

Handout 5

Whenever you want to design a programming language or implement a compiler for an existing language, the first task is to fix the basic “words” of the language, like what are the keywords or reserved words of the language, what are permitted identifiers, numbers and so on. One convenient way to do this is, of course, to use regular expressions. In this course we want to take a closer look at the WHILE-language. This is a simple imperative language consisting of arithmetic expressions, assignments and loops only. For example the Fibonacci program can be written in this language as follows

```
1  write "Input a number ";
2  read n;
3  x := 0;
4  y := 1;
5  while n > 0 do {
6    temp := y;
7    y := x + y;
8    x := temp;
9    n := n - 1
10 };
11 write "Result ";
12 write y
```

The keywords in this language will be

while, if, then, else, write, read

In addition we will have some typical operators, like $<$, $>$, $:=$ and so on; numbers and strings (which we however ignore for the moment). We also need to specify what the “white space” is in our programming language as well as what comments should look like. We might specify the regular expressions for our language roughly as follows

| | | |
|-------------------|------|--|
| <i>KEYWORD</i> | $:=$ | while + if + then + else + ... |
| <i>IDENT</i> | $:=$ | <i>LETTER</i> · (<i>LETTER</i> + <i>DIGIT</i> + $_$)* |
| <i>OP</i> | $:=$ | := + < + ... |
| <i>NUM</i> | $:=$ | <i>DIGIT</i> ⁺ |
| <i>WHITESPACE</i> | $:=$ | " " + \n |

with the usual meanings for the regular expressions *LETTER* and *DIGIT*.

The problem we have to solve is given a string of our programming language, which regular expression matches which part of the string. For example the input string

```
if true then x + 2 else x + 3
```

needs to be recognised as

```
i f _ t r u e _ t h e n _ x + 2 _ e l s e _ x + 3
```

Since `if` matches the *KEYWORD* regular expression, `_` is a whitespace and so on. This process of separating an input string into components is often called *lexing* or *scanning*. It is usually the first phase of a compiler. Note that the separation into words cannot, in general, be done by looking at whitespaces: while `if` and `true` are separated by a whitespace, the components in `x+2` are not. Another reason for recognising whitespaces explicitly is that in some languages, for example Python, whitespace matters. However in our small language we will eventually filter out all whitespaces and also comments.