



ECOLE
POLYTECHNIQUE
DE BRUXELLES

INFO-F403 | INTRODUCTION TO LANGUAGE THEORY AND COMPILING

Project Report

Part 1

Darius COUCHARD
Stefano DONNE

27th October 2020

1 Introduction

The project consists in the design and writing of a compiler for FORTR-S. This report aims on the work done on the first part of this project : the scanner, and thus will describes how the lexical analyzer was implemented.

2 Regular Expressions

Different regular expressions were implemented, in respect of the JFlex syntax, to allow pattern matching of the lexical units by the lexical analyzer.

-Variables Name = $[a - z][a - z0 - 9]^*$, starts with a lowercase character and can be followed by any combination of lowercase character or digits.

-Program Name = $[A - Z][A - Z a - z0 - 9]^*$, starts with a uppercase character and can be followed by any combination of lowercase/uppercase character or digits.

-Wrong Numbers = $[0][0 - 9]^+$, matches any combination of digits starting with 0.

-Numbers = $[1 - 9][0 - 9]^* | [0]$, matches any combination of digits starting by something else than 0, or 0 itself.

The idea for matching numbers is straightforward. Indeed, the Wrong Number RegEx has a higher priority in the lexical analyzer than the Number RegEx. By doing that we insure that an error is thrown immediately after passing by a wrong number pattern and that no partial matching of a wrong number is done.

-LineTerminator = $\backslash r | \backslash n | \backslash r \backslash n$, corresponds to carriage return, newline or the two. On Window *Enter* uses carriage return and newline, while on Unix only newline.

-InputCharacter = $[\wedge \backslash r \backslash n]$, corresponds to any pattern which doesn't include carriage return or newline.

3 Comments handling

Two types of comment exists :

- **Inline comment**, starting with `//`, until the end of line.
- **Block comment** , starting with `/*` and ending with `*/`.

The inline comments are handled by simply ignoring `//` and any **InputCharacter** following.

For the block comments it is a bit more tricky because we have to throw an error in case of comment nesting.

This is why another meta state called `BLOCK_COMMENT` was defined for managing this situation. This state is set with JFlex's `yybegin()` method when the **BlockCommentStart** = `/*` regular expression is matched in `YYINITIAL` state.

In this state only two behaviors are possible :

- **BlockCommentStart** RegEx is matched again, thus throwing an error because nested comments aren't supported.
- **BlockCommentEnd** = `*/` RegEx is matched, making the lexical analyze return to `YYINITIAL` state.

Nested comments are more difficult to handle because it's needed to keep track of the several start and end patterns that already matched. A solution for supporting them would be to implement a counter in the `BLOCK_COMMENT` state which is incremented each time **BlockCommentStart** is matched and decremented everytime **BlockCommentEnd** is found. When the counter is at 0 it returns at `YYINITIAL` state.

4 Example

Here's a code example and the output given by the lexical analyzer.

```
1 BEGINPROG Test
2
3  /* Test */
4
5  //this is a test
6  a := 0
7  b := 1
8  READ(input)
9  WHILE b > a DO //TEST
10     c := b
11     b := b + input
12     IF (b > c) THEN
13         /* test
14
15         */
16         PRINT(b)
17     ELSE
18         // test
19         input := input + 1
20     ENDIF
21 ENDWHILE
22 ENDPROG
```

Test.fs

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

41	token: c	lexical unit: VARNAME
42	token:)	lexical unit: RPAREN
43	token: THEN	lexical unit: THEN
44	token: \n	lexical unit: ENDLINE
45	token: \n	lexical unit: ENDLINE
46	token: PRINT	lexical unit: PRINT
47	token: (lexical unit: LPAREN
48	token: b	lexical unit: VARNAME
49	token:)	lexical unit: RPAREN
50	token: \n	lexical unit: ENDLINE
51	token: ELSE	lexical unit: ELSE
52	token: \n	lexical unit: ENDLINE
53	token: \n	lexical unit: ENDLINE
54	token: input	lexical unit: VARNAME
55	token: :=	lexical unit: ASSIGN
56	token: input	lexical unit: VARNAME
57	token: +	lexical unit: PLUS
58	token: 1	lexical unit: NUMBER
59	token: \n	lexical unit: ENDLINE
60	token: ENDIF	lexical unit: ENDIF
61	token: \n	lexical unit: ENDLINE
62	token: ENDWHILE	lexical unit: ENDWHILE
63	token: \n	lexical unit: ENDLINE
64	token: ENDPROG	lexical unit: ENDPROG
65	token: \n	lexical unit: ENDLINE
66		
67	Variables	
68	a 6	
69	b 7	
70	c 10	
71	input 8	

Test.output