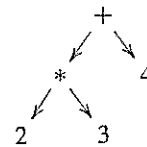# chapter 8

# Functional parsers

In this chapter we show how Haskell can be used to program simple parsers. We start by explaining what parsers are and why they are useful, show how parsers can naturally be viewed as functions, define a number of basic parsers and higher-order functions for combining parsers in various ways, and conclude by developing a parser for arithmetic expressions.

## 8.1 | Parsers

A *parser* is a program that takes a string of characters, and produces some form of tree that makes the syntactic structure of the string explicit. For example, given the string $2 * 3 + 4$, a parser for arithmetic expressions might produce a tree of the following form, in which the numbers appear at the leaves of the tree, and the operators appear at the branches:

$$
\begin{array}{c}
+ \\
\swarrow \quad \searrow \\
* \qquad 4 \\
\swarrow \quad \searrow \\
2 \qquad 3
\end{array}
$$

The structure of this tree makes explicit that $+$ and $*$ are operators with two arguments, and that $*$ has higher priority than $+$.

Parsers are an important topic in computing, because most real-life programs use a parser to pre-process their input. For example, a calculator program parses numeric expressions prior to evaluating them, the Hugs system parses Haskell programs prior to executing them, and a web browser parses hypertext documents prior to displaying them. In each case, making the structure of the input explicit considerably simplifies its further processing. For example, once a numeric expression has been parsed into a tree structure such as in the example above, evaluating the expression is straightforward.

## 8.2 | The parser type

In Haskell, a parser can naturally be viewed directly as a function that takes a string and produces a tree. Hence, given a suitable type *Tree* of trees, the notion of a parser can be represented as a function of type $String \rightarrow Tree$, which we abbreviate as *Parser* using the following declaration:

> **type** $Parser = String \rightarrow Tree$

In general, however, a parser might not always consume its entire argument string. For example, a parser for numbers might be applied to a string comprising a number followed by a word. For this reason, we generalise our type for parsers to also return any unconsumed part of the argument string:

> **type** $Parser = String \rightarrow (Tree, String)$

Similarly, a parser might not always succeed. For example, a parser for numbers might be applied to a string comprising a word. To handle this, we further generalise our type for parsers to return a list of results, with the convention that the empty list denotes failure, and a singleton list denotes success:

> **type** $Parser = String \rightarrow [(Tree, String)]$

Returning a list also opens up the possibility of returning more than one result if the argument string can be parsed in more than one way. For simplicity, however, we only consider parsers that return at most one result.

Finally, different parsers will likely return different kinds of trees, or more generally, any kind of value. For example, a parser for numbers might return an integer value. Hence, it is useful to abstract from the specific type *Tree* of result values, and make this into a parameter of the *Parser* type:

> **type** $Parser\ a = String \rightarrow [(a, String)]$

In summary, this declaration states that a parser of type $a$ is a function that takes an input string and produces a list of results, each of which is a pair comprising a result value of type $a$ and an output string. Alternatively, the parser type can also be read as a rhyme in the style of Dr Seuss!

> *A parser for things*
> *Is a function from strings*
> *To lists of pairs*
> *Of things and strings*

## 8.3 | Basic parsers

We now define three basic parsers that will be used as the building blocks for all other parsers. First of all, the parser *return v* always succeeds with the result value $v$, without consuming any of the input string:

> $return \quad :: \quad a \rightarrow Parser\ a$
> $return\ v \quad = \quad \lambda inp \rightarrow [(v, inp)]$

This function could equally well be defined by $result\ v\ inp = [(v, inp)]$. However, we prefer the above definition in which the second argument $inp$ is shunted to the body of the definition using a lambda expression, because it makes explicit that $return$ is a function that takes a single argument and returns a parser, as expressed by the type $a \rightarrow Parser\ a$.

Whereas $return\ v$ always succeeds, the dual parser $failure$ always fails, regardless of the contents of the input string:

$$
\begin{aligned}
failure &:: & Parser\ a \\
failure &= & \lambda inp \rightarrow [\,]
\end{aligned}
$$

Our final basic parser is $item$, which fails if the input string is empty, and succeeds with the first character as the result value otherwise:

$$
\begin{aligned}
item &:: & Parser\ Char \\
item &= & \lambda inp \rightarrow \textbf{case}\ inp\ \textbf{of} \\
& & \qquad [\,] \rightarrow [\,] \\
& & \qquad (x : xs) \rightarrow [(x, xs)]
\end{aligned}
$$

The **case** mechanism of Haskell used in this definition allows pattern matching to be used in the body of a definition, in this example by matching the string $inp$ against two patterns to choose between two possible results. The **case** mechanism is not used much in this book, but can sometimes be useful.

