

# COMPILER PROJECT I 2022 - Documentation

## Tokens & Regular Expressions

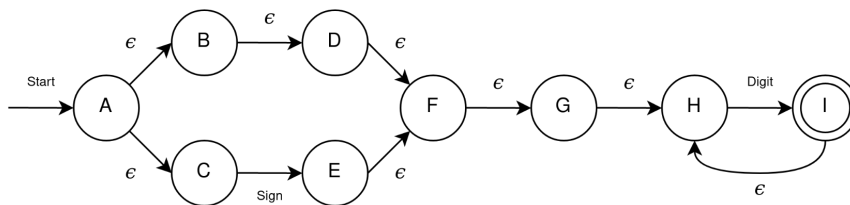
(arithmetic / assignment / comparison / symbols)

### Signed integer:

$\Sigma = \{0, 1, 2, \dots, 9, -\}$

Digit =  $0|1|2|\dots|9$

$r = (-|\epsilon)(\text{Digit})^+$



$T_0 = \epsilon\text{-closure}(A) = \{A\}$

$T_1 = \epsilon\text{-closure}(\delta(T_0, \text{Digit})) = \epsilon\text{-closure}(I) = \{I, H\}$

$T_2 = \epsilon\text{-closure}(\delta(T_0, \text{Sign})) = \epsilon\text{-closure}(E) = \{E, F, G, H\}$

$\epsilon\text{-closure}(\delta(T_1, \text{Digit})) = T_1$

$\epsilon\text{-closure}(\delta(T_1, \text{Sign})) = \emptyset$

$\epsilon\text{-closure}(\delta(T_2, \text{Digit})) = T_1$

$\epsilon\text{-closure}(\delta(T_2, \text{Sign})) = \emptyset$

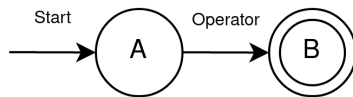
	Letter	Digit
T0	T1	T2
T1	T1	∅
T2	T1	∅

## Operators:

$\Sigma = \{+, -, *, /, =, <, >, ==, !=, <=, >=, :, \{, \}, (, ), ,\}$

Operator =  $+|-|*|/|=|<|>|=|!=|<=|>=|;|\{| \}|\(| \)|$ ,

$r = \text{Operator}$



	Operator
T0	T1
T1	$\emptyset$

$T0 = \epsilon\text{-closure}(A) = \{A\}$

$T1 = \epsilon\text{-closure}(\delta(T0, \text{Operator})) = \{B\}$

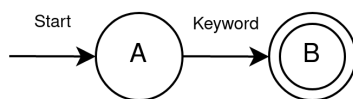
$\epsilon\text{-closure}(\delta(T1, \text{Operator})) = \emptyset$

## Keywords:

$\Sigma = \{\text{if, IF, else, ELSE, while, WHILE, return, RETURN, int, INT, char, CHAR}\}$

Keyword =  $\text{if|IF|else|ELSE|while|WHILE|return|RETURN|int|INT|char|CHAR}$

$r = \text{Keyword}$



$T0 = \epsilon\text{-closure}(A) = \{A\}$

$T1 = \epsilon\text{-closure}(\delta(T0, \text{Keyword})) = \{B\}$

$\epsilon\text{-closure}(\delta(T1, \text{Keyword})) = \emptyset$

	Keyword
T0	T1
T1	$\emptyset$

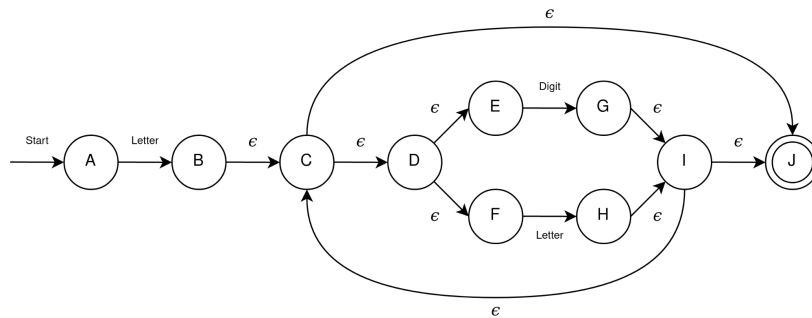
## Identifiers:

