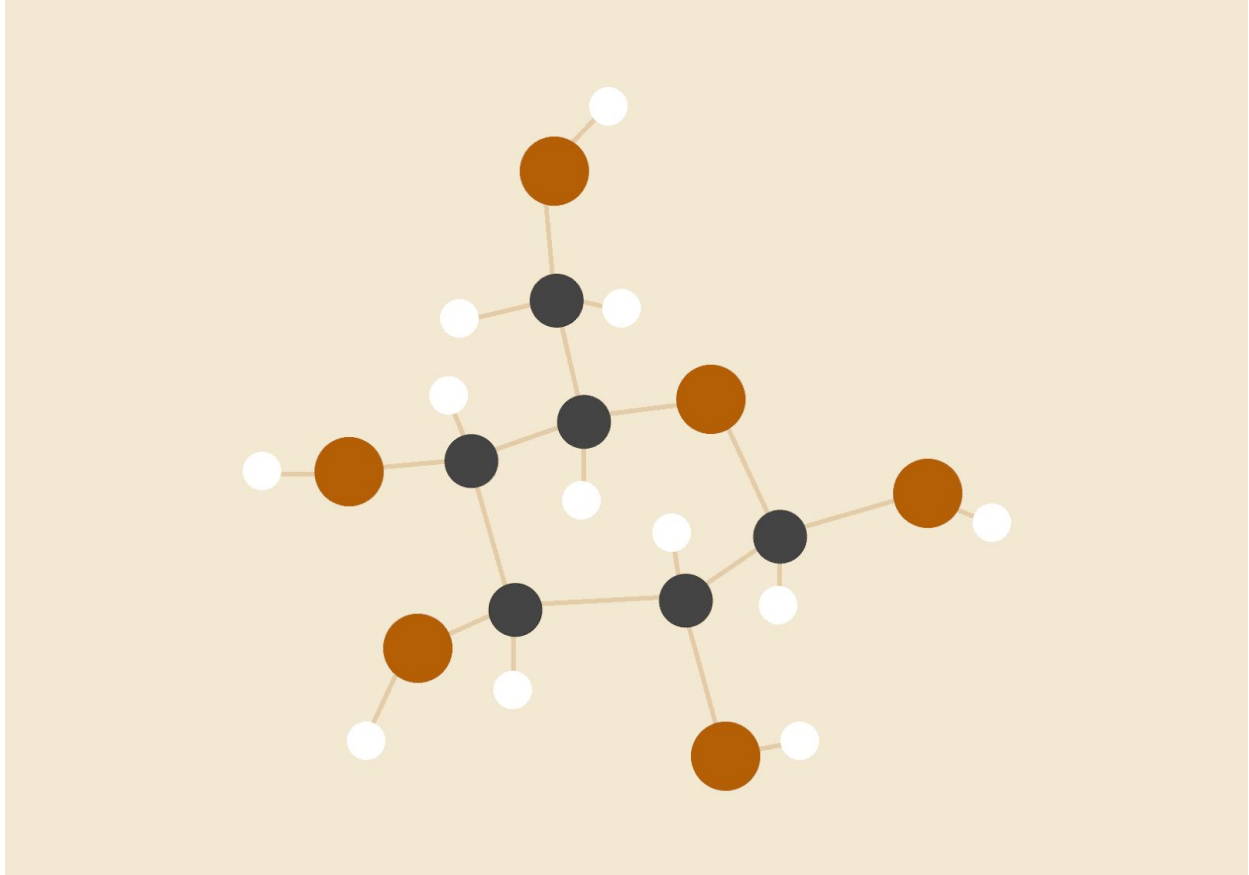# A Lexical and Syntax Analyser

*CA4003 Assignment One*

**Niall Lyons - 13493628**

12/11/18

CA4003 Compiler Construction

# Declaration

I declare that this material, which I/We now submit for assessment, is entirely my/our own work and has not been taken from the work of others, save and to the extent that such work has been cited and acknowledged within the text of my/our work. I/We understand that plagiarism, collusion, and copying are grave and serious offences in the university and accept the penalties that would be imposed should I engage in plagiarism, collusion or copying. I/We have read and understood the Assignment Regulations. I/We have identified and included the source of all facts, ideas, opinions, and viewpoints of others in the assignment references. Direct quotations from books, journal articles, internet sources, module text, or any other source whatsoever are acknowledged and the source cited are identified in the assignment references. This assignment, or any part of it, has not been previously submitted by me/us or any other person for assessment on this or any other course of study. I/We have read and understood the referencing guidelines found at http://www.dcu.ie/info/regulations/plagiarism.shtml , https://www4.dcu.ie/students/az/plagiarism and/or recommended in the assignment guidelines.


Name: Niall Lyons

Date: 12/11/18

## INTRODUCTION

The aim of this assignment is to implement a lexical and syntax analyser using JavaCC for a simple language called CAL.