Because parsers are functions, they could be applied to a string using normal function application, but we prefer to abstract from the representation of parsers by defining our own application function:

$$
\begin{aligned}
parse &:: & Parser\ a \rightarrow String \rightarrow [(a, String)] \\
parse\ p\ inp &= & p\ inp
\end{aligned}
$$

Using $parse$, we conclude this section with some examples that illustrate the behaviour of the three basic parsers defined above:

```
> parse (return 1) "abc"
[(1, "hello")]

> parse failure "abc"
[]

> parse item ""
[]

> parse item "abc"
[('a', "bc")]
```

## 8.4 | Sequencing

Perhaps the simplest way of combining two parsers is to apply one after the other in sequence, with the output string returned by the first parser becoming the input string to the second. But how should the result values from the two

parsers be handled? One approach would be to combine the two values as a pair, using a sequencing operator for parsers with the following type:

$$Parser\ a \to Parser\ b \to Parser\ (a, b)$$

In practice, however, it turns out to be more convenient to combine the sequencing of parsers with the processing of their result values, by means of a sequencing operator $\gg=$ (read as "then") defined as follows:

$$(\gg=) \quad :: \quad Parser\ a \to (a \to Parser\ b) \to Parser\ b$$
$$p \gg= f \quad = \quad \lambda inp \to \mathbf{case}\ parse\ p\ inp\ \mathbf{of}$$
$$[\,] \to [\,]$$
$$[(v, out)] \to parse\ (f\ v)\ out$$

That is, the parser $p \gg= f$ fails if the application of the parser $p$ to the input string fails, and otherwise applies the function $f$ to the result value to give a second parser, which is then applied to the output string to give the final result. In this manner, the result value produced by the first parser is made directly available for processing by the second.

A typical parser built using $\gg=$ has the following structure:

$$p1 \gg= \lambda v1 \to$$
$$p2 \gg= \lambda v2 \to$$
$$\vdots$$
$$pn \gg= \lambda vn \to$$
$$return\ (f\ v1\ v2\ ...\ vn)$$

That is, apply the parser $p1$ and call its result value $v1$; then apply the parser $p2$ and call its result value $v2$; ... ; then apply the parser $pn$ and call its result value $vn$; and, finally, combine all the results into a single value by applying the function $f$. Haskell provides a special syntax for such parsers, allowing them to be expressed in the following, more appealing, form:

$$\mathbf{do}\ v1 \leftarrow p1$$
$$v2 \leftarrow p2$$
$$\vdots$$
$$vn \leftarrow pn$$
$$return\ (f\ v1\ v2\ ...\ vn)$$

As with list comprehensions, the expressions $vi \leftarrow pi$ are called *generators*. If the result value produced by a generator $vi \leftarrow pi$ is not required, the generator can be abbreviated simply by $pi$. Note also that the layout rule applies to the **do** notation for sequencing parsers, in the sense that each parser in the sequence must begin in precisely the same column.

For example, a parser that consumes three characters, discards the second, and returns the first and third as a pair can now be defined as follows:

$$p \quad :: \quad Parser\ (Char,\ Char)$$
$$p \quad = \quad \mathbf{do}\ x \leftarrow item$$
$$item$$
$$y \leftarrow item$$
$$return\ (x, y)$$

Note that $p$ only succeeds if every parser in its defining sequence succeeds, which requires at least three characters in the input string:

```
> parse p "abcdef"
[(('a','c'),"def")]

> parse p "ab"
[]
```

## 8.5 | Choice

Another natural way of combining two parsers is to apply the first parser to the input string, and if this fails to apply the second instead. Such a choice operator +++ (read as "or else") can be defined as follows:

$$
\begin{array}{lll}
(+\!+\!+) & :: & Parser\ a \to Parser\ a \to Parser\ a \\
p +\!+\!+ q & = & \lambda inp \to \textbf{case } parse\ p\ inp\ \textbf{of} \\
& & \qquad\qquad [\,] \to parse\ q\ inp \\
& & \qquad\qquad [(v,\ out)] \to [(v,\ out)]
\end{array}
$$

