# Compiler Term Project

| Subject | Compiler |
|---|---|
| Class No. | 01 |
| Professor Name | Hyosu Kim |
| Team No. | 11 |
| Team Members | Hyomin Kim, Sanghwa Lee |
| Date | 2021.04.12 |

# < **Index** >

# A. Definition of Tokens and Their Regular Expressions

## 0) Predefined Tokens

### - letter :

Token name : letter

Regular Expression :

a | b | c | d | e | f | g | h | i | j | k | l | m | n | o | p | q | r | s | t | u | v | w | x | y | z | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y | Z

### - digit :

Token name : digit

Regular Expression : 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

### - positive :

Token name : positive

Regular Expression : 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

## 1) Variable type

Token Name : VTYPE

Regular Expression : int | char | boolean | String

## 2) Signed integer

Token Name : INT

Regular Expression : $0 \mid ((\ -\ \mid\ \epsilon\ )\ positive^{+}\ )$

### 3) Single character

Token Name : CHAR

Regular Expression : 'letter | digit | blank'

### 4) Boolean string

Token Name : BOOL

Regular Expression : true | false

### 5) Literal string

Token Name : STRING

Regular Expression : $"(digits \,|\, letter \,|\, blank)^+"$

### 6) An identifier of variables and functions

Token Name : ID

Regular Expression : ( letters | _ )( digits | letters | _ )*

### 7) Keywords for special statements

1. Token Name : IF

Regular Expression : if

2. Token Name : ELSE

Regular Expression : else

3. Token Name : WHILE

Regular Expression : while

4. Token Name : CLASS

Regular Expression : class

5. Token Name : RETURN

Regular Expression : return

➔ Merged Regular Expression : if | else | while | class | return

**8) Arithmetic operators**

Token Name : OP

Regular Expression : + | - | * | /

**9) Assignment operator**

Token Name : ASSIGN

Regular Expression : =

**10) Comparison operators**

Token Name : COMP

Regular Expression : (( < | > )( = | e )) | (( ! | = ) =)

**11) A terminating symbol of statements :**

Token Name : SEMI

Regular Expression : ;

**12) A pair of symbols for defining area/scope of variables and functions : { and }**

Token Name : BRACE    { : LBRACE                } : RBRACE

Regular Expression : { | }


**13) A pair of symbols for indicating a function/statement : ( and )**

Token Name : PAREN ( : LPAREN                ) : RPAREN

Regular Expression : ( | )


**14) A pair of symbols for using an array : [ and ]**

Token Name : BRACKET [ : LBRACKET           ] : RBRACKET

Regular Expression : [ | ]


**15) A symbol for separating input arguments in functions ; ,**

Token Name : COMMA
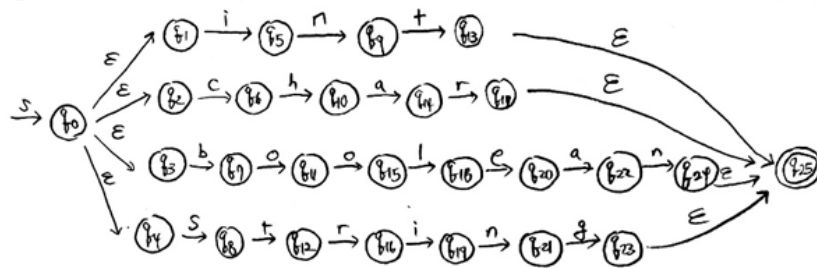
Regular Expression : ,


**16) Whitespaces :**

Token Name : WHITESPACE

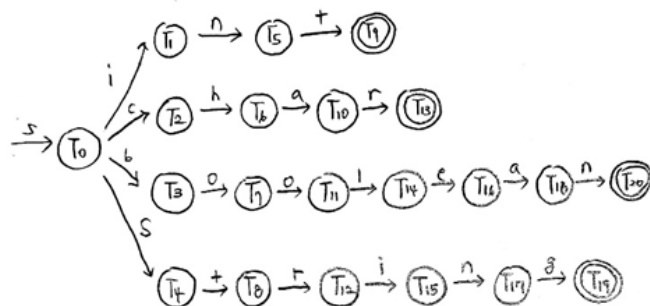Regular Expression : $( \backslash t \,|\, \backslash n \,|\, blank )^{+}$

# B. NFA, DFA transition graph

## 1) Variable type

(TYPE) (NFA)

$\varepsilon$: $q_0 \xrightarrow{\varepsilon} q_1 \xrightarrow{i} q_5 \xrightarrow{n} q_9 \xrightarrow{t} q_{13} \xrightarrow{\varepsilon} q_{25}$

$q_0 \xrightarrow{\varepsilon} q_2 \xrightarrow{c} q_6 \xrightarrow{h} q_{10} \xrightarrow{a} q_{14} \xrightarrow{r} q_{18} \xrightarrow{\varepsilon} q_{25}$

$q_0 \xrightarrow{\varepsilon} q_3 \xrightarrow{b} q_7 \xrightarrow{o} q_{11} \xrightarrow{o} q_{15} \xrightarrow{l} q_{16} \xrightarrow{e} q_{20} \xrightarrow{a} q_{22} \xrightarrow{n} q_{24} \xrightarrow{e} q_{25}$

$q_0 \xrightarrow{\varepsilon} q_4 \xrightarrow{s} q_8 \xrightarrow{t} q_{12} \xrightarrow{r} q_{16} \xrightarrow{i} q_{17} \xrightarrow{n} q_{21} \xrightarrow{g} q_{23} \xrightarrow{\varepsilon} q_{25}$

(TYPE) (DFA).

$T_0 \xrightarrow{i} T_1 \xrightarrow{n} T_5 \xrightarrow{t} T_9$

$T_0 \xrightarrow{c} T_2 \xrightarrow{h} T_6 \xrightarrow{a} T_{10} \xrightarrow{r} T_{13}$

$T_0 \xrightarrow{b} T_3 \xrightarrow{o} T_7 \xrightarrow{o} T_{11} \xrightarrow{l} T_{14} \xrightarrow{e} T_{16} \xrightarrow{a} T_{18} \xrightarrow{n} T_{20}$

$T_0 \xrightarrow{s} T_4 \xrightarrow{t} T_8 \xrightarrow{r} T_{12} \xrightarrow{i} T_{15} \xrightarrow{n} T_{17} \xrightarrow{g} T_{19}$

## 2) Signed integer

(int) (NFA)



(int) (DFA)

$T_0 = e\text{-closure} (8_0) = \{ 8_0, 8_1, 8_3, 8_4, 8_5, 8_7, 8_8, 8_9 \}$

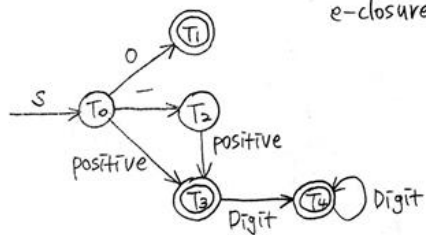$T_1 = e\text{-closure} ( \delta(T_0, 0)) = \{ 8_2, 8_{14} \}$

