

## Handout 5

Whenever you want to design a new programming language or implement a compiler for an existing language, the first task is to fix the basic “words” of the language. For example what are the keywords, or reserved words, of the language, what are permitted identifiers, numbers, expressions and so on. One convenient way to do this is, of course, by using regular expressions.

In this course we want to take a closer look at the WHILE programming language. This is a simple imperative programming language consisting of arithmetic expressions, assignments, if-statements and loops. 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 common operators, such as  $<$ ,  $>$ ,  $:=$  and so on, as well as numbers and strings (which we however ignore for the moment). We also need to specify what the “whitespace” is in our programming language and what comments should look like. As a first try, we might specify the regular expressions for our language roughly as follows

<i>LETTER</i>	$:=$	$a + A + b + B + \dots$
<i>DIGIT</i>	$:=$	$0 + 1 + 2 + \dots$
<i>KEYWORD</i>	$:=$	$\text{while} + \text{if} + \text{then} + \text{else} + \dots$
<i>IDENT</i>	$:=$	$LETTER \cdot (LETTER + DIGIT + \_)*$
<i>OP</i>	$:=$	$:= + < + \dots$
<i>NUM</i>	$:=$	$DIGIT^+$
<i>WHITESPACE</i>	$:=$	$(\text{ } + \backslash n)^+$

Having the regular expressions in place, the problem we have to solve is: given a string of our programming language, which regular expression matches which part of the string. By solving this problem, we can split up a string of our language into components. For example given the input string

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

we expect it will be split up as follows

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

This process of splitting 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 just looking at whitespaces: while **if** and **true** are separated by a whitespace, this is not always the case. As can be seen the three components in **x+2** are not separated by any whitespace. Another reason for recognising whitespaces explicitly is that in some languages, for example Python, whitespace matters. However in our small language we will eventually just filter out all whitespaces and also all comments.

Lexing will not just separate a string into its components, but also classify the components, that is explicitly record that **if** is a keyword, `_` a whitespace, **true** an identifier and so on. For the moment, though, we will only focus on the simpler problem of just splitting a string into components.

There are a few subtleties we need to consider first. For example, say the string is

```
i f f o o _ ...
```

then there are two possibilities for how it can be split up: either we regard the input as the keyword **if** followed by the identifier **foo** (both regular expressions match) or we regard **iffoo** as a single identifier. The choice that is often made in lexers is to look for the longest possible match. This leaves **iffoo** as the only match (since it is longer than **if**).

Unfortunately, the convention about the longest match does not yet make the whole process of lexing completely deterministic. Consider the string

```
t h e n _ ...
```

Clearly, this string should be identified as a keyword. The problem is that also the regular expression *IDENT* for identifiers matches this string. To overcome this ambiguity we need to rank our regular expressions. In our running example we just use the ranking

$$KEYWORD < IDENT < OP < \dots$$

and so on. So even if both regular expressions match in the example above, we can give the regular expression for **??** as follows

Let us see how our algorithm for lexing works in detail. The regular expressions and their ranking are shown above. For our algorithm it will be helpful to have a look at the function *zeroable* defined as follows:

$zeroable(\emptyset)$	$\stackrel{\text{def}}{=} true$
$zeroable(\epsilon)$	$\stackrel{\text{def}}{=} false$
$zeroable(c)$	$\stackrel{\text{def}}{=} false$
$zeroable(r_1 + r_2)$	$\stackrel{\text{def}}{=} zeroable(r_1) \wedge zeroable(r_2)$
$zeroable(r_1 \cdot r_2)$	$\stackrel{\text{def}}{=} zeroable(r_1) \vee zeroable(r_2)$
$zeroable(r^*)$	$\stackrel{\text{def}}{=} false$

In contrast to the function  $nullable(r)$ , which test whether a regular expression can match the empty string, the  $zeroable$  function identifies whether a regular expression cannot match anything at all. The mathematical way of stating this is

$$s \in zeroable(s) \text{ implies } L(r) = \emptyset$$

Let us fix a set of regular expressions  $rs$ . The crucial idea of the algorithm is to take the input string, say

$c_1 c_2 c_3 c_4 \dots$

and build the derivative of all regular expressions in  $rs$  with respect to the first character  $c_1$ . Then we take the result and continue with  $c_2$  until we have either exhausted our input string or all of the regular expressions are “zeroable”.