

# INFO-F403 - Project : Part1

Julien Baudru - N°000460130

October 2020

## 1 Introduction

The first part of this project consisted in building a lexical analyzer of a future compiler for the FORTR-S language.

## 2 Structure

The structure of the project is the following :

```
1 part1:
2 |
3 |---dist
4 |   |---part1.jar
5 |
6 |---doc
7 |   |---Javadoc
8 |   |   |--- ...
9 |   |
10 |   |---Rapport.pdf
11 |
12 |---more
13 |
14 |---src
15 |   |---LexicalAnalyzer.java
16 |   |---LexicalAnalyzer.flex
17 |   |---LexicalUnit.java
18 |   |---Main.java
19 |   |---Symbol.java
20 |   |---SymbolTable.java
21 |
22 |---test
23 |   |---Factorial_3.fs
24 |   |---Fibonacci.fs
25 |   |---Test.fs
```

**Note :** The file index.html as well as all other files of the type html constitute the JavaDoc of this project. To compile and execute this project the following instructions must be followed:

```
1 jflex src/LexicalAnalyzer.flex
2 javac -d bin -cp src/ src/Main.java
3 jar cfe dist/Part1.jar Main -C bin
4 java -jar dist/Part1.jar test/Factorial.fs
```

When the file *Main.java* is executed, a loop will iterate through all the tokens in the file until it encounters the symbol **EOS** (the symbol for the end of the file). The symbols are obtained thanks to the *nextToken()* function of the *LexicalAnalyzer.java* class, this class is generated thanks to the file *LexicalAnalyzer.flex* by executing the command below :

```
1 C:Part1> jflex LexicalAnalyzer.flex
```

In the class *LexicalAnalyzer.java*, the symbols are extracted from the tokens via a set of regular expressions detailed in the point 3 of this document. From the set of symbols received by the function *nextToken()*, a table of symbols will be created via the class *SymbolTable.java*

However, this symbol table still contains the line and block comments of the initial file, as these have to be ignored by the scanner, the function *DeleteComment()* from the class *SymbolTable.java* is called. This function scrolls through the symbol table and deletes everything between a starting X character and an ending Y character, for the line comments the signature of this function is :

```
1 st.DeleteComment("//", "\\n");
```

and for the block comments his signature is : .

```
1 st.DeleteComment ("/*", "*/");
```

The choice to delete comments after the construction of the symbol table has been made to simplify the writing of regular expressions.

Moreover, this technique would allow to handle nested comments, in fact, as explained above, the function *DeleteComment()* browses the entries of the symbol table looking for the part to delete, so if we have a comment of the following type :

```
1 /*
2 Frist line
3 // Second line
4 Last line
5 /*
```

So we could simply call this function with the same signature as for block comments and the nested comment will be ignored as it should be.

### 3 Regular expressions

Lexical Unit	Regular expression
BEGINPROG	"BEGINPROG"
PROGNAME	({AlphaUpperCase}({AlphaUpperCase}*{AlphaLowerCase}*{Numeric}*)+)
ENDLINE	"\n"
ENDPROG	"ENDPROG"
COMMA	","
VARNAME	({AlphaLowerCase}({AlphaLowerCase}*{Numeric}*)+)
ASSIGN	":="
NUMBER	{Integer}    {Decimal}
LPAREN	"("
RPAREN	")"
MINUS	"_"
PLUS	"+"
TIMES	"*"
DIVIDE	"/"
IF	"IF"
THEN	"THEN"
ENDIF	"ENDIF"
ELSE	"ELSE)"
EQ	"=="
GT	">"
WHILE	"WHILE"
DO	"DO"
ENDWHILE	"ENDWHILE"
PRINT	"PRINT"
READ	"READ"
EOS	null

**Note :** To allow more clarity during the reading, other regular expressions have been used to build the expressions above, here is how they were defined :

Alias for regular expression	Regular expression
{Numeric}	[0-9]
{Zero}	0
{ZeroPlus}	0[0-9]
{Integer}	(([1-9][0-9]*)    {Zero}
{Decimal}	{Integer}"."([0-9]+)    {Zero}"."([0-9]+)
{AlphaUpperCase}	[A-Z]
{AlphaLowerCase}	[a-z]
{OnlyUpperCase}	{AlphaUpperCase}+

The expression {ZeroPlus} is used to detect cases where numbers contain unnecessary zeros such as **002** or **047**. If such a case is recognized then it is changed from *YYINITIAL* to *ZEROCASE* state in the *LexicalAnalyzer.flex* file, so these numbers are ignored by the scanner.

