

**CS-402 Compiler Construction**  
**Parser Analysis and Design**  
**(Documenting the Design of Parser's Grammar)**

**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 “]”	]

**Note:**

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.

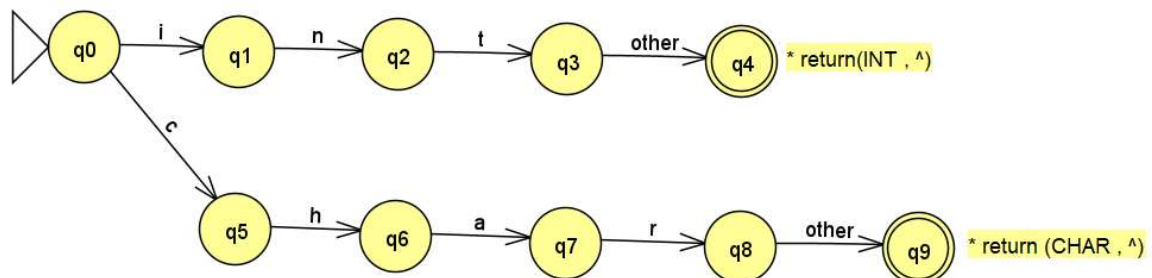
**i. Data Types: int, char**

These are basically sequence of letters.

Regular Definition

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

Transition Diagram



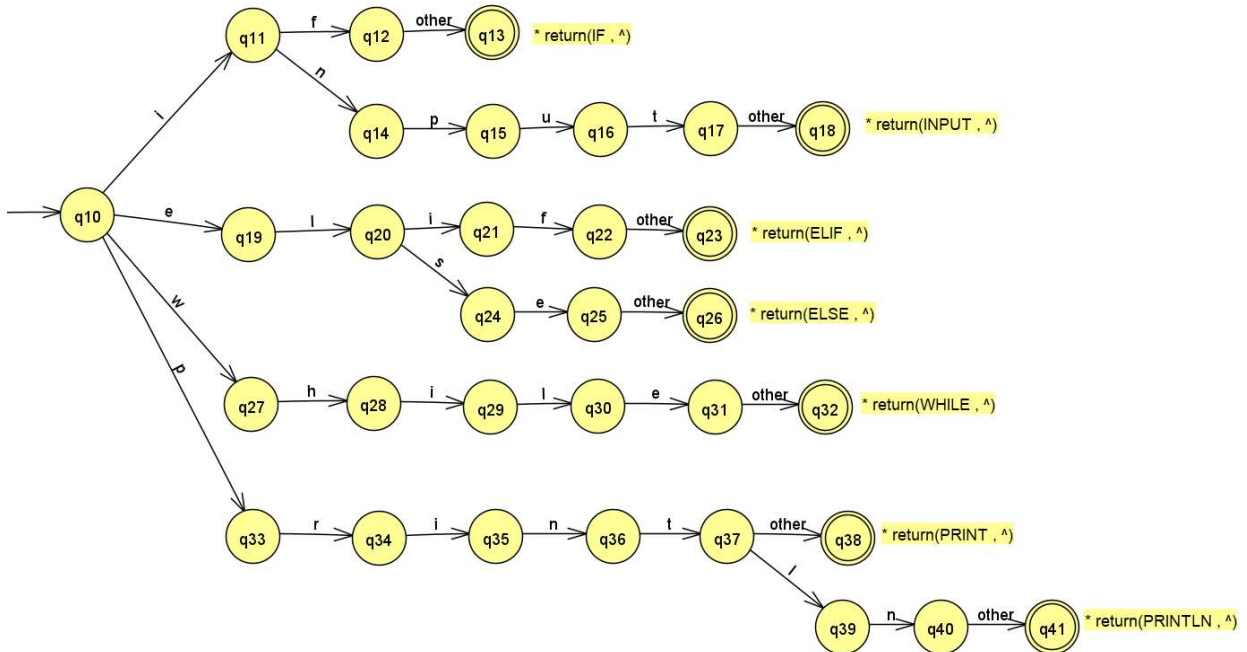
**ii. Keywords: if, elif, else, while, input, print, println**

These keywords are basically sequence of letters as well.