The details of the CAL language are given [here.](here.)

I will be going through each of the four main components and briefly explaining how they were implemented. The four main sections of this report are:

- Options
- User Code
- Token Definitions
- Grammar and Production Rules

## OPTIONS

The parser will have two options as seen below.

```
1   /*
2      SECTION ONE — OPTIONS
3   */
4
5   options {
6      IGNORE_CASE = true;
7      JAVA_UNICODE_ESCAPE = true;
8   }
```

Option one, IGNORE_CASE is used as it is said in the details of CAL language that the language will not be case sensitive. Option two is the JAVA_UNICODE_ESCAPE, when set to true the generated parser uses an input stream object that processes Java Unicode escapes before sending characters to the token manager. By default, Java Unicode escapes are not processed.

## USER CODE

This section holds the class definition of the parser, here we must enclose our code between the PARSER_BEGINS and PARSER_END keywords.

```
14     PARSER_BEGIN(CALParser)        60     PARSER_END (CALParser)
```

The main method of the CALParser is pretty much the same as the SLPTokeniser parser code which is in the notes. I started off by initializing my parser by simply calling it parser as seen below.

```
21          //Parser initialization
22          CALParser parser;
```

The code then goes on to initialise the parser that I have defined to read from the correct place. The parser will either begin to read input from standard input or from a file that is given as an argument from the command line. It also will also print out of the file name that was inputted cannot be found.

```
24        if(args.length == 0){
25          System.out.println ("CALParser: Reading input ...");
26          parser = new CALParser(System.in);
27        }
28        else if(args.length == 1){
29          System.out.println ("CALParser: Reading the file " + args[0] + " ..." );
30          try {
31            parser = new CALParser(new java.io.FileInputStream(args[0]));
32          }
33          catch(java.io.FileNotFoundException e) {
34            System.out.println ("CALParser: The file " + args[0] + " was not found.");
35            return;
36          }
37        }
38        else {
39          System.out.println ("CALParser:  You must use one of the following:");
40          System.out.println ("        java CALParser < file");
41          System.out.println ("Or");
42          System.out.println ("        java CALParser file");
43          return ;
44        }
```

When the parser is initialized, it will then begin to try and parse either the input file or standard input from programme() which is defined below in the grammar section of the code. The parser will then go through the production rules and grammar to parse the file. This will result in either the file being parsed successfully, which will print a successful message or it will result in the file not being parsed successfully and will print the ParserException error message to the command line. I also use TokenMgrError here,

to accommodate any errors that may have occurred in the token definitions.

```
45        try {
46            parser.programme();
47            System.out.println("CALParser: The input was read successfully.");
48        }
49        catch(ParseException e){
50            System.out.println ("CALParser: There was an error during the parse.");
51            System.out.println (e.getMessage());
52        }
53        catch(TokenMgrError e){
54            System.out.println ("CALParser: There was an error.");
55            System.out.println (e.getMessage());
56        }
```

## TOKEN DEFINITIONS

There are many different aspects of token definitions which must be thought about when developing a parser. I split this into four sections:

- **Comments**

  I first began by declaring the comment tokens which must be ignored. I started off by declaring a commentNesting variable which is defined in the JavaCC notes. This form of a comment begins with a " /* " and ends with a " */ ".

  ```
  66    TOKEN_MGR_DECLS :
  67    {
  68        static int commentNesting = 0;
  69    }
  ```

  This comment nesting variable is incremented each time we see " /* " resulting in the lexical state being altered. We then move to a new state where the same variable is either decremented when we see " */ "  and incremented when we see " /* ". After every decrementation, it must be checked if it is equal to zero, if the variable is zero we move back to the default lexical state. This is declared below:

  ```
  89    <IN_COMMENT> SKIP :
  90    {
  91        "/*" { commentNesting++;}
  92      | "*/" { commentNesting--;
  93                if(commentNesting == 0)
  94                    SwitchTo(DEFAULT);
  95            }
  96      |<~[]>
  97    }
  98
  ```

4

I then had to take into consideration single line comments. This token is defined by:

```
82    SKIP :
83    {
84        "/*" { commentNesting++; } : IN_COMMENT
85      | <"//" (~["\n"])* "\n">
86    }
```

Here we also take into consideration the commentNesting variable which we have defined before and increment it if we find "/* " and use it in <IN_COMMENT> SKIP.

On line 85 above, <"//" (~["\n"])* "\n"> is defined. This will match any comments that begin with " // " and is the followed by anything other than a new line. This will go on until a new line character is seen.