For example:

```
> parse (item +++ return 'd') "abc"
[('a',"bc")]

> parse (failure +++ return 'd') "abc"
[('d',"abc")]

> parse (failure +++ failure) "abc"
[]
```

## 8.6 | Derived primitives

Using the three basic parsers together with sequencing and choice, we can now define a number of other useful parsing primitives. First of all, we define a parser $sat\ p$ for single characters that satisfy the predicate $p$:

$$
\begin{array}{lll}
sat & :: & (Char \to Bool) \to Parser\ Char \\
sat\ p & = & \textbf{do } x \leftarrow item \\
& & \qquad \textbf{if } p\ x \textbf{ then } return\ x \textbf{ else } failure
\end{array}
$$

Using $sat$ and appropriate predicates from the standard prelude, we can define parsers for single digits, lower-case letters, upper-case letters, arbitrary letters, alphanumeric characters, and specific characters:

$$
\begin{array}{lll}
\textit{digit} & :: & \textit{Parser Char} \\
\textit{digit} & = & \textit{sat isDigit} \\
\textit{lower} & :: & \textit{Parser Char} \\
\textit{lower} & = & \textit{sat isLower} \\
\textit{upper} & :: & \textit{Parser Char} \\
\textit{upper} & = & \textit{sat isUpper} \\
\textit{letter} & :: & \textit{Parser Char} \\
\textit{letter} & = & \textit{sat isAlpha} \\
\textit{alphanum} & :: & \textit{Parser Char} \\
\textit{alphanum} & = & \textit{sat isAlphaNum} \\
\textit{char} & :: & \textit{Char} \to \textit{Parser Char} \\
\textit{char } x & = & \textit{sat } (== x)
\end{array}
$$

For example:

```
> parse digit "123"
[('1',"23")]
```

```
> parse digit "abc"
[]
```

```
> parse (char 'a') "abc"
[('a',"bc")]
```

```
> parse (char 'a') "123"
[]
```

In turn, using *char* we can define a parser *string xs* for the string of characters *xs*, with the string itself returned as the result value:

$$
\begin{array}{lll}
\textit{string} & :: & \textit{String} \to \textit{Parser String} \\
\textit{string } [\,] & = & \textbf{return } [\,] \\
\textit{string } (x : xs) & = & \textbf{do } \textit{char } x \\
& & \quad\quad \textit{string } xs \\
& & \quad\quad \textit{return } (x : xs)
\end{array}
$$

Note that *string* is defined using recursion, and only succeeds if the entire target string is consumed. The base case states that the empty string can always be parsed. The recursive case states that a non-empty string can be parsed by parsing the first character, parsing the remaining characters, and returning the entire string as the result value. For example:

```
> parse (string "abc") "abcdef"
[("abc","def")]
```

```
> parse (string "abc") "ab1234"
[]
```

Our next two parsers, *many p* and *many1 p*, apply a parser *p* as many times as possible until it fails, with the result values from each successful application of *p* being combined as a list. The difference between these two

repetition primitives is that *many* permits zero or more applications of $p$, whereas *many1* requires at least one successful application:

$$
\begin{array}{lcl}
\textit{many} & :: & \textit{Parser } a \rightarrow \textit{Parser } [a] \\
\textit{many } p & = & \textit{many1 } p \mathbin{+\!\!+\!\!+} \textit{return } [\,] \\
\textit{many1} & :: & \textit{Parser } a \rightarrow \textit{Parser } [a] \\
\textit{many1 } p & = & \textbf{do } v \leftarrow p \\
& & \quad\quad vs \leftarrow \textit{many } p \\
& & \quad\quad \textit{return } (v : vs)
\end{array}
$$

Note that *many* and *many1* are defined using mutual recursion. In particular, the definition for *many* $p$ states that $p$ can either be applied at least once or not at all, while the definition for *many1* $p$ states that $p$ can be applied once and then zero or more times. For example:

```
> parse (many digit) "123abc"
[("123", "abc")]


> parse (many digit) "abcdef"
[("", "abcdef")]


> parse (many1 digit) "abcdef"
[]
```

