

## Coursework 2 (Strand 1)

This coursework is worth 5% and is due on 6 November at 16:00. You are asked to implement the Sulzmann & Lu tokeniser for the WHILE language. You can do the implementation in any programming language you like, but you need to submit the source code with which you answered the questions, otherwise a mark of 0% will be awarded. You can submit your answers in a txt-file or as pdf.

### Disclaimer

It should be understood that the work you submit represents your own effort. You have not copied from anyone else. An exception is the Scala code from KEATS and the code I showed during the lectures, which you can both use. You can also use your own code from the CW 1.

### Question 1 (marked with 1%)

To implement a tokeniser for the WHILE language, you first need to design the appropriate regular expressions for the following eight syntactic entities:

1. keywords are

while, if, then, else, do, for, to, true, false, read, write,  
skip

2. operators are

+, -, \*, %, /, ==, !=, >, <, :=, &&, ||

3. strings are enclosed by "..."

4. parentheses are (, {, } and }

5. there are semicolons ;

6. whitespaces are either " " (one or more) or \n

7. identifiers are letters followed by underscores `_`, letters or digits

8. numbers are `0, 1, ...` and so on; give a regular expression that can recognise `0`, but not numbers with leading zeroes, such as `001`

You can use the basic regular expressions

$$\emptyset, \epsilon, c, r_1 + r_2, r_1 \cdot r_2, r^*$$

but also the following extended regular expressions

$[c_1c_2 \dots c_n]$	a range of characters
$r^+$	one or more times $r$
$r^?$	optional $r$
$r\{n\}$	n-times $r$

Try to design your regular expressions to be as small as possible. For example you should use character ranges for identifiers and numbers.

### Question 2 (marked with 3%)

Implement the Sulzmann & Lu tokeniser from the lectures. For this you need to implement the functions *nullable* and *der* (you can use your code from CW 1), as well as *mkeys* and *inj*. These functions need to be appropriately extended for the extended regular expressions from Q1. Write down the clauses for

$mkeys([c_1c_2 \dots c_n])$	$\stackrel{\text{def}}{=} ?$
$mkeys(r^+)$	$\stackrel{\text{def}}{=} ?$
$mkeys(r^?)$	$\stackrel{\text{def}}{=} ?$
$mkeys(r\{n\})$	$\stackrel{\text{def}}{=} ?$
$inj([c_1c_2 \dots c_n]) c \dots$	$\stackrel{\text{def}}{=} ?$
$inj(r^+) c \dots$	$\stackrel{\text{def}}{=} ?$
$inj(r^?) c \dots$	$\stackrel{\text{def}}{=} ?$
$inj(r\{n\}) c \dots$	$\stackrel{\text{def}}{=} ?$

where *inj* takes three arguments: a regular expression, a character and a value. Test your lexer code with at least the two small examples below:

regex:	string:
$a\{3\}$	<i>aaa</i>
$(a + \epsilon)\{3\}$	<i>aa</i>

Both strings should be successfully lexed by the respective regular expression, that means the lexer returns in both examples a value.

Also add the record regular expression from the lectures to your tokeniser and implement a function, say *env*, that returns all assignments from a value (such that you can extract easily the tokens from a value).

Finally give the tokens for your regular expressions from Q1 and the string

`"read n;"`

and use your *env* function to give the token sequence.

**Question 3 (marked with 1%)**

Extend your tokenizer from Q2 to also simplify regular expressions after each derivation step and rectify the computed values after each injection. Use this tokenizer to tokenize the programs in Figure ?? and ?. Give the tokens of these programs where whitespaces are filtered out.

```

1  write "Fib";
2  read n;
3  minus1 := 0;
4  minus2 := 1;
5  while n > 0 do {
6      temp := minus2;
7      minus2 := minus1 + minus2;
8      minus1 := temp;
9      n := n - 1
10 };
11 write "Result";
12 write minus2

```

Figure 1: Fibonacci program in the WHILE language.

```

1  start := 1000;
2  x := start;
3  y := start;
4  z := start;
5  while 0 < x do {
6      while 0 < y do {
7          while 0 < z do { z := z - 1 };
8          z := start;
9          y := y - 1
10 };
11 y := start;
12 x := x - 1
13 }

```

Figure 2: The three-nested-loops program in the WHILE language. Usually used for timing measurements.