mark "[" | [ |
| ']' | It is the punctuation mark "]" | ] |

All tokens having a symbol within single quotes represent ASCII values.

For single-line comments, the token isn't meant to be sent to parser rather whenever the LA reads // it will read the complete line till newline character ('\n') and ignore it or remove it by replacing by "".

Similarly, for multi-line comments the LA will everything after /* till */ and ignore or remove it.

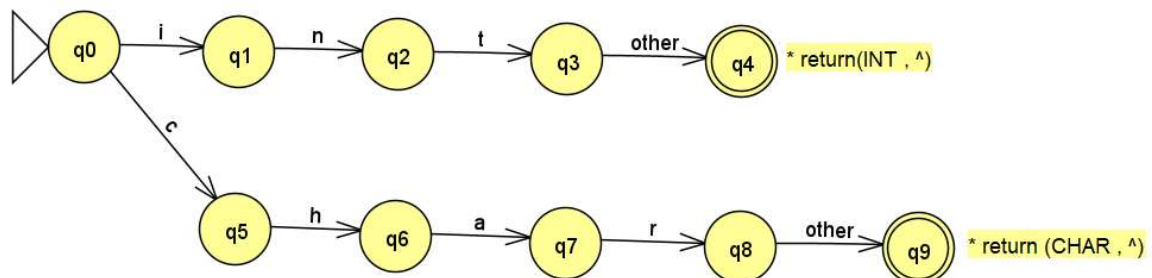For tokens having multiple lexemes, each lexeme will be given a attribute value or will be recorded in the symbol table.

## i.  Data Types: int, char

These are basically sequence of letters.

Regular Definition

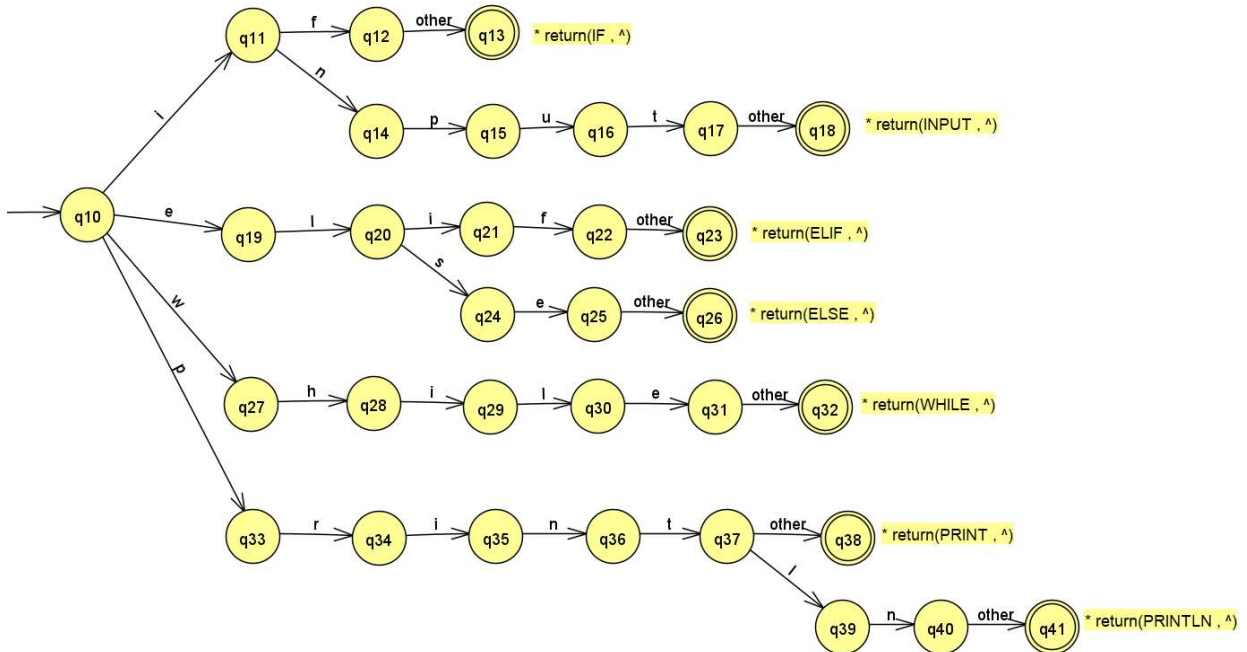$$DT \rightarrow \text{int} \mid \text{char}$$

Transition Diagram



## ii.  Keywords: if, elif, else, while, input, print, println

These keywords are basically sequence of letters as well.

