

## Handout 5 (Lexing)

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 "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
```

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>	$:=$	$\text{LETTER} \cdot (\text{LETTER} + \text{DIGIT} + \_)^*$
<i>OP</i>	$:=$	$:= + < + \dots$
<i>NUM</i>	$:=$	$\text{DIGIT}^+$
<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 is 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 up 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 in the example above, 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, whitespaces matters, that is carry meaning. However in our small language we will eventually just filter out all whitespaces and also all comments.

Lexing not just separates a string into its components, but also classifies the components, meaning it explicitly records 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 up 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 in this case (since it is longer than `if`).

Unfortunately, the convention about the longest match does not yet make the 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* < ...

So even if both regular expressions match in the example above, we give preference to the regular expression for keywords.

Let us see how our algorithm for lexing works in detail. In addition to the functions *nullable* and *der*, it will use the function *zeroable* defined as follows:

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

Recall that the function  $\text{nullable}(r)$  tests whether a regular expression can match the empty string. The function  $\text{zeroable}$ , on the other hand, tests whether a regular expression cannot match anything at all. The mathematical way of stating this property is

$$\text{zeroable}(r) \text{ if and only if } L(r) = \emptyset$$

For what follows let us fix a set of regular expressions  $rs$  as being

*KEYWORD, IDENT, WHITESPACE*

specifying the “words” of our programming language. The algorithm takes as input the  $rs$  and a string, say

c<sub>1</sub>c<sub>2</sub>c<sub>3</sub>c<sub>4</sub> ...

and tries to chop off one word from the beginning of the string. If none of the regular expression in  $rs$  matches, we will just return the empty string.

The crucial idea in the algorithm is to build the derivatives of all regular expressions in  $rs$  with respect to the first character  $c_1$ . Then we take the results and continue with building the derivatives with respect to  $c_2$  until we have either exhausted our input string or all of the regular expressions are “zeroable”. Suppose the input string is

i f 2 ...

then building the derivatives with respect to  $i$  gives

$\text{der } i \text{ (KEYWORD)}$	no
$\text{der } i \text{ (IDENT)}$	no
$\text{der } i \text{ (WHITESPACE)}$	yes

We can eliminate *WHITESPACE* as a potential candidate, because no derivative can go from  $\text{zeroable} = \text{yes}$  to no. That leaves the other two regular expressions as potential candidate and we have to consider the next character,  $f$ , from the input string

$\text{der } f \text{ (der } i \text{ (KEYWORD))}$	no
$\text{der } f \text{ (der } i \text{ (IDENT))}$	no

Since both are 'no', we have to continue with 2 from the input string

$der\ 2\ (der\ f\ (der\ i\ (KEYWORD)))$	<i>zeroable</i>
$der\ 2\ (der\ f\ (der\ i\ (IDENT)))$	yes
	no

Although we now know that the beginning is definitely an *IDENT*, we do not yet know how much of the input string should be considered as an *IDENT*. So we still have to continue and consider the next derivative.

$der\ \_ \ (der\ 2\ (der\ f\ (der\ i\ (IDENT))))$	<i>zeroable</i>
	yes

Since the answer is now 'yes' also in this case, we can stop: once all derivatives are zeroable, we know the regular expressions cannot match any more letters from the input. In this case we only have to go back to the derivative that is nullable. In this case it is

$$der\ 2\ (der\ f\ (der\ i\ (IDENT)))$$

which means we recognised an identifier. In case where there is a choice of more than one regular expressions that are nullable, then we choose the one with the highest precedence. You can try out such a case with the input string

t h e n \_ ...

which can both be recognised as a keyword, but also an identifier.

While in the example above the last nullable derivative is the one directly before the derivative turns zeroable, this is not always the case. Imagine, identifiers can be letters, as permuted by the regular expression *IDENT*, but must end with an underscore.

$$NEWIDENT := LETTER \cdot (LETTER + DIGIT + \_)^* \cdot \_$$

If we use *NEWIDENT* with the input string

i f f o o \_ ...

then it will only become *zeroable* after the \_ has been analysed. In this case we have to go back to the first *f* because only

$$der\ f\ (der\ i\ (KEYWORD))$$

is nullable. As a result we recognise successfully the keyword *if* and the remaining string needs to be consumed by other regular expressions or lead to a lexing error.