Using *many* and *many1*, we can define parsers for *identifiers* (variable names) comprising a lower-case letter followed by zero or more alphanumeric characters, natural numbers comprising one or more digits, and spacing comprising zero or more space, tab, and newline characters:

$$
\begin{array}{lcl}
\textit{ident} & :: & \textit{Parser String} \\
\textit{ident} & = & \textbf{do } x \leftarrow \textit{lower} \\
& & \quad\quad xs \leftarrow \textit{many alphanum} \\
& & \quad\quad \textit{return } (x : xs) \\
\\
\textit{nat} & :: & \textit{Parser Int} \\
\textit{nat} & = & \textbf{do } xs \leftarrow \textit{many1 digit} \\
& & \quad\quad \textit{return } (\textit{read } xs) \\
\\
\textit{space} & :: & \textit{Parser } () \\
\textit{space} & = & \textbf{do } \textit{many } (\textit{sat isSpace}) \\
& & \quad\quad \textit{return } ()
\end{array}
$$

For example:

```
> parse ident "abc def"
[("abc", " def")]


> parse nat "123 abc"
[(123, " abc")]


> parse space "   abc"
[((), "abc")]
```

Note that *space* returns the empty tuple () as a dummy result value, reflecting the fact that the details of spacing are not usually important.

## 8.7 | Handling spacing

Most real-life parsers allow spacing to be freely used around the basic *tokens* in their input string. For example, the strings 1+2 and 1 + 2 are both parsed in the same way by Hugs. To handle such spacing, we define a new primitive that ignores any space before and after applying a parser for a token:

$$
\begin{aligned}
token \quad &:: \quad Parser\ a \to Parser\ a \\
token\ p \quad &= \quad \textbf{do}\ space \\
&\qquad\quad v \leftarrow p \\
&\qquad\quad space \\
&\qquad\quad return\ v
\end{aligned}
$$

Using *token*, it is now easy to define parsers that ignore spacing around identifiers, natural numbers, and special symbols:

$$
\begin{aligned}
identifier \quad &:: \quad Parser\ String \\
identifier \quad &= \quad token\ ident \\[4pt]
natural \quad &:: \quad Parser\ Int \\
natural \quad &= \quad token\ nat \\[4pt]
symbol \quad &:: \quad String \to Parser\ String \\
symbol\ xs \quad &= \quad token\ (string\ xs)
\end{aligned}
$$

For example, a parser for a non-empty list of natural numbers that ignores spacing around tokens can be defined as follows:

$$
\begin{aligned}
p \quad &:: \quad Parser\ [Int] \\
p \quad &= \quad \textbf{do}\ symbol\ "["\\
&\qquad\quad n \leftarrow natural \\
&\qquad\quad ns \leftarrow many\ (\textbf{do}\ symbol\ ","\\
&\qquad\qquad\qquad\qquad\qquad\qquad natural) \\
&\qquad\quad symbol\ "]"\\
&\qquad\quad return\ (n:ns)
\end{aligned}
$$

This definition states that such a list begins with an opening square bracket and a natural number, followed by zero or more commas and natural numbers, and concludes with a closing square bracket. Note that $p$ only succeeds if a complete list in precisely this format is consumed:

```
> parse p "  [1,  2,  3]  "
[([1,2,3], "")]

> parse p "[1,2,]"
[]
```
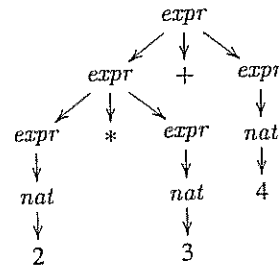
n particular,
t once or not
ed once and

## 8.8 | Arithmetic expressions

We conclude this chapter with an extended example. Consider a simple form of arithmetic expressions built up from natural numbers using addition, multiplication, and parentheses. We assume that addition and multiplication associate to the right, and that multiplication has higher priority than addition. For example, $2 + 3 + 4$ means $2 + (3 + 4)$, while $2 * 3 + 4$ means $(2 * 3) + 4$.
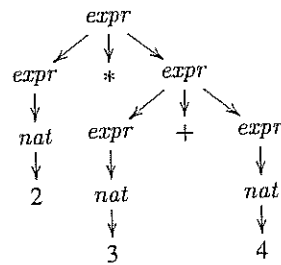
The syntactic structure of a language can be formalised using the mathematical notion of a *grammar*, which is a set of rules that describes how strings of the language can be constructed. For example, a grammar for our language of arithmetic expressions can be defined by the following two rules:

$$expr ::= expr + expr \mid expr * expr \mid (expr) \mid nat$$
$$nat ::= 0 \mid 1 \mid 2 \mid \cdots$$