Regular Definition

$$KW \rightarrow \text{if} \mid \text{elif} \mid \text{else} \mid \text{while} \mid \text{input} \mid ((\text{print})(\char94 \mid \text{ln}))$$
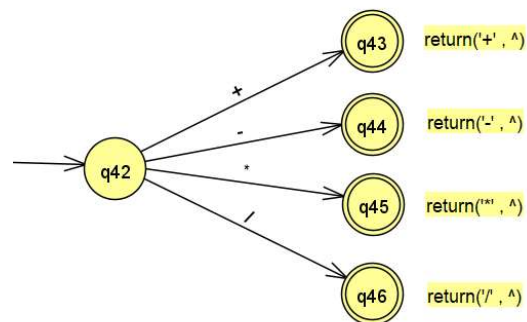
Transition Diagram

q11 —f→ q12 —other→ q13 * return(IF , ^)

q10 —i→ q11 —n→ q14 —p→ q15 —u→ q16 —t→ q17 —other→ q18 * return(INPUT , ^)

q10 —e→ q19 —l→ q20 —i→ q21 —f→ q22 —other→ q23 * return(ELIF , ^)

q20 —s→ q24 —e→ q25 —other→ q26 * return(ELSE , ^)

q10 —w→ q27 —h→ q28 —i→ q29 —l→ q30 —e→ q31 —other→ q32 * return(WHILE , ^)

q10 —p→ q33 —r→ q34 —i→ q35 —n→ q36 —t→ q37 —other→ q38 * return(PRINT , ^)

q37 —l→ q39 —n→ q40 —other→ q41 * return(PRINTLN , ^)

## iii. Arithmetic Operators: +, -, *, /

Regular Definition

$$arithOp \rightarrow + \mid - \mid * \mid /$$

Transition Diagram

q42 —+→ q43 return('+' , ^)

q42 —-→ q44 return('-' , ^)

q42 —*→ q45 return('*' , ^)

q42 —/→ q46 return('/' , ^)

## iv. Relational Operators: <, <=, >, >=, ==, ~=

Regular Definition

$$relOp \rightarrow < \mid <= \mid > \mid >= \mid == \mid \sim=$$

Transition Diagram

## v.  Comments

Single-line comments are // followed by a comment which can constitute any character till newline symbol.

Multi-line comments are comment enclosed in /* */ and can constitute any character even newline symbol.

<u>Regular Definition</u>

$$letter \rightarrow A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z$$

$$digit \rightarrow 0 \mid 1 \mid \dots \mid 9$$

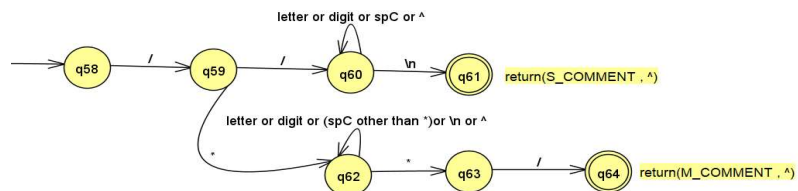$$newline \rightarrow \text{\textbackslash n}$$

$$spC \rightarrow \sim \mid @ \mid \$ \mid \% \mid \& \mid * \mid ( \mid ) \mid \{ \mid \} \mid [ \mid ] \mid + \mid = \mid \_ \mid - \mid \textbackslash \mid / \mid < \mid > \mid . \mid , \mid `` \mid " \mid ' \mid space \mid : \mid ; \mid ?$$
$$\mid \mid$$

$$S\_COMMENT \rightarrow // \; (letter \mid digit \mid spC \mid {}^\wedge)^* \; newline$$

$$M\_COMMENT \rightarrow /* \; (letter \mid digit \mid spC \mid newline \mid {}^\wedge)^* \; */$$

<u>Transition Diagram</u>

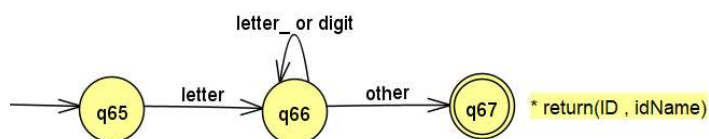**vi.** **Identifier:** a letter followed by any number of letters or digits or underscore symbol

<u>Regular Definition</u>

$$letter \rightarrow A \mid B \mid \ldots \mid Z \mid a \mid b \mid \ldots \mid z$$

$$letter\_ \rightarrow letter \mid \_$$

$$digit \rightarrow 0 \mid 1 \mid \ldots \mid 9$$

$$ID \rightarrow letter \ (letter\_ \mid digit)^*$$

<u>Transition Diagram</u>



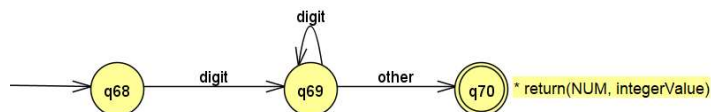Here, idName will be the name of the identifier recognized.

**vii.** **Numeric Constants:** only integers

These are sequence of digits that form only integers i.e., they don't include fraction, exponential parts etc.