Regular Definition

$$KW \rightarrow \text{if} \mid \text{elif} \mid \text{else} \mid \text{while} \mid \text{input} \mid ((\text{print})^{\wedge} \mid \text{ln}))$$

Transition Diagram

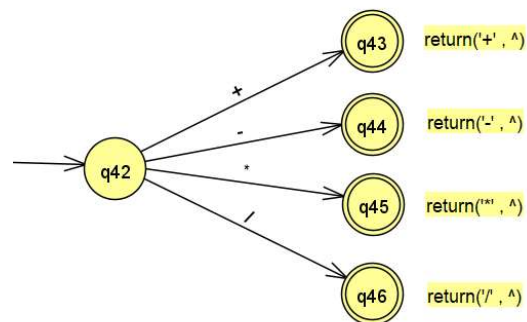


### iii. Arithmetic Operators: +, -, \*, /

Regular Definition

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

Transition Diagram

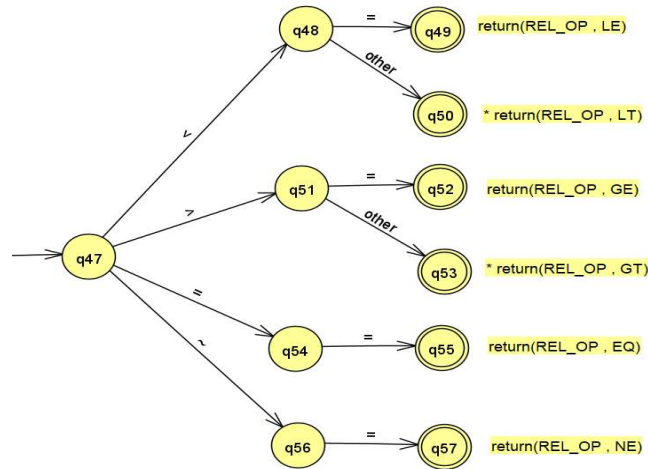


### iv. Relational Operators: <, <=, >, >=, ==, ~=

Regular Definition

$relOp \rightarrow < | <= | > | >= | == | ~=$

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.

### Regular Definition

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

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

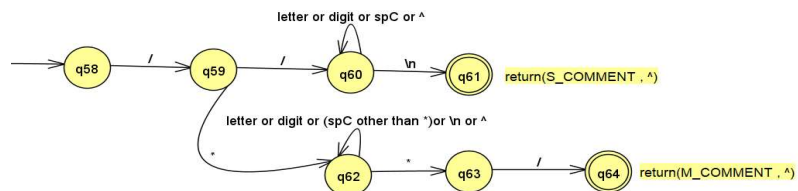
$$newline \rightarrow \backslash n$$

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

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

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

### Transition Diagram

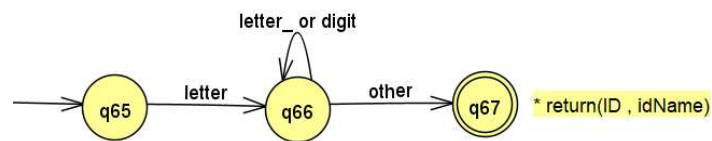


- vi. **Identifier:** a letter followed by any number of letters or digits or underscore symbol

Regular Definition

$$letter \rightarrow A | B | \dots | Z | a | b | \dots | z$$
$$letter\_ \rightarrow letter | \_$$
$$digit \rightarrow 0 | 1 | \dots | 9$$
$$ID \rightarrow letter (letter\_ | digit)^*$$

Transition Diagram



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

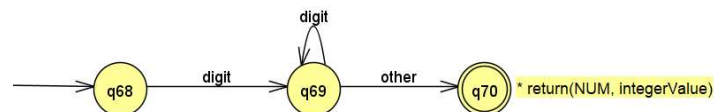
- vii. **Numeric Constants:** only integers

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

Regular Definition

$$digit \rightarrow 0 | 1 | \dots | 9$$
$$NUM \rightarrow digit^+$$

Transition Diagram



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

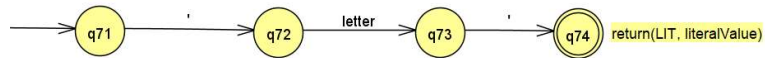
- viii. **Literal Constants:** a letter enclosed in single quotes