$T_2 = e\text{-closure} ( \delta(T_0, -)) = \{ 8_6, 8_8, 8_9 \}$

$T_3 = e\text{-closure} ( \delta(T_0, positive)) = \{ 8_{10}, 8_{11}, 8_{13}, 8_{14} \}$

$\quad e\text{-closure} ( \delta(T_2, positive)) = \{ 8_{10}, 8_{11}, 8_{13}, 8_{14} \} = T_3$

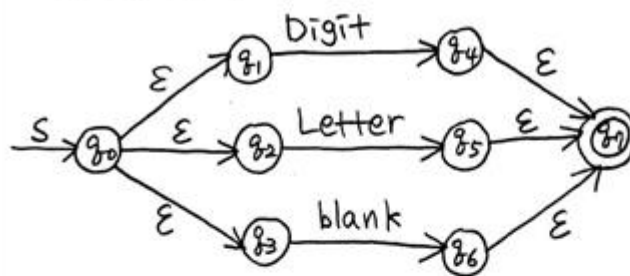$T_4 = e\text{-closure} ( \delta(T_3, Digit)) = \{ 8_{11}, 8_{12}, 8_{13}, 8_{14} \}$

$\qquad\qquad e\text{-closure} ( \delta(T_4, Digit)) = \{ 8_{11}, 8_{12}, 8_{13}, 8_{14} \} = T_4$

## 3) Single character

(char) (NFA)



(char) (DFA)
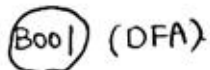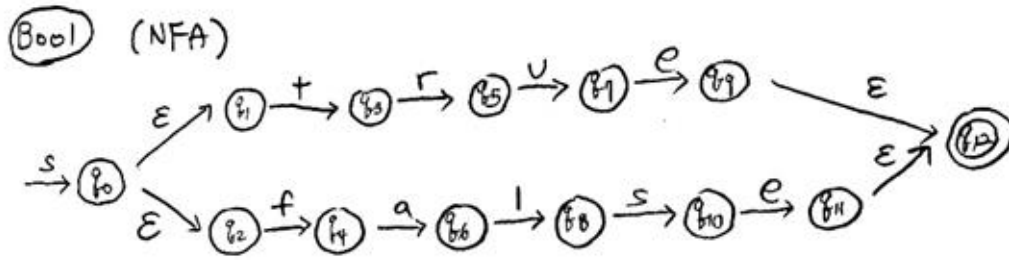
$T_0 = $ e-closure $(g_0) = \{ g_0, g_1, g_2, g_3 \}$

$T_1 = $ e-closure $(\delta(T_0, Digit)) = \{ g_4, g_7 \}$

$T_2 = $ e-closure $(\delta(T_0, Letter)) = \{ g_5, g_7 \}$

$T_3 = $ e-closure $(\delta(T_0, blank)) = \{ g_6, g_7 \}$

## 4) Boolean string

Bool (NFA)



Bool (DFA)

$T_0 = e\text{-}closure\,(g_0) = \{g_1, g_2\}$

$T_1 = e\text{-}closure\,(\delta(T_0, +)) = \{g_3\}$

$T_2 = e\text{-}closure\,(\delta(T_0, f)) = \{g_4\}$

$T_3 = e\text{-}closure\,(\delta(T_1, r)) = \{g_5\}$

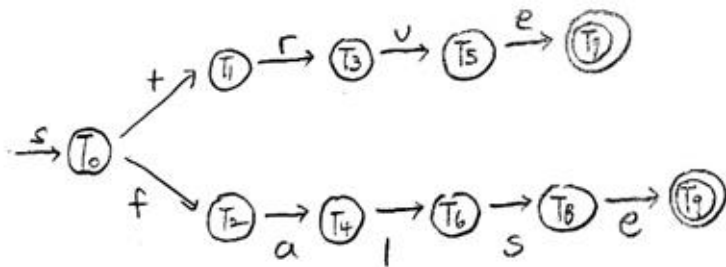$T_4 = e\text{-}closure\,(\delta(T_2, a)) = \{g_6\}$

$T_5 = e\text{-}closure\,(\delta(T_3, U)) = \{g_7\}$

$T_6 = e\text{-}closure\,(\delta(T_4, 1)) = \{g_8\}$
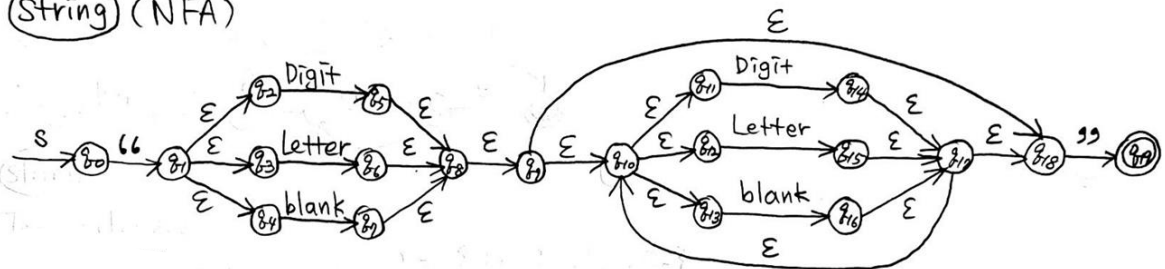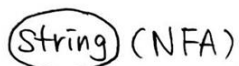
$T_7 = e\text{-}closure\,(\delta(T_5, e)) = \{g_9, g_{12}\}$

$T_8 = e\text{-}closure\,(\delta(T_6, S)) = \{g_{10}\}$

$T_9 = e\text{-}closure\,(\delta(T_8, e)) = \{g_{11}, g_{12}\}$



## 5) Literal String

String (NFA)

$\boxed{String}$ (DFA)

$T_0 = e\text{-closure}(q_0) = \{q_0\}$

$T_1 = e\text{-closure}(\delta(T_0, \text{``})) = \{q_1, q_2, q_3, q_4\}$

$T_2 = e\text{-closure}(\delta(T_1, Digit)) = \{q_5, q_8, q_9, q_{10}, q_{11}, q_{12}, q_{13}, q_{18}\}$

$T_3 = e\text{-closure}(\delta(T_1, Letter)) = \{q_6, q_8, q_9, q_{10}, q_{11}, q_{12}, q_{13}, q_{18}\}$

$T_4 = e\text{-closure}(\delta(T_1, blank)) = \{q_7, q_8, q_9, q_{10}, q_{11}, q_{12}, q_{13}, q_{18}\}$

$T_5 = e\text{-closure}(\delta(T_2, Digit)) = e\text{-closure}(\delta(T_3, Digit)) = e\text{-closure}(\delta(T_4, Digit))$
$\quad = \{q_{10}, q_{11}, q_{12}, q_{13}, q_{14}, q_{17}, q_{18}\}$