<u>Regular Definition</u>

$$digit \rightarrow 0 \mid 1 \mid \ldots \mid 9$$

$$NUM \rightarrow digit^+$$

<u>Transition Diagram</u>



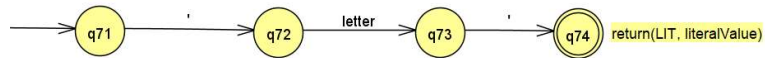Here, integerValue will be the value of the NUM recognized.

**viii.** **Literal Constants:** a letter enclosed in single quotes

<u>Regular Definition</u>

$$letter \rightarrow A \mid B \mid \ldots \mid Z \mid a \mid b \mid \ldots \mid z$$

$$LIT \rightarrow \text{ ' } (letter) \text{ '}$$

<u>Transition Diagram</u>



Here, literalValue will be the value of the LIT recognized.

ix. **Strings:** sequence of letters and white spaces enclosed in double quotes

<u>Regular Definition</u>

$$letter \rightarrow A \mid B \mid \ldots \mid Z \mid a \mid b \mid \ldots \mid z$$

$$digit \rightarrow 0 \mid 1 \mid \ldots \mid 9$$

$$newline \rightarrow \backslash n$$

$$ws \rightarrow (\text{blank} \mid \text{tab} \mid newline)^+$$

$$spC \rightarrow \sim \mid @ \mid \$ \mid \% \mid \& \mid * \mid ( \mid ) \mid \{ \mid \} \mid [ \mid ] \mid + \mid = \mid \_ \mid - \mid \backslash \mid / \mid < \mid > \mid . \mid , \mid \text{“} \mid \text{”} \mid \text{'} \mid \text{space} \mid : \mid \mid ; \mid ? \mid \mid$$

$$STR \rightarrow \text{“ } (letter \mid ws \mid spC \mid digit)^* \text{ ”}$$

<u>Transition Diagram</u>
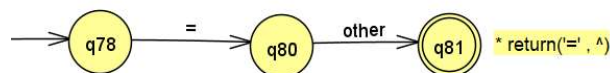


x. **Assignment Operator:** =

<u>Regular Definition</u>

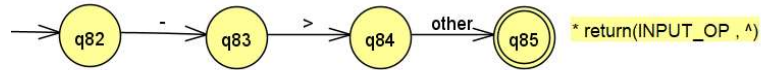$$assignOp \rightarrow =$$

<u>Transition Diagram</u>



xi. **Input Operator:** ->

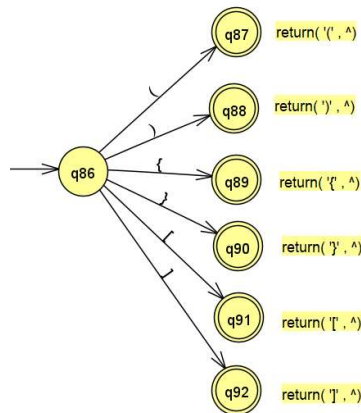$$INPUT\_OP \rightarrow \text{->}$$

Transition Diagram



## xii. Parenthesis, Braces, Square Brackets: (, ), {, }, [, ]

These are basically punctuation symbols being used as operators.

Regular Definition

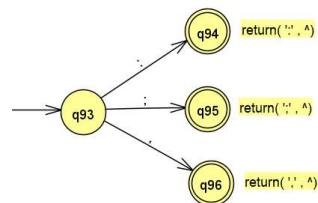$$PBC \rightarrow (\,|\,) \,|\, \{ \,|\, \} \,|\, [ \,|\, ]$$

Transition Diagram



## xiii. Colon, Semi Colon, Comma: :, ;, ,

These are basically punctuation symbols being used as operators.

Regular Definition

$$SCC \rightarrow : \,|\, ; \,|\, ,$$

Transition Diagram

## 2. Overall Regular Expressions

$letter \rightarrow$ A | B | … | Z | a | b | … | z

$letter\_ \rightarrow letter$ | _

$digit \rightarrow$ 0 | 1 | … | 9

$NUM \rightarrow digit^+$

$ID \rightarrow letter\ (letter\_\ |\ digit)^*$

$newline \rightarrow$ \n

$ws \rightarrow$ (blank | tab | $newline)^+$

$LIT \rightarrow$ ' ($letter$) '

$STR \rightarrow$ " ($letter$ | $ws$ | $spC$ | $digit)^*$ "

$INT \rightarrow$ int

$CHAR \rightarrow$ char

$IF \rightarrow$ if

$ELIF \rightarrow$ elif

$ELSE \rightarrow$ else

$WHILE \rightarrow$ while

$INPUT \rightarrow$ input

$PRINT \rightarrow$ print

$PRINTLN \rightarrow$ println

$INPUT\_OP \rightarrow$ ->

$spC \rightarrow$ ~ | @ | $ | % | & | * | ( | ) | { | } | [ | ] | + | = | _ | - | \ | / | < | > | . | , | " | " | ' | space | : | ; | ?
| |

