



# Tiny Lexical Analyzer Documentation

Tiny Programming Language Compiler

---

Design And Implemente a Lexical Analyzer for Tiny  
programming Language using java Language

---

BY ISMAILEL KABOURI

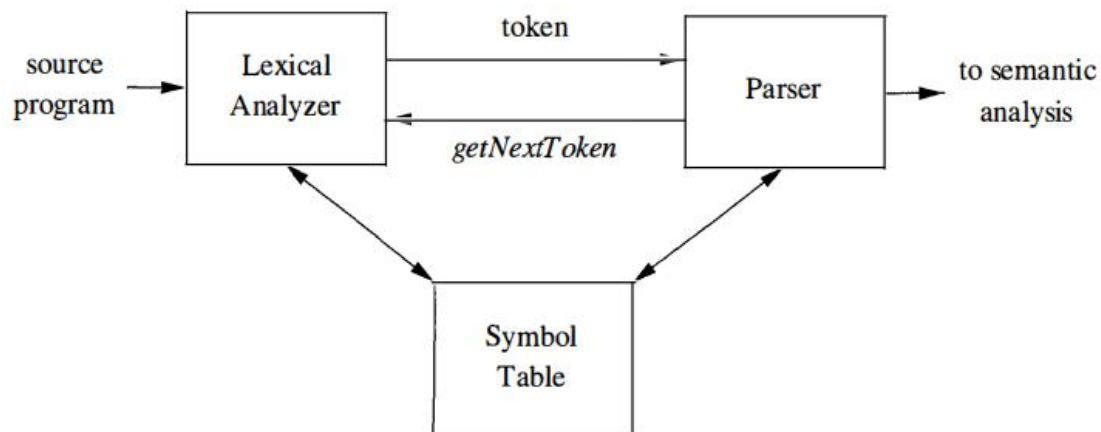
伊斯玛

## Catalog

<b>1 - Introduction .....</b>	<b>1</b>
<b>2 - Tiny Programing Language .....</b>	<b>2</b>
1 - Introduction .....	2
2 - Tiny Language Structure and specifications .....	2
3 -Tiny Language Token .....	5
4 - TINY Language Regular Expressions .....	6
<b>3 - Design the Lexical Analyzer .....</b>	<b>6</b>
3-1 Regular Expressions to NFA .....	6
3-2 NFA to DFA .....	7
3-3 DFA for Tiny Language .....	7
3-4 DFA implementation with Jflax. ....	9
<b>4 Tiny language sample .....</b>	<b>11</b>
4 - 1 Sample one .....	11
4-2 Sample tow .....	12
<b>5 Ackowlegement .....</b>	<b>13</b>
<b>6 - References .....</b>	<b>13</b>

# 1 - Introduction

At the heart of every compiler lies the crucial role of the lexical analyzer which reads the input characters of the source program, groups them into lexemes, and produces as output a sequence of tokens for each lexeme in the source program. The stream of tokens is sent to the parser for syntax analysis. It is common for the lexical analyzer to interact with the symbol table as well. When the lexical analyzer discovers a lexeme constituting an identifier, it needs to enter that lexeme into the symbol table



Since the lexical analyzer is the part of the compiler that reads the source text, it may perform certain other tasks besides identification of lexemes. 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). Another task is correlating error messages generated by the compiler with the source program.

Therefore lexical analyzers consist of

- *Scanning* consists of the simple processes that do not require tokenization of the input, such as deletion of comments and compaction of consecutive whitespace characters into one.
- *Lexical analysis* proper is the more complex portion, where the scanner produces the sequence of tokens as output.

## 2 - Tiny Programing Language

### 1 - Introduction

TINY is a simplified programming language often used for educational purposes and compiler construction projects. Its minimalist syntax and structure make it an ideal platform for learning fundamental concepts such as lexical analysis, parsing, and code generation

A program in TINY consists of a set of functions (any number of functions and ends with a main function), each function is a sequence of statements including (declaration, assignment, write, read, if, repeat, function, comment, ...) each statement consists of (number, string, identifier, expression, condition, ...).