$T_6 = e\text{-closure}(\delta(T_2, Letter)) = e\text{-closure}(\delta(T_3, Letter)) = e\text{-closure}(\delta(T_4, Letter))$
$\quad = \{q_{10}, q_{11}, q_{12}, q_{13}, q_{15}, q_{17}, q_{18}\}$

$T_7 = e\text{-closure}(\delta(T_2, blank)) = e\text{-closure}(\delta(T_3, blank)) = e\text{-closure}(\delta(T_4, blank))$
$\quad = \{q_{10}, q_{11}, q_{12}, q_{13}, q_{16}, q_{17}, q_{18}\}$

$T_8 = e\text{-closure}(\delta(T_2, \text{''})) = e\text{-closure}(\delta(T_3, \text{''})) = e\text{-closure}(\delta(T_4, \text{''})) = \{q_{19}\}$

$\quad e\text{-closure}(\delta(T_5, Digit)) = e\text{-closure}(\delta(T_6, Digit)) = e\text{-closure}(\delta(T_7, Digit)) = T_5$

$\quad e\text{-closure}(\delta(T_5, Letter)) = e\text{-closure}(\delta(T_6, Letter)) = e\text{-closure}(\delta(T_7, Letter)) = T_6$

$\quad e\text{-closure}(\delta(T_5, blank)) = e\text{-closure}(\delta(T_6, blank)) = e\text{-closure}(\delta(T_7, blank)) = T_7$
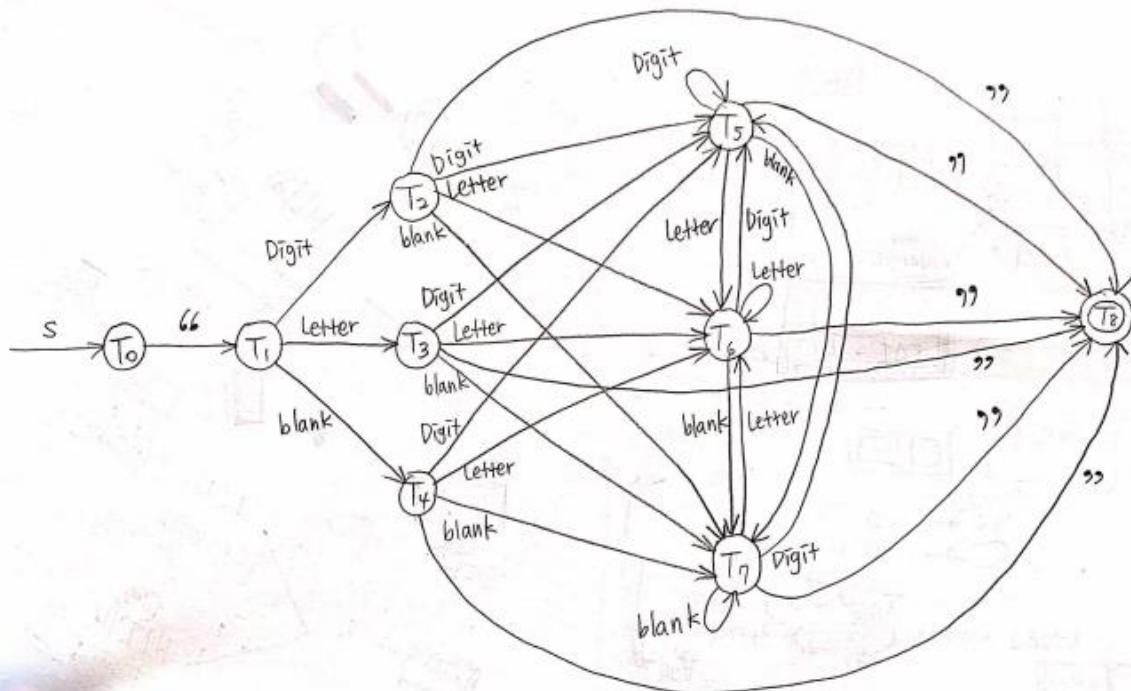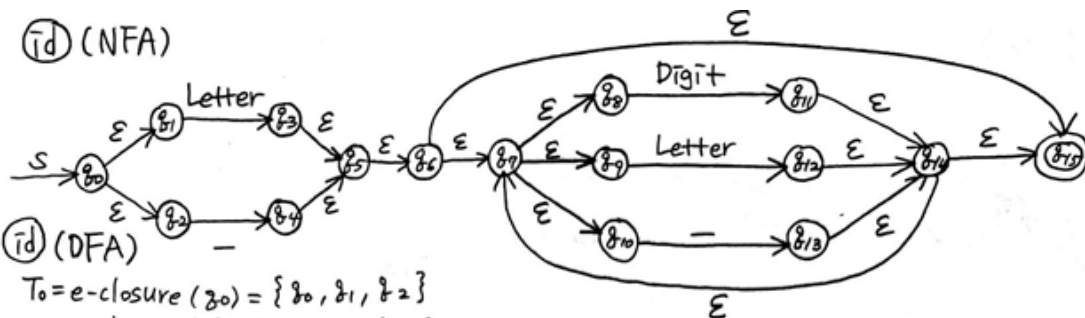
$\quad e\text{-closure}(\delta(T_5, \text{''})) = e\text{-closure}(\delta(T_6, \text{''})) = e\text{-closure}(\delta(T_7, \text{''})) = \{q_{19}\} = T_8$

## 6) An identifier of variables and functions

**(id) (NFA)**



**(id) (DFA)**

$T_0 = $ e-closure $(z_0) = \{ z_0, z_1, z_2 \}$

$T_1 = $ e-closure $(\delta(T_0, \text{Letter})) = \{ z_3, z_5, z_6, z_7, z_8, z_9, z_{10}, z_{15} \}$

$T_2 = $ e-closure $(\delta(T_0, \_)) = \{ z_4, z_5, z_6, z_7, z_8, z_9, z_{10}, z_{15} \}$

$T_3 = $ e-closure $(\delta(T_1, \text{Digit})) = \{ z_7, z_8, z_9, z_{10}, z_{11}, z_{14}, z_{15} \}$

$T_4 = $ e-closure $(\delta(T_1, \text{Letter})) = \{ z_7, z_8, z_9, z_{10}, z_{12}, z_{14}, z_{15} \}$