Regular Definition

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

$LIT \rightarrow ' (letter) '$

#### Transition Diagram



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

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

#### Regular Definition

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

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

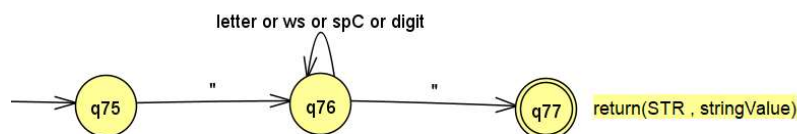
$newline \rightarrow \backslash n$

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

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

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

#### Transition Diagram

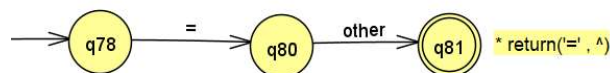


**x. Assignment Operator: =**

#### Regular Definition

$assignOp \rightarrow =$

#### Transition Diagram



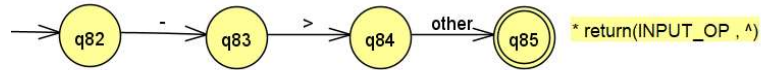
**xi. Input Operator: ->**



### Regular Definition

$INPUT\_OP \rightarrow ->$

### Transition Diagram



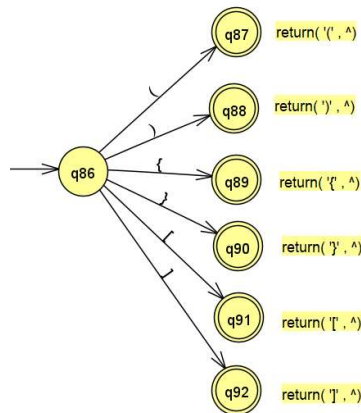
## xii. **Parenthesis, Braces, Square Brackets: (, ), {, }, [, ]**

These are basically punctuation symbols being used as operators.

### Regular Definition

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

### Transition Diagram



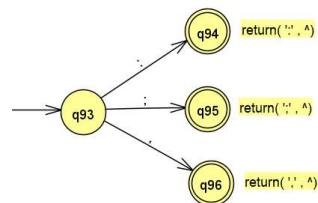
## xiii. **Colon, Semi Colon, Comma: :, ;, ,**

These are basically punctuation symbols being used as operators.