### 2 - Tiny Language Structure and specifications

- *Number*: any sequence of digits (e.g. 123 | 554 | 0.23 | ...)
- *String*: starts with double quotes followed by any combination of characters and digits then ends with double quotes (e.g. "Hello" | "2nd + 3rd" | ...)
- *Reserved\_Keywords*: int | float | string | read | write | repeat | until | if | elseif | else | then | return | endl
- *Comment\_Statement*: starts with /\* followed by any combination of characters and digits then ends with / (e.g. / comment\*/ | ...)
- *Identifiers*: starts with letter then any combination of letters and digits. (e.g. x | val | counter1 | str1 | s2 | ...)
- *Function\_Call*: starts with Identifier then left bracket "(" followed by zero or more Identifier separated by "," and ends with right bracket ")". (e.g. sum(a,b) | factorial(c) | ...)

- **Arithmetic\_Operator:** any arithmetic operation (+ | - | \* | /)
- **Assignment\_Statement:** starts with Identifier then assignment operator “:=” followed by Expression (e.g. x := 1 | y:= 2+3 | ...)
- **Datatype:** set of reserved keywords (int, float, string)
- **Declaration\_Statement:** starts with Datatype then one or more identifiers (assignment statement might exist) separated by coma and ends with semi-colon. (e.g. int x; | ...)
- **Write\_Statement:** starts with reserved keyword “write” followed by an Expression or endl and ends with semi-colon (e.g. write x; | write 5; | write 3+5; | write “Hello World”; | ...)
- **Read\_Statement:** starts with reserved keyword “read” followed by an Identifier and ends with semi-colon (e.g. read x; | ...)
- **Return\_Statement:** starts with reserved keyword “return” followed by Expression then ends with semi-colon (e.g. return a+b; | return 5; | return “Hi”; | ...)
- **Condition\_Operator:** ( less than “<” | greater than “>” | is equal “=” | not equal “<>”)
- **Condition:** starts with Identifier then Condition\_Operator then Term (e.g. z1 < 10)
- **Boolean\_Operator:** AND operator “&&” and OR operator “||”
- **Condition\_Statement:** starts with Condition followed by zero or more Boolean\_Operator and Condition (e.g. x < 5 && x > 1)
- **If\_Statement:** starts with reserved keyword “if” followed by Condition\_Statement then reserved keyword “then” followed by set of Statements (i.e. any type of statement: write, read, assignment, declaration, ...) then Else\_If\_Statment or Else\_Statment or reserved keyword “end”
- **Else\_If\_Statement:** same as if statement but starts with reserved keyword “elseif”
- **Else\_Statement:** starts with reserved keyword “else” followed by a set of Statements then ends with reserved keyword “end”

- *Repeat\_Statement*: starts with reserved keyword “repeat” followed by a set of Statements then reserved keyword “until” followed by Condition\_Statement
- *Parameter*: starts with Datatype followed by Identifier (e.g. int x)
- *Function\_Declaration*: starts with Datatype followed by FunctionName followed by “(“ then zero or more Parameter separated by “,” then “)” (e.g. int sum(int a, int b) | ...)
- *Function\_Body*: starts with curly bracket “{” then a set of Statements followed by Return\_Statement and ends with “}”
- *Function\_Statement*: starts with Function\_Declaration followed by Function\_Body
- *Main\_Function*: starts with Datatype followed by reserved keyword “main” then “()” followed by Function\_Body
- *Program*: has zero or more Function\_Statement followed by Main\_Function