$\Sigma = \{A, B, C, \dots, Z, a, b, c, \dots, z, 0, 1, 2, \dots, 9\}$

Letter =  $A|B|C|\dots|Z|a|b|c|\dots|z$

Digit =  $0|1|2|\dots|9$

$r = \text{Letter}(\text{Letter}|\text{Digit})^*$



$T_0 = \epsilon\text{-closure}(A) = \{A\}$

$T_1 = \epsilon\text{-closure}(\delta(T_0, \text{Letter})) = \epsilon\text{-closure}(B) = \{B, C, D, E, F, J\}$ ,  $\epsilon\text{-closure}(\delta(T_0, \text{Digit})) = \emptyset$

$T_2 = \epsilon\text{-closure}(\delta(T_1, \text{Letter})) = \epsilon\text{-closure}(H) = \{D, E, F, H, I, J\}$

$T_3 = \epsilon\text{-closure}(\delta(T_1, \text{Digit})) = \epsilon\text{-closure}(G) = \{D, E, F, G, I, J\}$

$\epsilon\text{-closure}(\delta(T_2, \text{Letter})) = \{D, E, F, H, I, J\} = T_2$

$\epsilon\text{-closure}(\delta(T_2, \text{Digit})) = \{D, E, F, G, I, J\} = T_3$

$\epsilon\text{-closure}(\delta(T_3, \text{Letter})) = \{D, E, F, H, I, J\} = T_2$

$\epsilon\text{-closure}(\delta(T_3, \text{Digit})) = \{D, E, F, G, I, J\} = T_3$

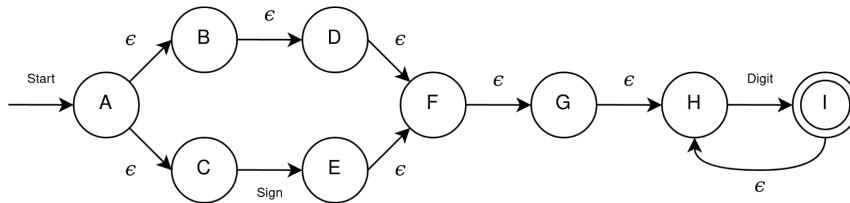
	Letter	Digit
T0	T1	$\emptyset$
T1	T2	T3
T2	T2	T3
T3	T2	T3

## Signed integer:

$\Sigma = \{0, 1, 2, \dots, 9, -\}$

Digit =  $0|1|2|\dots|9$

$r = (-|\epsilon)(\text{Digit})^+$



$T_0 = \epsilon\text{-closure}(A) = \{A\}$

$T_1 = \epsilon\text{-closure}(\delta(T_0, \text{Digit})) = \epsilon\text{-closure}(I) = \{I, H\}$

$T_2 = \epsilon\text{-closure}(\delta(T_0, \text{Sign})) = \epsilon\text{-closure}(E) = \{E, F, G, H\}$

$\epsilon\text{-closure}(\delta(T_1, \text{Digit})) = T_1$

$\epsilon\text{-closure}(\delta(T_1, \text{Sign})) = \emptyset$

$\epsilon\text{-closure}(\delta(T_2, \text{Digit})) = T_1$

$\epsilon\text{-closure}(\delta(T_2, \text{Sign})) = \emptyset$

	Letter	Digit
T0	T1	T2
T1	T1	$\emptyset$
T2	T1	$\emptyset$

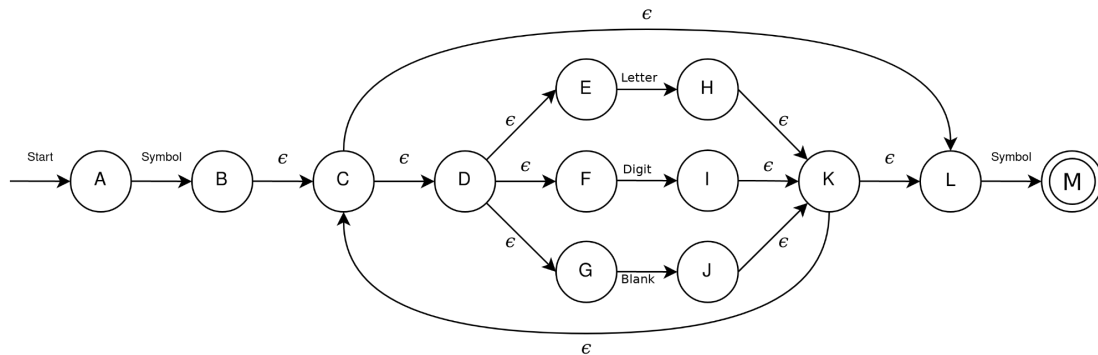
## Literal string:

$\Sigma = \{A, B, C, \dots, Z, a, b, c, \dots, z, 0, 1, 2, \dots, 9, , "\}$

Letter =  $A|B|C|\dots|Z|a|b|c|\dots|z$

Digit =  $0|1|2|\dots|9$

$r = \text{"(Letter|Digit| )*"}$



$T_0 = \epsilon\text{-closure}(A) = \{A\}$

$T_1 = \epsilon\text{-closure}(\delta(T_0, \text{Symbol})) = \epsilon\text{-closure}(B) = \{B, C, D, E, F, L\}$

$\epsilon\text{-closure}(\delta(T_0, \text{Symbol})) = \emptyset$

$\epsilon\text{-closure}(\delta(T_0, \text{Symbol})) = \emptyset$

$\epsilon\text{-closure}(\delta(T_0, \text{Symbol})) = \emptyset$

$T_2 = \epsilon\text{-closure}(\delta(T_1, \text{Letter})) = \epsilon\text{-closure}(H) = \{H, K, L, C, D, E, F, G\}$

$T_3 = \epsilon\text{-closure}(\delta(T_1, \text{Digit})) = \epsilon\text{-closure}(I) = \{I, K, L, C, D, E, F, G\}$

$T_4 = \epsilon\text{-closure}(\delta(T_1, \text{Blank})) = \epsilon\text{-closure}(J) = \{J, K, L, C, D, E, F, G\}$

$T_5 = \epsilon\text{-closure}(\delta(T_1, \text{Symbol})) = \epsilon\text{-closure}(M) = \{M\}$

$\epsilon\text{-closure}(\delta(T_2, \text{Letter})) = T_2$

$\epsilon\text{-closure}(\delta(T_2, \text{Digit})) = T_3$

$\epsilon\text{-closure}(\delta(T_2, \text{Blank})) = T_4$

$\epsilon\text{-closure}(\delta(T_2, \text{Symbol})) = T_5$

$\epsilon\text{-closure}(\delta(T_3, \text{Letter})) = T_2$

$\epsilon\text{-closure}(\delta(T_3, \text{Digit})) = T_3$

$\epsilon\text{-closure}(\delta(T_3, \text{Blank})) = T_4$

$\epsilon\text{-closure}(\delta(T_3, \text{Symbol})) = T_5$

$\epsilon\text{-closure}(\delta(T_4, \text{Letter})) = T_2$

$\epsilon\text{-closure}(\delta(T_4, \text{Digit})) = T_3$

$\epsilon\text{-closure}(\delta(T_4, \text{Blank})) = T_4$

$\epsilon\text{-closure}(\delta(T_4, \text{Symbol})) = T_5$

$\epsilon\text{-closure}(\delta(T_5, \text{Letter})) = \emptyset$

$\epsilon\text{-closure}(\delta(T_5, \text{Digit})) = \emptyset$

$\epsilon\text{-closure}(\delta(T_5, \text{Blank})) = \emptyset$

$\epsilon\text{-closure}(\delta(T_5, \text{Symbol})) = \emptyset$

	Letter	Digit	Blank	Symbol
T0	Ø	Ø	Ø	T1
T1	T2	T3	T4	T5
T2	T2	T3	T4	T5
T3	T2	T3	T4	T5
T4	T2	T3	T4	T5
T5	Ø	Ø	Ø	Ø

## Implementation details

We have made a diagram (projectRoot/diagram.png) showing the overall procedures of our lexical analysis. This will help in understanding the flow of the program and how we proceed to recognize tokens.

In the file 'lexical\_analyser.cpp', contains the loop in order to check for tokens among the string from the input file given to the program. Instead of calling every function (string\_analyser, integer\_analyser, assignment\_analyser, etc...) that were made to check the type of the token, we decided to create an array of pointer functions for better visibility of the code.

The pointer function was implemented like that:

```
bool (*func_ptrs[]) (
    const std::string &,
    std::size_t &,
    Scanner &
) = {
    vtype_analyzer,
    integer_analyzer,
    string_analyzer,
    identifier_analyzer,
    keyword_analyzer,
    arithmetic_analyzer,
    comparison_analyzer,
    assignment_analyzer,
    semi_analyzer,
    brace_analyzer,
    paren_analyzer,
    comma_analyzer
};
```

As you can see, this array will only contains function that follow this specific function signature : **bool funcname(const std::string &, std::size\_t &, Scanner &).**

We will iterate over this array while sending to each of these functions: our input string, the current index, and our Scanner structure containing a list of our tokens found during the lexical analysis. The loop will continue until one of the analyzers returns a true boolean, which means that a new token has been found and pushed to our list of tokens, or until we reach the end of the array meaning it will result in throwing an error because no token has been found.