$T_5 = $ e-closure $(\delta(T_1, \_)) = \{ z_7, z_8, z_9, z_{10}, z_{13}, z_{14}, z_{15} \}$

  e-closure $(\delta(T_2, \text{Digit})) = T_3$
  e-closure $(\delta(T_2, \text{Letter})) = T_4$
  e-closure $(\delta(T_2, \_)) = T_5$
  e-closure $(\delta(T_3, \text{Digit})) = T_3$
  e-closure $(\delta(T_3, \text{Letter})) = T_4$
  e-closure $(\delta(T_3, \_)) = T_5$
  e-closure $(\delta(T_4, \text{Digit})) = T_3$
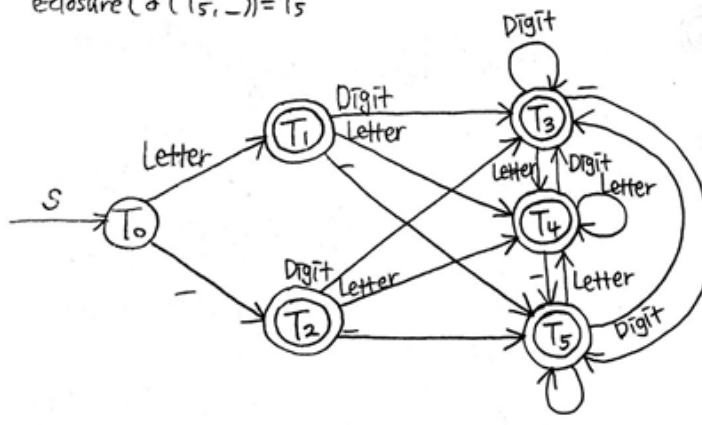  e-closure $(\delta(T_4, \text{Letter})) = T_4$
  e-closure $(\delta(T_4, \_)) = T_5$
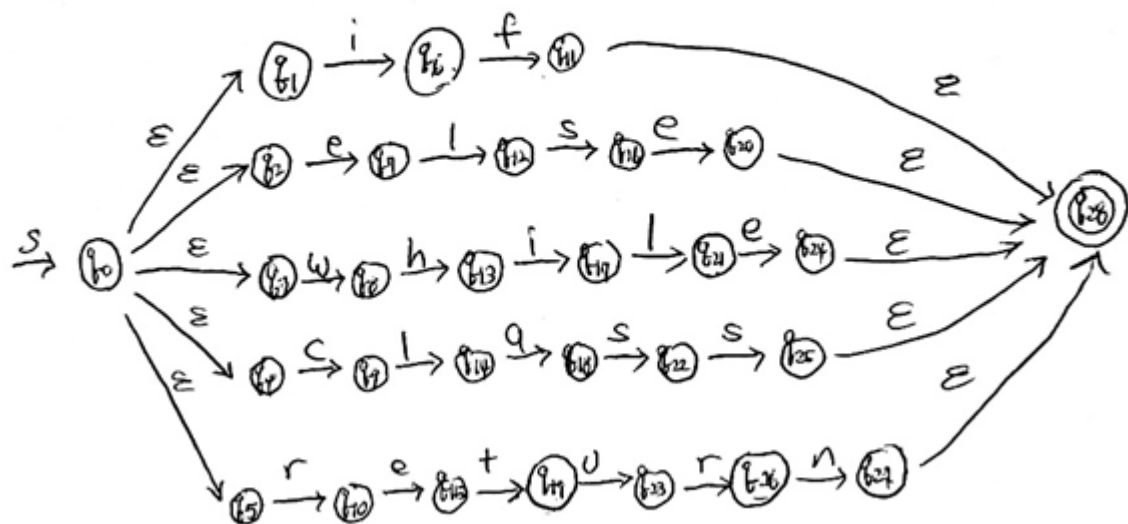  e-closure $(\delta(T_5, \text{Digit})) = T_3$
  e-closure $(\delta(T_5, \text{Letter})) = T_4$
  e-closure $(\delta(T_5, \_)) = T_5$

## 7) Keywords for special statements

(Keyword) (NFA)

(keyword) (DFA)

$T_0 = \text{e-closure}(q_0) = \{q_0, q_1, q_2, q_3, q_4, q_5\}$

$T_1 = \text{e-closure}(\delta(T_0, i)) = \{q_6\}$

$T_2 = \text{e-closure}(\delta(T_0, e)) = \{q_7\}$

$T_3 = \text{e-closure}(\delta(T_0, w)) = \{q_8\}$

$T_4 = \text{e-closure}(\delta(T_0, c)) = \{q_9\}$

$T_5 = \text{e-closure}(\delta(T_0, r)) = \{q_{10}\}$

$T_6 = \text{e-closure}(\delta(T_1, f)) = \{q_{11}, q_{28}\}$

$T_7 = \text{e-closure}(\delta(T_2, l)) = \{q_{12}\}$

$T_8 = \text{e-closure}(\delta(T_3, h)) = \{q_{13}\}$

$T_9 = \text{e-closure}(\delta(T_4, l)) = \{q_{14}\}$

$T_{10} = \text{e-closure}(\delta(T_5, e)) = \{q_{15}\}$

$T_{11} = \text{e-closure}(\delta(T_7, s)) = \{q_{16}\}$

$T_{12} = \text{e-closure}(\delta(T_8, i)) = \{q_{17}\}$

$T_{13} = \text{e-closure}(\delta(T_9, a)) = \{q_{18}\}$

$T_{14} = \text{e-closure}(\delta(T_{10}, t)) = \{q_{19}\}$

$T_{15} = \text{e-closure}(\delta(T_{11}, e)) = \{q_{20}, q_{28}\}$

$T_{16} = \text{e-closure}(\delta(T_{12}, l)) = \{q_{21}\}$

$T_{17} = \text{e-closure}(\delta(T_{13}, s)) = \{q_{22}\}$

$T_{18} = \text{e-closure}(\delta(T_{14}, u)) = \{q_{23}\}$

$T_{19} = \text{e-closure}(\delta(T_{16}, e)) = \{q_{24}, q_{28}\}$

$T_{20} = \text{e-closure}(\delta(T_{17}, s)) = \{q_{25}, q_{28}\}$

$T_{21} = \text{e-closure}(\delta(T_{18}, r)) = \{q_{26}\}$

$T_{22} = \text{e-closure}(\delta(T_{21}, n)) = \{q_{27}, q_{28}\}$

```
                    T1 --f--> (T6)

             i
          -------> T2 --l--> T7 --s--> T11 --e--> (T15)
          e
   s   -------
 --> T0 ---------> T3 --h--> T8 --i--> T12 --l--> T16 --e--> (T19)
          w
          c
          ------->  T4 --l--> T9 --a--> T13 --s--> T17 --s--> (T20)
          r
          ------->  T5 --e--> T10 --t--> T14 --u--> T18 --r--> T21 --n--> (T22)
```

**8) Arithmetic operators**

Arith NFA



Arith DFA

$$T_0 = \text{e-closure}(q_0) = \{q_0, q_1, q_2, q_3, q_4\}$$

$$T_1 = \text{e-closure}(\delta(T_0, +)) = \{q_5, q_9\}$$

$$T_2 = \text{e-closure}(\delta(T_0, -)) = \{q_6, q_9\}$$

$$T_3 = \text{e-closure}(\delta(T_0, *)) = \{q_7, q_9\}$$

$$T_4 = \text{e-closure}(\delta(T_0, /)) = \{q_8, q_9\}$$

## 9) Assignment operator (NFA and DFA is equal)

Assign) NFA , DFA

$S \rightarrow (q_0) \xrightarrow{=} (q_1)$

## 10) Comparison operators

Comp (NFA)