### Regular Definition

$SCC \rightarrow : | ; | ,$

### Transition Diagram



## 2. Overall Regular Expressions

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

$letter\_ \rightarrow letter | \_$

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

$NUM \rightarrow digit^+$

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

$newline \rightarrow \backslash n$

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

$LIT \rightarrow ' (letter) '$

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

$INT \rightarrow int$

$CHAR \rightarrow char$

$IF \rightarrow if$

$ELIF \rightarrow elif$

$ELSE \rightarrow else$

$WHILE \rightarrow while$

$INPUT \rightarrow input$

$PRINT \rightarrow print$

$PRINTLN \rightarrow println$

$INPUT\_OP \rightarrow ->$

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

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

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

$assignOp \rightarrow =$

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

$relOp \rightarrow < | <= | > | >= | == | ~=$

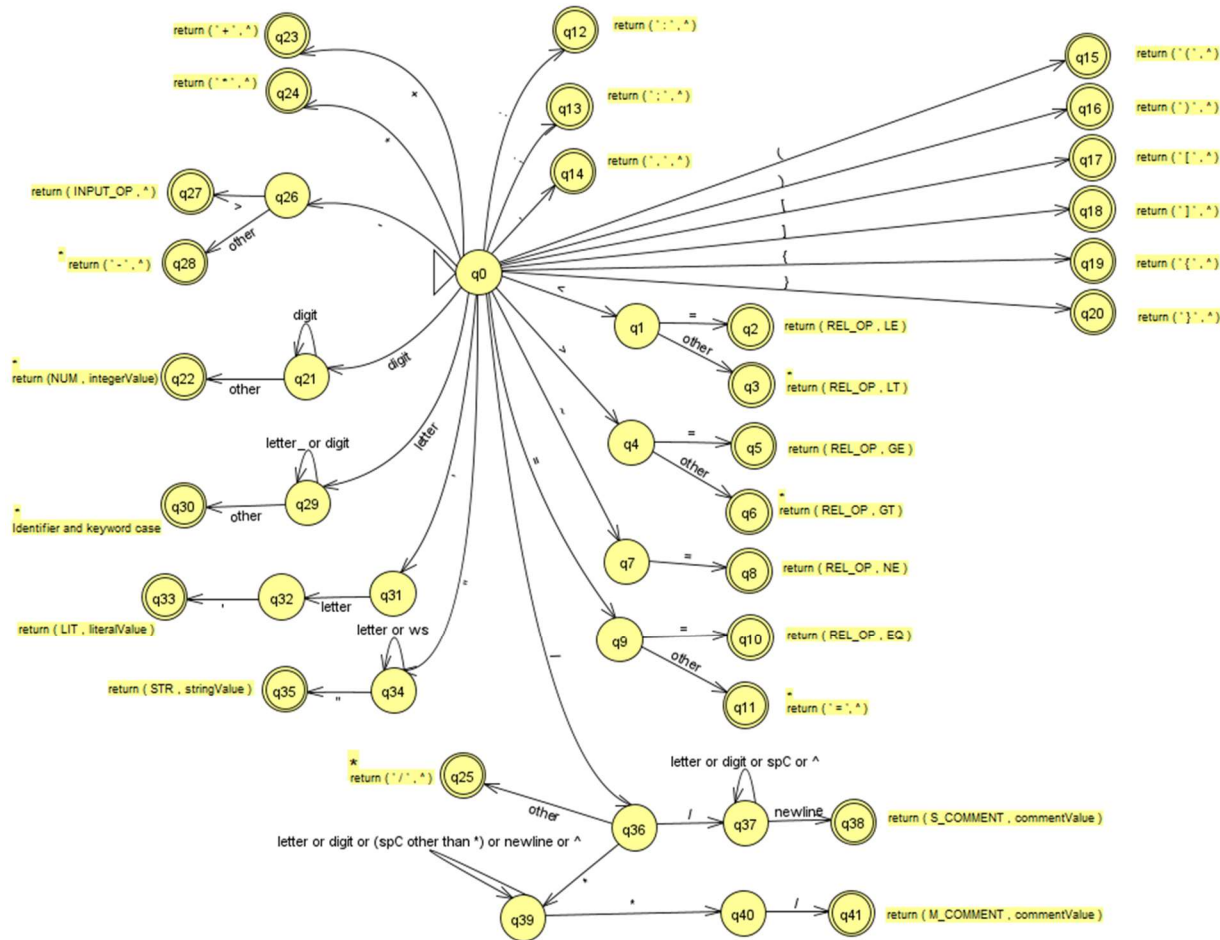
$$PBC \rightarrow ( ) | \{ | \} | [ | ]$$

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

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

### 3. Overall Finite State Machine for Lexer

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.

## 4. CFGs for Parser

Following context free grammars (CFGs) constitute our language which will be used to implement its parser.

### i. Arithmetic Operations

$$Expression \rightarrow Term R$$

$$R \rightarrow + Term R \mid - Term R \mid ^$$

$$Term \rightarrow Factor R'$$

$$R' \rightarrow * Factor R' \mid / Factor R' \mid ^$$

$$Factor \rightarrow ID \mid NUM \mid ( Expression )$$

### ii. Relational Operations

$$Condition \rightarrow Expression relOp Expression$$

### iii. Variable Declaration

$$DT \rightarrow int \mid char$$

$$Variable \rightarrow DT : ID VariableDelimiter$$

$$VariableDelimiter \rightarrow ; \mid , nextVariable$$

$$nextVariable \rightarrow ID VariableDelimiter$$

### iv. Variable Assignment

$$AssignmentStatement \rightarrow ID assignOp Value ;$$

$$Value \rightarrow ID \mid NUM \mid LIT \mid Expression$$

### v. Variable Declaration and Assignment

$$VariableDA \rightarrow intDA \mid charDA$$

$$intDA \rightarrow int : ID intOptionAssign intVariableDelimiterDA$$

$$intVariableDelimiterDA \rightarrow ; \mid , intNextVariableDA$$

$intNextVariableDA \rightarrow ID \ intOptionAssign \ intVariableDelimiterDA$

$intOptionAssign \rightarrow assignOp \ (NUM \mid Expression) \mid ^$

$charDA \rightarrow char : ID \ charOptionAssign \ charVariableDelimiterDA$

$charVariableDelimiterDA \rightarrow ; \mid , \ charNextVariableDA$

$charNextVariableDA \rightarrow ID \ charOptionAssign \ charVariableDelimiterDA$

$charOptionAssign \rightarrow assignOp \ LIT \mid ^$

**vi. Variable Input**

$VariableInput \rightarrow input \ INPUT\_OP \ ID \ inputDelimiter$

$inputDelimiter \rightarrow ; \mid , \ nextInput$

$nextInput \rightarrow ID \ inputDelimiter$

**vii. If-Elif-Else Statement**

$Statement \rightarrow if \ Condition : \{ Statements \} \ ElifOrElse$

$ElifOrElse \rightarrow elif \ Condition : \{ Statements \} \ ElifOrElse \mid Else$