$S\_COMMENT \rightarrow$ // ($letter$ | $digit$ | $spC$ | ^$)^*$ $newline$

$M\_COMMENT \rightarrow$ /* ($letter$ | $digit$ | $spC$ | $newline$ | ^$)^*$ */

$assignOp \rightarrow$ =

$arithOp \rightarrow$ + | - | * | /
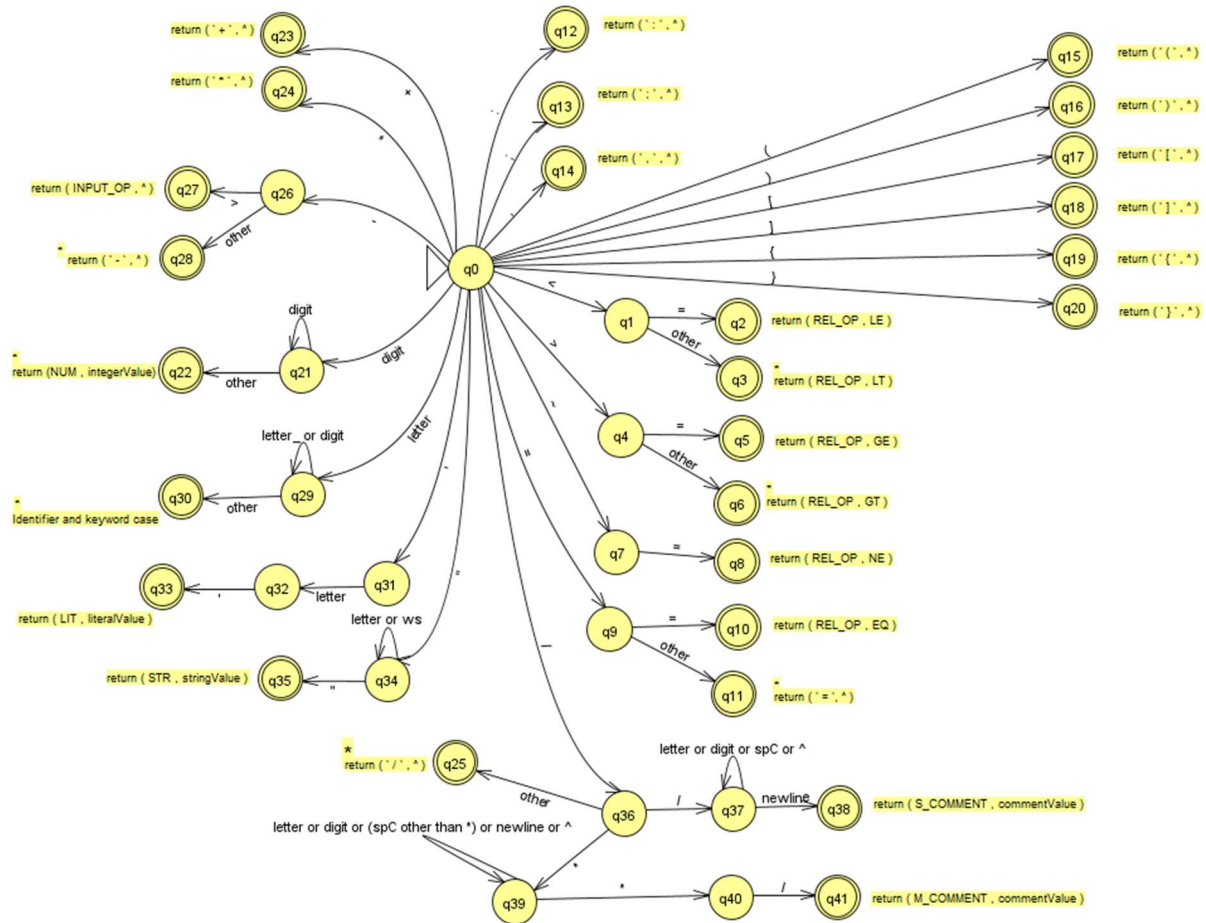
$relOp \rightarrow$ < | <= | > | >= | == | ~=

$$PBC \rightarrow ( | ) | \{ | \} | [ | ]$$

$$SCC \rightarrow : | ; | ,$$

Here, we have expanded the data types and keywords as individual tokens for clarity.

## 3. Overall Finite State Machine

Now, we will combine all the transition diagrams into one big finite state machine. It is optimized by minimizing the states to necessary ones and combines some token cases like that of identifiers and keywords.



Here, all keywords are placed into the identifier case as they all qualify as valid identifiers. Later on, keywords will be distinguished from identifiers by looking up the symbol table for reserved words.