Comp (DFA)

$T_0 = \text{e-closure}(q_0) = \{ q_0, q_1, q_2, q_3, q_4, q_5, q_6 \}$

$T_1 = \text{e-closure}(\delta(T_0, <)) = \{ q_7, q_{11}, q_{13}, q_{14}, q_{17}, q_{18}, q_{19} \}$

$T_2 = \text{e-closure}(\delta(T_0, >)) = \{ q_8, q_{11}, q_{13}, q_{14}, q_{17}, q_{18}, q_{19} \}$

$T_3 = \text{e-closure}(\delta(T_0, !)) = \{ q_9, q_{12} \}$

$T_4 = \text{e-closure}(\delta(T_0, =)) = \{ q_{10}, q_{12} \}$

$T_5 = \text{e-closure}(\delta(T_1, =)) = \{ q_{16}, q_{18}, q_{19} \}$

$\quad \text{e-closure}(\delta(T_2, =)) = \{ q_{16}, q_{18}, q_{19} \} = T_5$

$T_6 = \text{e-closure}(\delta(T_3, =)) = \{ q_{15}, q_{19} \}$

$\quad \text{e-closure}(\delta(T_4, =)) = \{ q_{15}, q_{19} \} = T_6$

**11) A terminating symbol of statements : ; (NFA and DFA is equal)**

(Semi) NFA , DFA

$\xrightarrow{;}$ ($q_0$) $\xrightarrow{j}$ (($q_1$))

**12) A pair of symbols for defining area/scope of variables and functions : { and }**

(Braces) NFA

$\xrightarrow{S}$ ($q_0$)

$\xrightarrow{\varepsilon}$ ($q_1$) $\xrightarrow{\{}$ ($q_3$) $\xrightarrow{\varepsilon}$ (($q_5$)) LBRACE

$\xrightarrow{\varepsilon}$ ($q_2$) $\xrightarrow{\}}$ ($q_4$) $\xrightarrow{\varepsilon}$ (($q_6$)) RBRACE

(Braces) DFA

$T_0$ = e-closure ($q_0$) = $\{q_1, q_2\}$

$T_1$ = e-closure ($\delta(T_0, \{)$) = $\{q_3, q_5\}$

$T_2$ = e-closure ($\delta(T_0, \})$) = $\{q_4, q_6\}$

$\xrightarrow{S}$ ($T_0$)

$\xrightarrow{\{}$ (($T_1$)) LBRACE

$\xrightarrow{\}}$ (($T_2$)) RBRACE

**13) A pair of symbols for indicating a function / statement : ( and )**

PAREN NFA



PAREN DFA

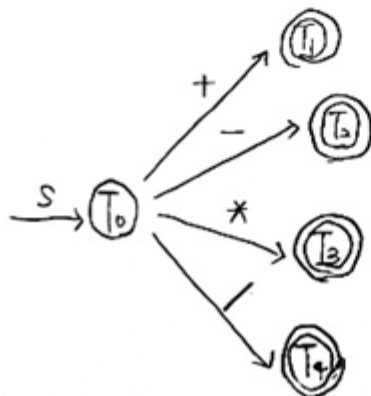$T_0 = \text{e-closure} (q_0) = \{q_1, q_2\}$

$T_1 = \text{e-closure} (\delta(T_0, ()) = \{q_3, q_5\}$
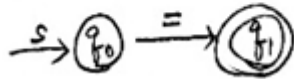
$T_2 = \text{e-closure} (\delta(T_0, ))) = \{q_4, q_6\}$

## 14) A pair of symbols for using an array : [ and ]

**BRACKET ) NFA**

$$S \to q_0$$
$$\varepsilon \to q_1 \xrightarrow{[} q_3 \xrightarrow{\varepsilon} q_5 \quad \text{LBRACKET}$$
$$\varepsilon \to q_2 \xrightarrow{]} q_4 \xrightarrow{\varepsilon} q_6 \quad \text{RBRACKET}$$

**BRACKET DFA**

$$T_0 = e\text{-closure}(q_0) = \{q_1, q_2\}$$
$$T_1 = e\text{-closure}(\delta(T_0, [\,)) = \{q_3, q_5\}$$
$$T_2 = e\text{-closure}(\delta(T_0, ]\,)) = \{q_4, q_6\}$$

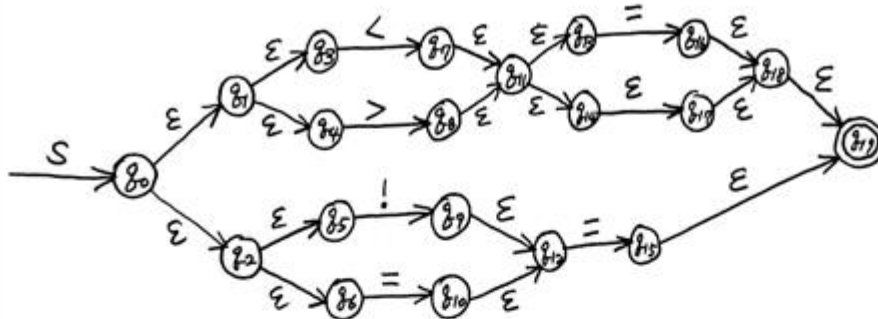$$S \to T_0 \xrightarrow{[} T_1 \quad \text{LBRACKET}$$
$$T_0 \xrightarrow{]} T_2 \quad \text{RBRACKET}$$