- **Characters to be ignored**

Another SKIP block is then created in order to ignore tabs, spaces and newlines as we dont want to read these in as tokens. This is shown below:

```
72    SKIP:
73    {
74        "\n"
75      | "\r"
76      | "\r\n"
77      | "\t"
78      | " "
79      | "\f"
80    }
81
```

- **Keyword and Operators**

To declare keywords and operators you first need to define it as a token. This is done by defining the token, and followed by the matching string. This can be all done in one token block, but for the ease of reading the code I decided to split them up into two as you can see below. All of the tokens that are seen below have been defined in the CAL language definition.

```
102    TOKEN :                              123    TOKEN :
103    {                                    124    {
104        <COMMA:(",")>                     125        <VAR: "variable">
105      | <SCOLON:(";")>                    126      | <CONST: "constant">
106      | <COLON:(":")>                     127      | <RET: "return">
107      | <ASIGN:(":=")>                    128      | <INT: "integer">
108      | <OPBRA :("(")>                    129      | <BOOL: "boolean">
109      | <CLBRA:(")")>                     130      | <VOID: "void">
110      | <PLUS:("+")>                      131      | <MAIN: "main">
111      | <MINUS:("-")>                     132      | <IF: "if">
112      | <TILDA:("~")>                     133      | <ELSE: "else">
113      | <OR:("|")>                        134      | <TRUE: "true">
114      | <AND:("&")>                       135      | <FALSE: "false">
115      | <EQ:("=")>                        136      | <WHILE: "while">
116      | <NOTEQ:("!=")>                    137      | <BEGIN: "begin">
117      | <LT:("<")>                        138      | <END: "end">
118      | <LTEQ:("<=")>                     139      | <IS: "is">
119      | <GT:(">")>                        140      | <SKI: ("skip")>
120      | <GTEQ:(">=")>                     141    }
121    }
```

As you can see above, I used <SKI> as it would not conflict with SKIP that I defined in above for the comments and skipping of characters.

- **Numbers and Identifiers**

In the language definition it reads *" Identifiers are represented by a string of letters, digits or underscore character ('_ ') beginning with a letter. "* and integers were defined as *" Integers are represented by a string of one or more digits ('0'-'9') and may start with a minus sign ('-'), e.g. 123, -456. Unless it is the number '0', numbers may not start with leading '0's "* .

Below is how I implemented this:

```
143    TOKEN :
144    {
145        <IDENTIFIER: <LETTER>(<LETTER>|<DIGIT>|<SPECHAR>)*>
146      | <DIGIT: "0" | <MINUS> ["1" - "9"] (["0" - "9"])* | ["1" - "9"] (["0" - "9"])*>
147      | <LETTER: (["a"-"z","A"-"Z"])>
148      | <SPECHAR : "_"> //special characters
149    }
```

As you can see I created a letter token which takes any letter from a-z and also in capital letters. I also created a special character to accommodate the identifier. I called all integers digits, here I accommodated all numbers represented by strings of one or more digits and can also start with a minus sign, unless it is 0. Here I reused the minus token I made before. I then created the identifier, this shows an identifier is a letter followed by zero or more letters, digits or special characters,

this is because of the " * " which represents Kleene Closure.

## GRAMMAR AND PRODUCTION RULES

From the CAL language definition given I began to write out the grammar for my syntax analyser. I started by declaring my prog which the parser calls in the main method.

```
155    void programme () : {}
156    {
157      declartionList()functionList()mainMethod(){}
158    }
```

I quickly noticed after writing these out that there was a few changes that I had to make. After compiling my JavaCC file I found that I had some left recursion errors. From the original grammar that was given there was left recursion present in expression and fragment non terminals. Originally fragment and expression looked like this:

$$\langle expression \rangle \models \langle fragment \rangle \langle binary\_arith\_op \rangle \langle fragment \rangle \mid$$
$$( \langle expression \rangle ) \mid$$
$$identifier ( \langle arg\_list \rangle ) \mid$$

$$\langle fragment \rangle$$

$$\langle fragment \rangle \models identifier \mid - identifier \mid number \mid \textbf{true} \mid \textbf{false} \mid$$
$$\langle expression \rangle$$

I began to develop a different approach to this problem by firstly creating a new simple_expression non terminal and adding this to the top line of my expression non terminal, meaning I would not be calling fragment twice. I then also took out the expression non terminal from fragment, and added in an optional binOp() expression() to the end of each production in expression() using my new simple_expression non terminal. The result of this can be seen in the code snippet below:

```
263    void simple_expression() : {}
264    {
265        binOp() expression()
266        | {}
267    }
268
269    void expression () : {}
270    {
271
272        fragment()simple_expression()
273        | <OPBRA>expression()<CLBRA> simple_expression()
274        | <IDENTIFIER>[<OPBRA>argList()<CLBRA>] simple_expression()
275
276    }
277    |
278    void binOp(): {}
279    {
280        <PLUS>
281        | <MINUS>
282    }
283
284    void fragment(): {}
285    {
286        <DIGIT>
287        | boolOp()
288        | <MINUS> <IDENTIFIER>
289    }
290
```