### 3 -Tiny Language Token

type	lexeme	Token Name (code)	Attribute
keywords	if	<b>if</b> (261)	-
	then	<b>then</b> (262)	-
	else	<b>else</b> (263)	-
	end	<b>end</b> (264)	-
	repeat	<b>repeat</b> (265)	-
	until	<b>until</b> (266)	-
	read	<b>read</b> (267)	-
	write	<b>write</b> (268)	-
relation operators	<	<b>relop</b> (270)	LT (271)
	<=	<b>relop</b>	LE (272)
	=	<b>relop</b>	EQ (273)
	>	<b>relop</b>	NE (274)
	>	<b>relop</b>	GT (275)
	>=	<b>relop</b>	GE (276)
arithmetical operators	+	<b>addop</b> (280)	ADD (281)
	-	<b>addop</b> (280)	MINUS (282)
	*	<b>mulop</b> (285)	MUL(286)
	/	<b>mulop</b> (285)	DIV(287)
	(	( (294)	-
	)	) (295)	-
assignment	:=	:= (296)	-
Segment	;	; (297)	-
numbers	Such as <i>12342</i>	<b>num</b> (298)	Symbol Table Entry
identifiers	Such as <i>student1</i>	<b>id</b> (299)	Symbol Table Entry
white spaces (ws)	<i>blank, tab, and newline</i>	-	-
Comments	{bla, bala, bala }	-	-

## 4 - TINY Language Regular Expressions

```
digit ::= 0/1/2/3/4...../9
letter ::= [a-z][A-Z]
Number ::= digit.?digit
String ::= "(letter/digit)*"
Reserved_Keywords ::= int | float | string | read | write | repeat
| until | if | elseif | else | then | return | endl
Comment ::= /*String*\
Identifier ::= letter (letter | digit)*
Arithmetic_Operator ::= + | - | * | /
Datatype ::= int | float | string
Condition_Operator ::= < | > | = | <>
Boolean_Operator ::= && | \/\
```

## 3 - Design the Lexical Analyzer

### 3 -1 Regular Expressions to NFA

Each regular expression is first converted into an NFA (Non-Deterministic Finite Automaton). This step involves creating states and transitions for each character in the regular expression, we can do that by hand or use some tools which can generate the automaton for us like "Gragpize" with the "Dot Language" which accept states and their transitions, and we can add some more attributes for good appearance like "color" or "style". Even we get the NFA representation it still a little difficult during implementation, coz the NFA can make the "Epsilon" transitions to navigate between states, and that not efficient at all we need to defined unique path from start state to accepting state for each regex, to avoid confusion during scanning and analysing of lexemes.



## 3-2 NFA to DFA

The NFA is then converted into a DFA (Deterministic Finite Automaton) which has a unique transition for each character in the alphabet from each state. and again we can do that by hand by following subset construction algorithm, which consider state of the constructed DFA corresponds to a set of NFA states. And remove “epsilon” transition, in DFA each token has a unique path from start state to accepting state, then minimize DFA to make it more simple and easy to implement. Or by using some tools like “Graphviz “ and “Dot language “ which generate well done DFA representation, with nice appearance.