## 15) A symbol for separating input arguments in functions : , (NFA and DFA is equal)

**Comma** NFA , DFA

$$S \to q_0 \xrightarrow{,} q_1$$

## 16) Whitespaces

(Space) **NFA**



(Space) **DFA**

$$T_0 = \text{e-closure}(q_0) = \{q_0, q_1, q_2, q_3\}$$

$$T_1 = \text{e-closure}(\delta(T_0, \backslash t)) = \{q_4, q_7\}$$

$$T_2 = \text{e-closure}(\delta(T_1, \backslash n)) = \{q_5, q_7\}$$

$$T_3 = \text{e-closure}(\delta(T_2,\ )) = \{q_6, q_7\}$$

# C. All About Our Lexical Analyzer

## 1) Overall Flow

```
                    ┌─────────────────┐
                    │ Input Text File │
                    └─────────────────┘
                             │
     ┌───────────────────────┼──────────────────────────────┐
     │ lexical_analyzer.py   │                              │
     │                       ▼                              │         lexemes
     │     ┌──────────────────────────────────┐             │
     │     │ Divide raw input string to lexemes│            │
     │     └──────────────────────────────────┘             │            ┌──────────────────────┐
     │                       │                              │            │ DFA.py               │
     │     ┌──────────────────────────────────┐             │            │         ▼            │
     │     │ Put divided lexemes into each DFA │─────────────┼───────────▶│ ┌──────────────────┐ │
     │     └──────────────────────────────────┘             │            │ │ Defined DFAs     │ │
     │                       │                              │  return    │ │ Information      │ │
     │     ┌──────────────────────────────────┐             │  state info└─┤                  │ │
     │     │ Put divided lexemes into each DFA │◀────────────┼──────────────┘                  │ │
     │     └──────────────────────────────────┘             │            └──────────────────────┘
     │     ┌──────────────────────────────────┐             │
     │     │ If DFA matches with lexeme,       │             │
     │     │ name token with DFA name          │             │
     │     └──────────────────────────────────┘             │
     │     ┌──────────────────────────────────┐             │
     │     │ Make output file                  │             │
     │     └──────────────────────────────────┘             │
     └───────────────────────┼──────────────────────────────┘
                             │
                    ┌─────────────────┐
                    │ Output Text File│
                    └─────────────────┘
```

**2) About Each File**

**1. DFA.py**

```
1   class Transition:
2       def __init__(self, current_state, input_symbol, next_state):
3           self.current_state = current_state
4           self.input_symbol = input_symbol
5           self.next_state = next_state
6
```

Define transition function class. If input_symbol is entered current_state moves to next_state.

```
class DFA:
    def __init__(self, name, start_state, trans_functions, final_state):
        self.name = name
        self.start_state = start_state
        self.trans_functions = trans_functions
        self.final_state = final_state
        self.token = None
```

DFA class initialization. It consists name, start_state, trans_functions(transition functions of each DFA), final_state. Token is initialized as "None" at first.

```
def is_not_error(self, s):        # DFA에 넣어서 error가 나오면 False
    current = [0]
    for c in s:
        destination = []
        f = 0
        for trans in self.trans_functions:
            if (c in trans.input_symbol) and (trans.current_state in current):
                destination.append(trans.next_state)
                f = 1
        if f == 0:
            return False
        if len(destination):
            current = destination
        else:
            return False
    return current
```

Function to check whether current input string 's' can be accepted at current DFA or not. Input string s is a part of "input.txt" string data. Variable 'f' is to check whether current input is accepted to current DFA. If f is 0, it returns false. If there is any input at destination list, it returns current which is a list of destination.

```python
def is_accept(self, s):          # DFA의 final까지 도달하는지 검사
    temp = self.is_not_error(s)
    if type(temp) is list:
        return self.is_final(temp)
    else:
        return False


def is_final(self, temp):         # final state에 있는지 검사
    for i in temp:
        if i in self.final_state:
            self.token = i
            return True
    return False
```

Function to check if input 's' reaches to current DFA's final state. If temp is list it goes to is_final function and check if temp list has final_state of current DFA.

```python
# define DFA
# DFA's final state is candidate token name of the string
digits = '0123456789'
positive = '123456789'
letter = 'ABCDEFGHIJKLMNOPQRSTUVWXYZabcdefghijklmnopqrstuvwxyz'
TYPE = DFA('TYPE', 0, [Transition(0, 'i', 1), Transition(1, 'n', 5), Transition(5, 't', 'VTYPE'),
                       Transition(0, 'c', 2), Transition(2, 'h', 6), Transition(6, 'a', 10), Transition(10, 'r', 'VTYPE'),
                       Transition(0, 'b', 3), Transition(3, 'o', 7), Transition(7, 'o', 11), Transition(11, 'l', 14),
                       Transition(14, 'e', 16), Transition(16, 'a', 18), Transition(18, 'n', 'VTYPE'),
                       Transition(0, 's', 4), Transition(4, 't', 8), Transition(8, 'r', 12),
                       Transition(12, 'i', 15), Transition(15, 'n', 17), Transition(17, 'g', 'VTYPE')], ['VTYPE'])
ZERO = DFA('ZERO', 0, [Transition(0, '0', 'INTEGER')], ['INTEGER'])
INT = DFA('INT', 0, [Transition(0, '-', 1), Transition(0, positive, 'INTEGER'),
                     Transition(1, positive, 'INTEGER'), Transition('INTEGER', digits, 'INTEGER')], ['INTEGER'])
CHAR = DFA('CHAR', 0, [Transition(0, "'", 1), Transition(1, digits, 2), Transition(1, letter, 3), Transition(1, ' ', 4),
                       Transition(2, "'", 'CHAR'), Transition(3, "'", 'CHAR'), Transition(4, "'", 'CHAR')], ['CHAR'])
BOOL = DFA('BOOL', 0, [Transition(0, 't', 1), Transition(1, 'r', 3), Transition(3, 'u', 5), Transition(5, 'e', 'BOOL'),
                       Transition(0, 'f', 2), Transition(2, 'a', 4), Transition(4, 'l', 6),
                       Transition(6, 's', 8), Transition(8, 'e', 'BOOL')], ['BOOL'])
STRING = DFA('STRING', 0, [Transition(0, '"', 1), Transition(1, digits, 2), Transition(1, letter, 3), Transition(1, ' ', 4),
                           Transition(2, digits, 5), Transition(2, letter, 6), Transition(2, ' ', 7),
                           Transition(3, digits, 5), Transition(3, letter, 6), Transition(3, ' ', 7),
                           Transition(4, digits, 5), Transition(4, letter, 6), Transition(4, ' ', 7),
                           Transition(2, '"', 'STRING'), Transition(3, '"', 'STRING'), Transition(4, '"', 'STRING'),
                           Transition(5, digits, 5), Transition(5, letter, 6), Transition(5, ' ', 7),
                           Transition(6, digits, 5), Transition(6, letter, 6), Transition(6, ' ', 7),
                           Transition(7, digits, 5), Transition(7, letter, 6), Transition(7, ' ', 7),
                           Transition(5, '"', 'STRING'), Transition(6, '"', 'STRING'), Transition(7, '"', 'STRING')], ['STRING'])
ID = DFA('ID', 0, [Transition(0, letter, 'ID'), Transition(0, '_', 'ID'),
                   Transition('ID', digits, 'ID'), Transition('ID', letter, 'ID'), Transition('ID', '_', 'ID')], ['ID'])
KEY = DFA('KEY', 0, [Transition(0, 'i', 1), Transition(1, 'f', 'IF'), Transition(0, 'e', 2), Transition(2, 'l', 7),
                     Transition(7, 's', 11), Transition(11, 'e', 'ELSE'), Transition(0, 'w', 3), Transition(3, 'h', 8),
                     Transition(8, 'i', 12), Transition(12, 'l', 16), Transition(16, 'e', 'WHILE'), Transition(0, 'c', 4),
                     Transition(4, 'l', 9), Transition(9, 'a', 13), Transition(13, 's', 17), Transition(17, 's', 'CLASS'),
                     Transition(0, 'r', 5), Transition(5, 'e', 10), Transition(10, 't', 14), Transition(14, 'u', 18),
                     Transition(18, 'r', 21), Transition(21, 'n', 'RETURN')], ['IF', 'ELSE', 'WHILE', 'CLASS', 'RETURN'])
ARITH = DFA('ARITH', 0, [Transition(0, '+', 'OP'), Transition(0, '-', 'OP'), Transition(0, '*', 'OP'),
                         Transition(0, '/', 'OP')], ['OP'])
ASSIGN = DFA('ASSIGN', 0, [Transition(0, '=', 'ASSIGN')], ['ASSIGN'])
COMP = DFA('COMP', 0, [Transition(0, '<', 'COMP'), Transition(0, '>', 'COMP'), Transition(0, '!', 3), Transition(0, '=', 4),
                       Transition('COMP', '=', 'COMP'), Transition(3, '=', 'COMP'), Transition(4, '=', 'COMP')], ['COMP'])
SEMI = DFA('SEMI', 0, [Transition(0, ';', 'SEMI')], ['SEMI'])
BRACKET = DFA('BRACKET', 0, [Transition(0, '[', 'LBRACKET'), Transition(0, ']', 'RBRACKET')], ['LBRACKET', 'RBRACKET'])
PAREN = DFA('PAREN', 0, [Transition(0, '(', 'LPAREN'), Transition(0, ')', 'RPAREN')], ['LPAREN', 'RPAREN'])
BRACE = DFA('BRACE', 0, [Transition(0, '{', 'LBRACE'), Transition(0, '}', 'RBRACE')], ['LBRACE', 'RBRACE'])
COMMA = DFA('COMMA', 0, [Transition(0, ',', 'COMMA')], ['COMMA'])
WHITESPACE = DFA('WHITESPACE', 0, [Transition(0, '\t', 'WHITESPACE'), Transition(0, '\n', 'WHITESPACE'),
                                   Transition(0, ' ', 'WHITESPACE')], ['WHITESPACE'])
```