I then went on to fix the next left recursion error. I used a similar approach to fix the condition non terminal, by creating a new simple_condition non terminal and using the same method as above.

```
291    void condition(): {}
292    {
293        simple_condition()(<OR>condition()| <AND>condition() |{})
294    }
295
296
297    //Because of the two conditions in test6 each side of and
298    void simple_condition() : {}
299    {
300        <TILDA> condition()(<OR>condition()| <AND>condition() |{})
301        | LOOKAHEAD(3) <OPBRA> condition() <CLBRA> (<OR> condition() | <AND> condition() |{})
302        | expression() compOp() expression() (<OR> condition()| <AND> condition() |{})
303
304    }
```

I also had to fix my statement non terminal as I had a choice conflict which original looked like this:

$$\langle statement\rangle \ \vDash \ \text{identifier} := \langle expression\rangle \ ; \ |$$

identifier **(** ⟨arg_list⟩ **)** ; |

**begin** ⟨statement_block⟩ **end** |

**if** ⟨condition⟩ **begin** ⟨statement_block⟩ **end**

**else begin** ⟨statement_block⟩ **end** |

**while** ⟨condition⟩ **begin** ⟨statement_block⟩ **end** |

**skip** ;

Here I made a new non terminal simple statement which satisfied statement. The resulting code looked like this:

```
246   void statement (): {}
247   {
248           simple_statement()
249         | <BEGIN> statementBlock() <END>
250         | <IF> condition() <BEGIN> statementBlock()<END>
251           <ELSE> <BEGIN> statementBlock() <END>
252         | <WHILE> condition() <BEGIN> statementBlock() <END>
253         | <SKI> <SCOLON>
254
255   }
256
257   void simple_statement() : {}
258   {
259     <IDENTIFIER>(<ASIGN> expression() <SCOLON> | <OPBRA>argList()<CLBRA> <SCOLON>)
260   }
261
```

Where I also changed the non terminal nemp_argument_list()  which uses argList(), this can be seen here:

```
323   void argList () : {}
324   {
325       nemp_argument_list() | {}
326   }
327
328   void nemp_argument_list() : {}
329   {
330       <IDENTIFIER>(<COMMA>nemp_argument_list() |{} )
331   }
332
```

I also had to do this for nemp_parameter_list(), creating a new non terminal simple_nemp_parameter_list() as seen below:

```
215    void paramList () : {}
216    {
217      nemp_parameter_list() | {}
218    }
219
220    void nemp_parameter_list() : {}
221    {
222      simple_nemp_param_list() (<COMMA> nemp_parameter_list() |{} )
223
224    }
225
226    void simple_nemp_param_list(): {}
227    {
228      <IDENTIFIER><COLON>type()
229
230    }
231
```

I also split up the function non terminal. The reason behind this was not only that I would understand what was happening more and also a better approach for ease of reading the code. I also thought that this would be better for assignment two. I did this by:

```
186    void function() : {}
187    {
188      type() <IDENTIFIER><OPBRA>paramList()<CLBRA>
189      functionBody()
190    }
191
192    void functionBody() : {}
193    {
194      <IS>
195      declartionList()
196      <BEGIN>
197      statementBlock()
198      returnStatement()
199      <END>
200    }
201
202    void returnStatement() : {}
203    {
204      <RET><OPBRA>(expression() |{} )<CLBRA><SCOLON>
205    }
206
```

## PARSING THE EXAMPLE CAL CODE

After I completed all the steps above I tried to compile the example CAL code tests that were given in the language definition file. The first two files parsed correctly, but then I had many unsuccessful parsing errors. Many were due to syntax errors in my code and

spelling corrections. I made good use of the " javacc -debug_parser " option on the command line. This helped me pinpoint where I was going wrong and the slight changes I had to make. In most cases I had forgotten to add a specific rule. Running the debug parser looks like this after passing in a file.

```
[Nialls-MacBook-Pro:CompilerAssignmentOne niall$ java CALParser test2.cal
CALParser: Reading the file test2.cal ...
Call:   programme
  Call:   declartionList
  Return: declartionList
  Call:   functionList
  Return: functionList
  Call:   mainMethod
    Consumed token: <"main" at line 1 column 1>
    Consumed token: <"begin" at line 2 column 1>
    Call:   declartionList
    Return: declartionList
    Call:   statementBlock
    Return: statementBlock
    Consumed token: <"end" at line 5 column 1>
  Return: mainMethod
Return: programme
CALParser: The input was read successfully.
Nialls-MacBook-Pro:CompilerAssignmentOne niall$ ▊
```

## HOW TO RUN MY PARSER

In order to compile this parser, you will need to run the following commands:

- Javacc assignment1.jj
- Javac *.java

To run my parser you will then need to run the following command:

- Java CALParser inputfile