## 3-3 DFA for Tiny Language

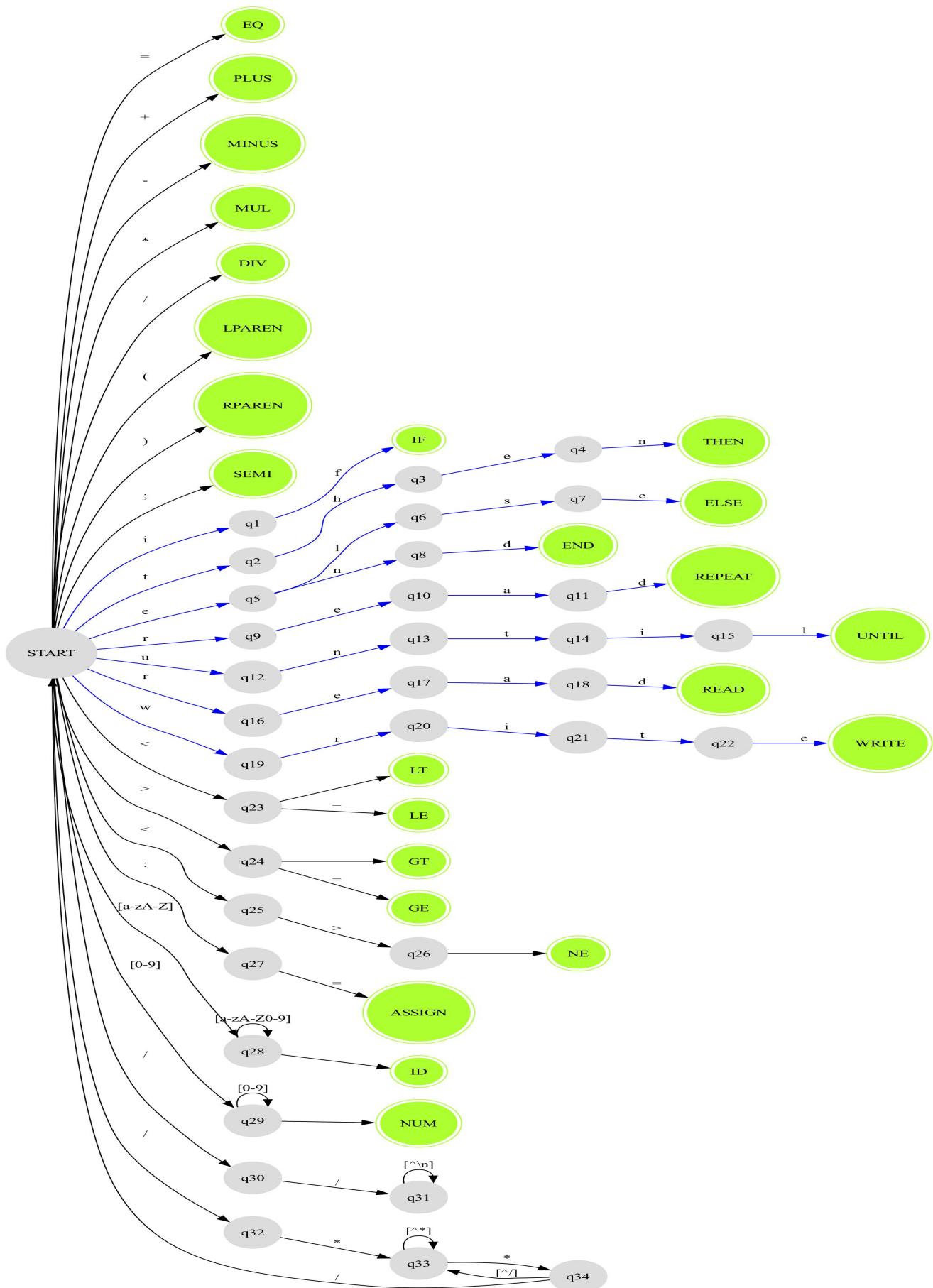
We use “Graphviz” and Dot Language to generate the DFA for common Tiny lexemes.

Dot Language Structure

```
-----  
digraph DFA {  
  
size = 10.8;  
node["Attributes"] // which represent the final states .  
node["Attributes"] // which represent other states.  
  
s0["Attributes"] // initail state  
  
// Transitions goeas here  
  
// Keywords  
.  
.  
.  
  
// Identifiers  
.  
.  
.  
-----
```

To generate the DFA graph, use Graphviz with the following command:

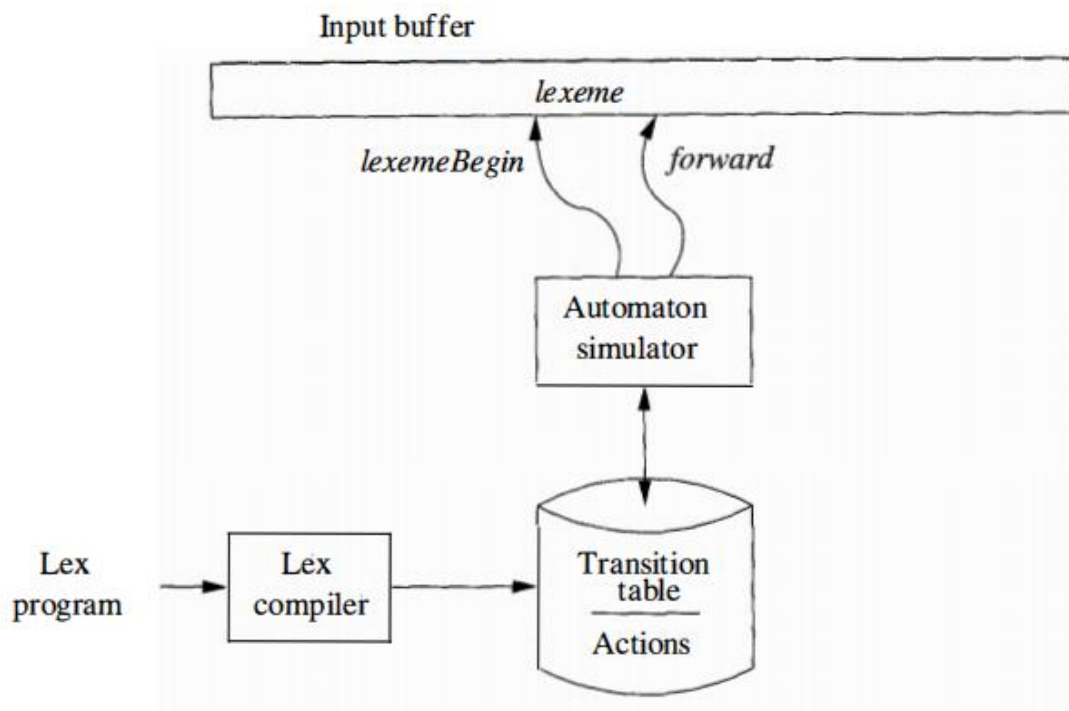
```
-----  
$$ --> dot -Tpng -Gdpi=1000 dfa.dot -o dfa.png  
-----
```



### 3-4 DFA implementation with Jflex.

Implementing DFA by hand can be challenging specially when we deal with complex language with wide range of tokens and specifications. To avoid complexity we can use a program to automatically transform a set of regular expressions into code for a scanner. Such a program is known as a scanner generator. As you know I choose java language to do the implementation of that lexer. The appropriate generator for that is "Jflex" which is a lexical analyser generator for Java written in Java. It takes as input a specification with a set of regular expressions and corresponding actions. It generates a program (a lexer) that reads input, matches the input against the regular expressions in the spec file, and runs the corresponding action if a regular expression matched.

#### - The Structure of the Generated Analyzer



## Jflex structure

```
// some import stm...

%%-----

%class className

    %unicode

    %public

    %type Token

    %{-----

                                // java code

    %}-----

    %state COMMENT                // decleration section

%% -----

<YYINITIAL> {

    "if"                { return new Token(sym.IF, yytext()); }

    "else"              { return new Token(sym.ELSE, yytext()); }

    "read"              { return new Token(sym.READ, yytext()); }

    "write"             { return new Token(sym.WRITE, yytext()); }

    "<"                { return new Token(sym.LT, yytext()); }

    ">"                { return new Token(sym.GT, yytext()); }

    "+"                 { return new Token(sym.PLUS, yytext()); }

    "-"                 { return new Token(sym.MINUS, yytext()); }

    [a-zA-Z][a-zA-Z0-9]* { return new Token(sym.ID, yytext()); }

    "/"                { yybegin(COMMENT); } { /* Ignore whitespaces */ }

<COMMENT> {

    \n                  { yybegin(YYINITIAL); newline(); }

    .                   { /* Ignore characters in single-line comments */ }

    -----
}
```