Define each DFA.

| DFA | Available Inputs | Final States ( Token Name ) |
|---|---|---|
| TYPE | i n t c h a r b o l e s g | VTYPE |
| ZERO | 0 | INTEGER |
| INT | 0 - positive digits | INTEGER |
| CHAR | ' '  letter digits blank | CHAR |
| BOOL | t r u e f a l s | BOOL |
| STRING | " " letter digits blank | STRING |
| ID | letter _ digits | ID |
| KEY | i f e l s w h i c a r t u n | IF, ELSE, WHILE, CLASS, RETURN |
| ARITH | + - * / | OP |
| ASSIGN | = | ASSIGN |
| COMP | < > = ! | COMP |
| SEMI | ; | SEMI |
| BRACKET | [ ] | LBRACKET, RBRACKET |
| PAREN | ( ) | LPAREN, RPAREN |
| BRACE | { } | LBRACE, RBRACE |
| COMMA | , | COMMA |
| WHITESPACE | ₩t, ₩n, blank | WHITESPACE |

This table is information of our DFAs. (We added 'ZERO' DFA to deal with 0 integers.)

```
INT = DFA('INT', 0, [Transition(0, '-', 1), Transition(0, positive, 'INTEGER'),
                     Transition(1, positive, 'INTEGER'), Transition('INTEGER', digits, 'INTEGER')], ['INTEGER'])
```

Let's see this(↑) example.

'INT' : name of DFA, which will be used as token name

0 : start state

[ Transition (current_state, input_symbol, next_state) ] : List of transition functions

['INT'] : list of DFA's final states

## 2. lexical_analyzer.py

```
import argparse
from DFA import *

DFAs = [TYPE, ZERO, INT, CHAR, BOOL, STRING, ID, KEY, ARITH, ASSIGN, COMP, SEMI, BRACKET, BRACE, PAREN, COMMA, WHITESPACE]
```

Import argparse : for command argument

from DFA import * : import DFA.py

DFAs = [ TYPE, INT, ..... WHITESPACE ] : list of DFAs

```
def determineDFA(candidateDFAs, finalDFAs):
    if len(finalDFAs) == 2 and ID in finalDFAs:
        finalDFAs.remove(ID)
        return list(finalDFAs)
    elif len(candidateDFAs) == 0 and len(finalDFAs) == 1:
        return list(finalDFAs)
    elif len(candidateDFAs) == 0 and len(finalDFAs) == 0:
        return False
```

We use two list to determine lexeme's DFA.

First Condition : If input such as "if", "true" are entered, "KEY" and "ID" DFA or "BOOL" and "ID" DFAs are remained at finalDFAs. However in that case, we need to determine them as "KEY", or "BOOL" instead of "ID". So we removed ID DFA from DFAs list and return DFAs list which only contains "KEY" or "BOOL".

Second Condition : if candidateDFAs list is empty and finalDFAs has one composition, we use DFA in finalDFAs list.

Third Condition : if both candidateDFAs and finalDFAs are empty we don't it returns false.

```python
def lexicalAnalysis(rawString, input_file_name):          # main algorithm of lexical analyzer
    fw = open(input_file_name + "_output.txt", 'w')       # output file
    word_start, word_final = 0, 1                         # index of the string for lexical analysis
    lexemes = []                                          # list of lexeme(string)
    tokens = []                                           # list of token(string : dfa.token)

    # set candidate DFA list and final DFA
    candidateDFAs = [TYPE, ZERO, INT, CHAR, BOOL, STRING, ID, KEY, ARITH, ASSIGN, COMP, SEMI, BRACKET, BRACE, PAREN, COMMA, WHITESPACE]
    finalDFAs = set()
```

Function for lexical analyzing. First create output file to write output information. "word_start" and "word_final" are index for string. We use "token" string variable, "lexemes" list to store lexemes and "tokens" list to store tokens.

```python
while word_final <= len(rawString):
    for dfa in candidateDFAs[:]:                          # filter candidate DFA to check error in DFA
        if not dfa.is_not_error(rawString[word_start:word_final]):
            candidateDFAs.remove(dfa)
    if not candidateDFAs:
        correctDFA = determineDFA(candidateDFAs, finalDFAs)    # determine string's correct DFA
        if not correctDFA:                               # print input code has error
            fw.write('ERROR')
            return 0
        lexeme = rawString[word_start:word_final - 1]    # define lexeme
        lexemes.append(lexeme)                           # append lexeme to list
        token = correctDFA[0].token                      # define token name
        tokens.append(token)                             # append token to list
        word_start = word_final - 1                      # reset index of string and candidate, final DFA
        word_final -= 1
        candidateDFAs = [TYPE, ZERO, INT, CHAR, BOOL, STRING, ID, KEY, ARITH, ASSIGN, COMP, SEMI, BRACKET, BRACE, PAREN,
                         COMMA, WHITESPACE]
        finalDFAs = set()
    else:
        finalDFAs = set()
        for dfa in DFAs[:]:                              # filter final DFA to check if string is accepted in DFA
            if dfa.is_accept(rawString[word_start:word_final]):
                finalDFAs.add(dfa)
```

Until "word_final" index variable reaches the end of the string following codes are repeatedly executed. First we check DFAs in "candidateDFAs" and remove DFAs which don't accept current input.