In the same logic, the expression *OnlyUpperCase* is used to detect symbols of type PROGNAME which are only in upper case, these will also be ignored by the scanner.

In addition to this, although it is not specified in the statement, the scanner is also able to recognize comma numbers because almost all programming languages are able to do so. This is why the regular expression allowing to recognize the lexical unit of type NUMBER is composed of the regular expressions {Integer} and {Decimal}.

## 4 Example of FORTR-S

### 4.1 Fibonacci program

The file given as an argument to our lexical analyzer is the following, this program print the  $n$  first numbers of the Fibonacci suite. This example was chosen because it uses almost all the types of lexical units provided by the FORTR-S language.

```

1 BEGINPROG Fibonacci
2 READ(n)
3 a := 0
4 olda := a
5 b := 1
6 i := 0
7 WHILE(n > i) DO
8     olda := a
9     a := b
10    b := olda + b
11    i := i + 1
12    PRINT(a)
13 ENDWHILE
14 ENDPROG

```

### 4.2 Symbol table of Fibonacci program

Below you will find the symbol table returned by the lexical analyzer for the program *Fibonacci.fs*.

**Note :** To have exactly the same symbol table as below, the line separator of the file must be \n (i.e. LF) or \r (i.e. CR) but not both at the same time (i.e. CRLF).

1 token: BEGINPROG	lexical unit: BEGINPROG
2 token: Fibonacci	lexical unit: PROGNAME
3 token: \n	lexical unit: ENDLINE
4 token: READ	lexical unit: READ
5 token: (	lexical unit: LPAREN
6 token: n	lexical unit: VARNAME
7 token: )	lexical unit: RPAREN
8 token: \n	lexical unit: ENDLINE
9 token: a	lexical unit: VARNAME
10 token: :=	lexical unit: ASSIGN
11 token: 0	lexical unit: NUMBER
12 token: \n	lexical unit: ENDLINE
13 token: olda	lexical unit: VARNAME
14 token: :=	lexical unit: ASSIGN
15 token: a	lexical unit: VARNAME
16 token: \n	lexical unit: ENDLINE
17 token: b	lexical unit: VARNAME
18 token: :=	lexical unit: ASSIGN
19 token: 1	lexical unit: NUMBER
20 token: \n	lexical unit: ENDLINE
21 token: i	lexical unit: VARNAME
22 token: :=	lexical unit: ASSIGN
23 token: 0	lexical unit: NUMBER
24 token: \n	lexical unit: ENDLINE
25 token: WHILE	lexical unit: WHILE
26 token: (	lexical unit: LPAREN
27 token: n	lexical unit: VARNAME
28 token: >	lexical unit: GT
29 token: i	lexical unit: VARNAME
30 token: )	lexical unit: RPAREN
31 token: DO	lexical unit: DO
32 token: \n	lexical unit: ENDLINE
33 token: olda	lexical unit: VARNAME
34 token: :=	lexical unit: ASSIGN
35 token: a	lexical unit: VARNAME
36 token: \n	lexical unit: ENDLINE

```

37 token: a          lexical unit: VARNAME
38 token: :=         lexical unit: ASSIGN
39 token: b          lexical unit: VARNAME
40 token: \n         lexical unit: ENDLINE
41 token: b          lexical unit: VARNAME
42 token: :=         lexical unit: ASSIGN
43 token: olda       lexical unit: VARNAME
44 token: +          lexical unit: PLUS
45 token: b          lexical unit: VARNAME
46 token: \n         lexical unit: ENDLINE
47 token: i          lexical unit: VARNAME
48 token: :=         lexical unit: ASSIGN
49 token: i          lexical unit: VARNAME
50 token: +          lexical unit: PLUS
51 token: 1          lexical unit: NUMBER
52 token: \n         lexical unit: ENDLINE
53 token: PRINT      lexical unit: PRINT
54 token: (          lexical unit: LPAREN
55 token: a          lexical unit: VARNAME
56 token: )          lexical unit: RPAREN
57 token: \n         lexical unit: ENDLINE
58 token: ENDWHILE   lexical unit: ENDWHILE
59 token: \n         lexical unit: ENDLINE
60 token: ENDPROG    lexical unit: ENDPROG
61 token: \n         lexical unit: ENDLINE
62
63 Variables :
64 a 3
65 b 5
66 i 6
67 n 2
68 olda 4

```