## 4 Tiny language sample

### 4 - 1 Sample one

INPUTE	OUTPUT	
<b>if x &lt; 10 then</b>  <b>read y;</b>  <b>write y + 1;</b>  <b>end</b>	Token: <READ, 'read'>	
	Token: <ID, 'x'>	Token: <ID, 'x'>
	Token: <SEMI, ';'>	Token: <SEMI, ';'>
	Token: <IF, 'if'>	Token: <ID, 'x'>
	Token: <ID, 'x'>	Token: <ASSIGN, ':=>
	Token: <GT, '>'>	Token: <ID, 'x'>
	Token: <NUM, '0'>	Token: <MINUS, '-'>
	Token: <THEN, 'then'>	Token: <NUM, '1'>
	Token: <ID, 'fact'>	Token: <UNTIL, 'until'>
	Token: <ASSIGN, ':=>	Token: <ID, 'x'>
	Token: <NUM, '1'>	Token: <EQ, '='>
	Token: <SEMI, ';'>	Token: <NUM, '0'>
	Token: <REPEAT, 'repeat'>	Token: <SEMI, ';'>
	Token: <ID, 'fact'>	Token: <WRITE, 'write'>
	Token: <ASSIGN, ':=>	Token: <ID, 'fact'>
	Token: <ID, 'fact'>	Token: <END, 'end'>
	Token: <MUL, '*'>	

## 4-2 Sample tow

INPUTE	OUTPUT	
<pre> int main() {     int x;     read x;     if x &gt; 0     then         int fact := 1;          repeat             fact := fact * x;             x := x - 1;         until x = 0          write fact;     end     return 0; } </pre>	Token: <ID, 'int'>	Token: <ASSIGN, ':= '>
	Token: <ID, 'main'>	Token: <ID, 'fact'>
	Token: <LPAREN, '('>	Token: <MULT, '*'>
	Token: <RPAREN, ')'>	Token: <ID, 'x'>
	Token: <LBRACE, '{'>	Token: <SEMI, ';'>
	Token: <ID, 'int'>	Token: <ID, 'x'>
	Token: <ID, 'x'>	Token: <ASSIGN, ':= '>
	Token: <SEMI, ';'>	Token: <ID, 'x'>
	Token: <READ, 'read'>	Token: <MINUS, '-'>
	Token: <ID, 'x'>	Token: <NUM, '1'>
	Token: <SEMI, ';'>	Token: <SEMI, ';'>
	Token: <IF, 'if'>	Token: <UNTIL, 'until'>
	Token: <ID, 'x'>	Token: <ID, 'x'>
	Token: <GT, '>'>	Token: <EQ, '='>
	Token: <NUM, '0'>	Token: <NUM, '0'>
	Token: <THEN, 'then'>	Token: <WRITE, 'write'>
	Token: <ID, 'int'>	Token: <ID, 'fact'>
	Token: <ID, 'fact'>	Token: <SEMI, ';'>
	Token: <ASSIGN, ':= '>	Token: <END, 'end'>
	Token: <NUM, '1'>	Token: <RETURN, 'return'>
	Token: <SEMI, ';'>	Token: <NUM, '0'>
	Token: <REPEAT, 'repeat'>	Token: <SEMI, ';'>
	Token: <ID, 'fact'>	Token: <RBRACE, '}'>

## 5 Acknowledgement

I would like to express my gratitude to everyone who supported and guided me throughout the development of this project.

Special thanks to:

- Professors and Mentors: For their invaluable guidance and insightful feedback.
- Open Source Community: For providing tools like JFlex and resources that were crucial in building this lexical analyzer.

This project would not have been possible without the contributions and support from these individuals and communities.

## 6 - References

- <https://www3.nd.edu/~dthain/compilerbook/>
- <https://www.jflex.de/manual.html>
- <https://a7medayman6.github.io/Tiny-Compiler/>
- <https://www.graphviz.org/>
- Compiler Principles Techniques and Tools (2<sup>nd</sup> Edition) Book.

Made by : ISMAIL EL KABOURI  
ID : 92115855E124  
University: NJUST  
Contractor: Kedong Yan  
Location : China