Then if "candidateDFAs" is not empty we check both "candidateDFAs" and "finalDFAs" to determine "correctDFA" through "determineDFA" function. "correctDFA" variable will be used as final DFA for current string data.

After "determineDFA", if there aren't any "correctDFA" it will return "ERROR". However if "correctDFA" exists we add current string as "lexeme" and add it to "lexemes" list. Then we use "correctDFA"'s DFA name as token and add it to "tokens" list.

Then we initialize "word_start" and "word_final" index variable to the next string index. We also initialize "candidateDFAs" list and "finalDFAs" set.

However, if "candidateDFAs" is empty we use "finalDFAs" and check each DFA in "DFAs" list. If current string reaches to the final state of current DFA it will be added to "finalDFAs" list.

```python
                    finalDFAs.add(dfa)
    # exception handling : - can be OP or negative sign
    # - is OP when previous token is INT or ID, and negative sign in other situations
    # also the immediately token of - can be blank
    if rawString[word_start:word_final] == '-':
        if lexemes[-1] == ' ':
            if tokens[-2] == 'ID' or tokens[-2] == 'INTEGER':
                candidateDFAs.remove(INT)
            else:
                candidateDFAs.remove(ARITH)
        else:
            if tokens[-1] == 'ID' or tokens[-1] == 'INTEGER':
                candidateDFAs.remove(INT)
            else:
                candidateDFAs.remove(ARITH)
    word_final += 1
```

Continuously, we used this code to deal with '-' symbol. To determine token of number with '-' symbol correctly, we need to see previous token. We divided into some cases to deal with this symbol.

1) 1 + -1 ➔ -1 (INT) , Previous token : OP

      1-1) a + -1 ➔ -1 (INT) , Previous token : OP

3) a = -1 ➔ -1 (INT), Previous token : ASSIGN

4) 1 – 1 ➔ - (OP), Previous token : INT

5) a – 1 ➔ - (OP), Previous token : ID

It means that if previous token is "INT" or "ID", '-' symbol will be classified as "OP" token and if previous token is **not** "INT" or "ID", '-' symbol will be classified as "INT" token. So we used if statements as image.

"if tokens[-1] == ' '" is to deal with blank space between minus symbol and previous token.

After that we increase "word_final" index to see the next character data of input string data.

```
fw.write("----------------------\n")
# Print lexemes+tokens info (except for whitespace)
for i in range(len(lexemes)):
    if not tokens[i] == 'WHITESPACE':
        strFormat = '%-12s%-12s'
        strOut = strFormat % (tokens[i], lexemes[i])
        fw.write(strOut + '\n')
fw.write("----------------------\n")
```

This part is for writing our lexeme and token result into output file.

```
parser = argparse.ArgumentParser()
parser.add_argument("input_file_name", help="file name of your input code")
args = parser.parse_args()
fr = open(args.input_file_name, 'r')
data = fr.read()
lexicalAnalysis(data + ' ', args.input_file_name)
```

This part is for execution command and start point of the code.

We use " **_python3 lexical_analyzer.py input.txt_** " command to run this code.

## 3) Result Screenshots

### Command Line

```
root@ubuntu:/home/ubuntu/Documents/compiler-main# python3 lexical_analyzer.py in
put.txt
root@ubuntu:/home/ubuntu/Documents/compiler-main#
```

### Input File ( Input.txt )

```
 1 001
 2 0010
 3 0010a0010
 4 0010-10
 5 0010--10
 6 int main(){char if123='1';int 0a=a+-1;return -0;}
 7 a-1
 8 -0, 0abc0
 9 123if, 123if0
10 int char boolean string
11 0 -1 123
12 ' '
13 true false
14 " "
15 _id
16 if else while class return
17 + - * /
18 =
19 < > <= >= == !=
20 ;
21 []
22 ()
23 {}
24 ,
```

Plain Text ▾    Tab Width: 8 ▾        Ln 24, Col 2    ▾    INS

## Output File ( input.txt_output.txt )

```
 1 --------------------------
 2 INTEGER      0
 3 INTEGER      0
 4 INTEGER      1
 5 INTEGER      0
 6 INTEGER      0
 7 INTEGER      10
 8 INTEGER      0
 9 INTEGER      0
10 INTEGER      10
11 ID           a0010
12 INTEGER      0
13 INTEGER      0
14 INTEGER      10
15 OP           -
16 INTEGER      10
17 INTEGER      0
18 INTEGER      0
19 INTEGER      10
20 OP           -
21 INTEGER      -10
22 VTYPE        int
23 ID           main
24 LPAREN       (
25 RPAREN       )
26 LBRACE       {
27 VTYPE        char
28 ID           if123
29 ASSIGN       =
30 CHAR         '1'
31 SEMI         ;
32 VTYPE        int
33 INTEGER      0
34 ID           a
35 ASSIGN       =
36 ID           a
37 OP           +
38 INTEGER      -1
39 SEMI         ;
40 RETURN       return
41 OP           -
42 INTEGER      0
43 SEMI         ;
44 RBRACE       }
45 ID           a
46 OP           -
47 INTEGER      1
48 OP           -
49 INTEGER      0
50 COMMA        ,
51 INTEGER      0
52 ID           abc0
53 INTEGER      123
54 IF           if
```

```
42 INTEGER      0
43 SEMI         ;
44 RBRACE       }
45 ID           a
46 OP           -
47 INTEGER      1
48 OP           -
49 INTEGER      0
50 COMMA        ,
51 INTEGER      0
52 ID           abc0
53 INTEGER      123
54 IF           if
55 COMMA        ,
56 INTEGER      123
57 ID           if0
58 VTYPE        int
59 VTYPE        char
60 VTYPE        boolean
61 VTYPE        string
62 INTEGER      0
63 OP           -
64 INTEGER      1
65 INTEGER      123
66 CHAR         ' '
67 BOOL         true
68 BOOL         false
69 STRING       " "
70 ID           _id
71 IF           if
72 ELSE         else
73 WHILE        while
74 CLASS        class
75 RETURN       return
76 OP           +
77 OP           -
78 OP           *
79 OP           /
80 ASSIGN       =
81 COMP         <
82 COMP         >
83 COMP         <=
84 COMP         >=
85 COMP         ==
86 COMP         !=
87 SEMI         ;
88 LBRACKET     [
89 RBRACKET     ]
90 LPAREN       (
91 RPAREN       )
92 LBRACE       {
93 RBRACE       }
94 COMMA        ,
95 --------------------------
```