The first rule states that an expression is either the addition or multiplication of two expressions, a parenthesised expression, or a natural number. In turn, the second rule states that a natural number is either zero, one, two, etc. For example, using this grammar the construction of the expression $2 * 3 + 4$ can be represented by the following *parse tree*, in which the tokens in the expression appear at the leaves, and the grammatical rules applied to construct the expression give rise to the branching structure:



The structure of this tree makes explicit that $2 * 3 + 4$ can be constructed from the addition of two expressions, the first given by the multiplication of two further expressions which are in turn given by the numbers two and three, and the second expression given by the number four. However, the grammar also permits another possible parse tree for this example, which corresponds to the erroneous interpretation of the expression as $2 * (3 + 4)$:
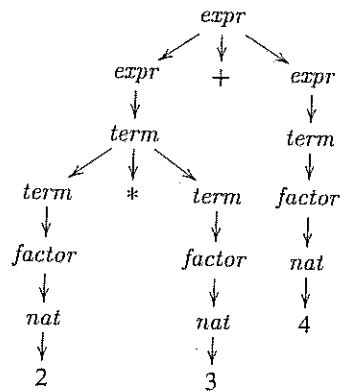


The problem is that our grammar for expressions does not take account of the fact that multiplication has higher priority than addition. However, this can

easily be fixed by modifying the grammar to have a separate rule for each level of priority, with addition at the lowest level of priority, multiplication at the middle level, and parentheses and numbers at the highest level:

$$
\begin{aligned}
expr &::= expr + expr \mid term \\
term &::= term * term \mid factor \\
factor &::= (expr) \mid nat \\
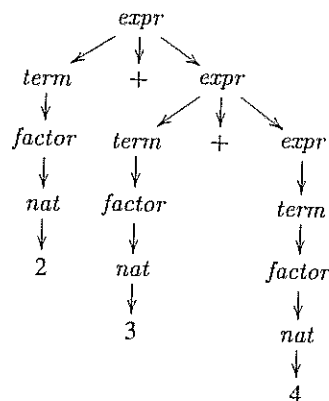nat &::= 0 \mid 1 \mid 2 \mid \cdots
\end{aligned}
$$

Using this new grammar, $2 * 3 + 4$ indeed has a single parse tree, which corresponds to the correct interpretation of the expression as $(2 * 3) + 4$:

```
                        expr
                      ↙  ↓  ↘
                 expr    +    expr
                   ↓           ↓
                 term        term
               ↙  ↓  ↘         ↓
           term  *   term    factor
             ↓         ↓        ↓
          factor    factor    nat
             ↓         ↓        ↓
            nat       nat       4
             ↓         ↓
             2         3
```

We have now dealt with the issue of priority, but our grammar does not yet take account of the fact that addition and multiplication associate to the right. For example, the expression $2 + 3 + 4$ currently has two possible parse trees, corresponding to $(2 + 3) + 4$ and $2 + (3 + 4)$. However, this can also easily be fixed by modifying the grammatical rules for addition and multiplication to be recursive in their right argument only, rather than both arguments:

$$
\begin{aligned}
expr &::= term + expr \mid term \\
term &::= factor * term \mid factor
\end{aligned}
$$

Using these new rules, $2 + 3 + 4$ now has a single parse tree, which corresponds to the correct interpretation of the expression as $2 + (3 + 4)$:

```
                    expr
                  ↙  ↓  ↘
             term    +    expr
               ↓        ↙  ↓  ↘
            factor   term  +   expr
               ↓       ↓         ↓
              nat   factor     term
               ↓       ↓         ↓
               2      nat      factor
                       ↓         ↓
                       3        nat
                                 ↓
                                 4
```

In fact, our grammar for expressions is now unambiguous, in the sense that every well-formed expression has precisely one parse tree.

Our final modification to the grammar is one of simplification. For example, consider the rule $expr ::= term + expr \mid term$, which states that an expression is either the addition of a term and an expression, or is a term. In other words, an expression always begins with a term, which can then be followed by the addition of an expression or by nothing. Hence, the rule for expressions can be simplified to $expr ::= term\ (+ expr \mid \epsilon)$, in which the symbol $\epsilon$ denotes the empty string. Simplifying the rule for terms in a similar manner gives our final grammar for arithmetic expressions:

$$
\begin{array}{lcl}
expr & ::= & term\ (+ expr \mid \epsilon) \\
term & ::= & factor\ (* term \mid \epsilon) \\
factor & ::= & (expr) \mid nat \\
nat & ::= & 0 \mid 1 \mid 2 \mid \cdots
\end{array}
$$