$Else \rightarrow else \{ Statements \} \mid ^$

**viii. Loop Statement**

$Statement \rightarrow while \ Condition : \{ Statements \}$

**ix. Output Statement**

$Statement \rightarrow print \ ( \ OutputOptions \ ) ; \mid println \ ( \ OutputOptions \ ) ;$

$OutputOptions \rightarrow ID \mid NUM \mid LIT \mid STR \mid Expression$

**x. Increment-Decrement Statement**

$Statement \rightarrow ID \ (IncOp \mid DecOp) ;$

$IncOp \rightarrow ++$

$DecOp \rightarrow --$

## 5. Overall Grammar of Parser

In context of the Regular Expressions defined in section 2 and CFGs defined in section 4, the overall grammar of our language will be as following:

$Start \rightarrow Statements$

$Statements \rightarrow Statement\ Statements \mid \wedge$

$Statement \rightarrow \text{if Condition} : \{ Statements \} \text{ ElifOrElse}$

$\mid \text{while Condition} : \{ Statements \}$

$\mid \text{print} ( OutputOptions ) ; \mid \text{println} ( OutputOptions ) ;$

$\mid ID (IncOp \mid DecOp) ;$

$\mid VariableInput$

$\mid Variable$

$\mid AssignmentStatement$

$\mid VariableDA$

$\mid S\_COMMENT \mid M\_COMMENT$

$Expression \rightarrow Term\ R$

$R \rightarrow + Term\ R \mid - Term\ R \mid \wedge$

$Term \rightarrow Factor\ R'$

$R' \rightarrow * Factor\ R' \mid / Factor\ R' \mid \wedge$

$Factor \rightarrow ID \mid NUM \mid ( Expression )$

$Condition \rightarrow Expression\ relOp\ Expression$

$ElifOrElse \rightarrow \text{elif Condition} : \{ Statements \} \text{ ElifOrElse} \mid \text{Else}$

$\text{Else} \rightarrow \text{else} \{ Statements \} \mid \wedge$

$OutputOptions \rightarrow ID \mid NUM \mid LIT \mid STR \mid Expression$

$IncOp \rightarrow ++$

$DecOp \rightarrow --$

$VariableInput \rightarrow \text{input INPUT\_OP ID inputDelimiter}$

$\text{inputDelimiter} \rightarrow ; \mid , \text{nextInput}$

$\text{nextInput} \rightarrow ID\ \text{inputDelimiter}$

$$\begin{aligned}
Variable &\rightarrow DT : ID \ VariableDelimiter \\
VariableDelimiter &\rightarrow ; | , \ nextVariable \\
nextVariable &\rightarrow ID \ VariableDelimiter \\
AssignmentStatement &\rightarrow ID \ assignOp \ Value ; \\
Value &\rightarrow ID | NUM | LIT | Expression \\
VariableDA &\rightarrow intDA | charDA \\
intDA &\rightarrow int : ID \ intOptionAssign \ intVariableDelimiterDA \\
intVariableDelimiterDA &\rightarrow ; | , \ intNextVariableDA \\
intNextVariableDA &\rightarrow ID \ intOptionAssign \ intVariableDelimiterDA \\
intOptionAssign &\rightarrow assignOp \ (NUM | Expression) | ^ \\
charDA &\rightarrow char : ID \ charOptionAssign \ charVariableDelimiterDA \\
charVariableDelimiterDA &\rightarrow ; | , \ charNextVariableDA \\
charNextVariableDA &\rightarrow ID \ charOptionAssign \ charVariableDelimiterDA \\
charOptionAssign &\rightarrow assignOp \ LIT | ^
\end{aligned}$$

**Note:**

Type mismatch is avoided when variable is declared and assigned at the same time i.e., in case of *VariableDA* by using separate trees for int and char types but for simple declaration and later on assignment, symbol table will be referred for information about the identifier.