It is now straightforward to translate this grammar into a parser for expressions, by simply rewriting the rules using our parsing primitives. In fact, we choose to have the parser itself evaluate the expression being parsed to its integer value, rather than returning some form of tree:

```
expr  ::  Parser Int
expr  =  do t ← term
             do symbol "+"
                e ← expr
                return (t + e)
             +++ return t


term  ::  Parser Int
term  =  do f ← factor
             do symbol "*"
                t ← term
                return (f * t)
             +++ return f


factor  ::  Parser Int
factor  =  do symbol "("
               e ← expr
               symbol ")"
               return e
             +++ natural
```

For example, the parser $expr$ first parses a term with value $t$, then parses a plus symbol followed by an expression with value $e$ and returns the value $t + e$, or else parses nothing further and simply returns the value $t$. The parsers $term$ and $factor$ can be read in a similar manner.

Finally, using $expr$ we define a function $eval :: String \rightarrow Int$ that evaluates an arithmetic expression to its integer value. To handle the cases of unconsumed and invalid input, we use the library function $error :: String \rightarrow a$ that displays an error message and then terminates the program:

$$eval \quad :: \quad String \rightarrow Int$$
$$eval \; xs \quad = \quad \textbf{case} \; parse \; expr \; xs \; \textbf{of}$$
$$[(n,[\,])] \rightarrow n$$
$$[(\_, out)] \rightarrow error \; (\texttt{"unused input "} +\!\!+ out)$$
$$[\,] \rightarrow error \; \texttt{"invalid input"}$$

For example:

```
>  eval "2*3+4"
10
```

```
>  eval "2*(3+4)"
14
```

```
>  eval "2  *  (3 + 4)"
14
```

```
>  eval "2*3-4"
Error : unused input − 4
```

```
>  eval "-1"
Error : invalid input
```

## 8.9 | Chapter remarks

A library file comprising the parsing primitives from this chapter is available from the book's website. For technical reasons concerning the monadic nature of parsers, a number of the basic definitions in this library are slightly different to those given here. Further details are available in (16; 17), upon which this chapter is based. More information concerning grammars can be found in (27), and more advanced approaches to building parsers in Haskell are given in (22; 9). The reading of the parser type as a rhyme is due to Fritz Ruehr.

## 8.10 | Exercises

1. The library file also defines a parser *int* :: *Parser Int* for an integer. Without looking at this definition, define *int*. Hint: an integer is either a minus symbol followed by a natural number, or a natural number.

2. Define a parser *comment* :: *Parser* () for ordinary Haskell comments that begin with the symbol -- and extend to the end of the current line, which is represented by the control character ' \n' .

3. Using our second grammar for arithmetic expressions, draw the two possible parse trees for the expression $2 + 3 + 4$.

4. Using our third grammar for arithmetic expressions, draw the parse trees for the expressions $2 + 3$, $2 * 3 * 4$ and $(2 + 3) + 4$.

---

*(left margin fragments)*

ie sense that

For example,
n expression
other words,
owed by the
ssions can be
denotes the
ives our final

parser for ex-
tives. In fact,
parsed to its

en parses a plus
value $t + e$, or
ie parsers *term*

it that evaluates
of unconsumed
a that displays

5. Explain why the final simplification of the grammar for arithmetic expressions has a dramatic effect on the efficiency of the resulting parser. Hint: begin by considering how an expression comprising a single number would be parsed if this step had not been made.

6. Extend the parser for arithmetic expressions to support subtraction and division, based upon the following extensions to the grammar:

$$expr \quad ::= \quad term\ (+\ expr\ |\ -\ expr\ |\ \epsilon)$$
$$term \quad ::= \quad factor\ (*\ term\ |\ /\ term\ |\ \epsilon)$$

7. Further extend the grammar and parser for arithmetic expressions to support exponentiation, which is assumed to associate to the right and have higher priority than multiplication and division, but lower priority than parentheses and numbers. For example, $2 \uparrow 3 * 4$ means $(2 \uparrow 3) * 4$. Hint: the new level of priority requires a new rule in the grammar.

8. Consider expressions built up from natural numbers using a subtraction operator that is assumed to associate to the left.

   (a) Define a natural grammar for such expressions.
   (b) Translate this grammar into a parser $expr :: Parser\ Int$.
   (c) What is the problem with this parser?
   (d) Show how it can be fixed. Hint: rewrite the parser using the repetition primitive *many* and the library function *foldl*.