

Monads for functional programming

Philip Wadler, University of Glasgow*

Department of Computing Science, University of Glasgow, G12 8QQ, Scotland
(wadler@dcs.glasgow.ac.uk)

Abstract. The use of monads to structure functional programs is described. Monads provide a convenient framework for simulating effects found in other languages, such as global state, exception handling, output, or non-determinism. Three case studies are looked at in detail: how monads ease the modification of a simple evaluator; how monads act as the basis of a datatype of arrays subject to in-place update; and how monads can be used to build parsers.

1 Introduction

Shall I be pure or impure?

The functional programming community divides into two camps. *Pure* languages, such as Miranda⁰ and Haskell, are lambda calculus pure and simple. *Impure* languages, such as Scheme and Standard ML, augment lambda calculus with a number of possible *effects*, such as assignment, exceptions, or continuations. Pure languages are easier to reason about and may benefit from lazy evaluation, while impure languages offer efficiency benefits and sometimes make possible a more compact mode of expression.

Recent advances in theoretical computing science, notably in the areas of type theory and category theory, have suggested new approaches that may integrate the benefits of the pure and impure schools. These notes describe one, the use of *monads* to integrate impure effects into pure functional languages.

The concept of a monad, which arises from category theory, has been applied by Moggi to structure the denotational semantics of programming languages [13, 14]. The same technique can be applied to structure functional programs [21, 23].

The applications of monads will be illustrated with three case studies. Sec-

tion 2 introduces monads by showing how they can be used to structure a simple evaluator so that it is easy to modify. Section 3 describes the laws satisfied by

* Appears in: J. Jeuring and E. Meijer, editors, *Advanced Functional Programming*, Proceedings of the Båstad Spring School, May 1995, Springer Verlag Lecture Notes in Computer Science 925. A previous version of this note appeared in: M. Broy, editor, *Program Design Calculi*, Proceedings of the Marktoberdorf Summer School, 30 July–8 August 1992. Some errata fixed August 2001; thanks to Dan Friedman for pointing these out.

^o Miranda is a trademark of Research Software Limited.

monads. Section 4 shows how monads provide a new solution to the old problem of providing updatable state in pure functional languages. Section 5 applies monads to the problem of building recursive descent parsers; this is of interest in its own right, and because it provides a paradigm for sequencing and alternation, two of the central concepts of computing.

It is doubtful that the structuring methods presented here would have been discovered without the insight afforded by category theory. But once discovered they are easily expressed without any reference to things categorical. No knowledge of category theory is required to read these notes.

The examples will be given in Haskell, but no knowledge of that is required either. What is required is a passing familiarity with the basics of pure and impure functional programming; for general background see [3, 12]. The languages referred to are Haskell [4], Miranda [20], Standard ML [11], and Scheme [17].

2 Evaluating monads

Pure functional languages have this advantage: all flow of data is made explicit. And this disadvantage: sometimes it is painfully explicit.

A program in a pure functional language is written as a set of equations. Explicit data flow ensures that the value of an expression depends only on its free variables. Hence substitution of equals for equals is always valid, making such programs especially easy to reason about. Explicit data flow also ensures that the order of computation is irrelevant, making such programs susceptible to lazy evaluation.

It is with regard to modularity that explicit data flow becomes both a blessing and a curse. On the one hand, it is the ultimate in modularity. All data in and all data out are rendered manifest and accessible, providing a maximum of flexibility. On the other hand, it is the nadir of modularity. The essence of an algorithm can become buried under the plumbing required to carry data from

its point of creation to its point of use.

Say I write an evaluator in a pure functional language.

- To add error handling to it, I need to modify each recursive call to check for and handle errors appropriately. Had I used an impure language with exceptions, no such restructuring would be needed.
- To add a count of operations performed to it, I need to modify each recursive call to pass around such counts appropriately. Had I used an impure language with a global variable that could be incremented, no such restructuring would be needed.
- To add an execution trace to it, I need to modify each recursive call to pass around such traces appropriately. Had I used an impure language that performed output as a side effect, no such restructuring would be needed.

Or I could use a *monad*.

These notes show how to use monads to structure an evaluator so that the changes mentioned above are simple to make. In each case, all that is required is to redefine the monad and to make a few local changes.

This programming style regains some of the flexibility provided by various features of impure languages. It also may apply when there is no corresponding impure feature. It does not eliminate the tension between the flexibility afforded by explicit data and the brevity afforded by implicit plumbing; but it does ameliorate it to some extent.

The technique applies not just to evaluators, but to a wide range of functional programs. For a number of years, Glasgow has been involved in constructing a compiler for the functional language Haskell. The compiler is itself written in Haskell, and uses the structuring technique described here. Though this paper describes the use of monads in a program tens of lines long, we have experience of using them in a program three orders of magnitude larger.

We begin with the basic evaluator for simple terms, then consider variations that mimic exceptions, state, and output. We analyse these for commonalities, and abstract from these the concept of a monad. We then show how each of the variations fits into the monadic framework.

2.1 Variation zero: The basic evaluator

The evaluator acts on terms, which for purposes of illustration have been taken to be excessively simple.

data *Term* = *Con Int* | *Div Term Term*

A term is either a constant $Con\ a$, where a is an integer, or a quotient, $Div\ t\ u$, where t and u are terms.

The basic evaluator is simplicity itself.

$$\begin{aligned} eval & \quad \quad \quad :: Term \rightarrow Int \\ eval (Con\ a) & = a \\ eval (Div\ t\ u) & = eval\ t \div eval\ u \end{aligned}$$

The function $eval$ takes a term to an integer. If the term is a constant, the constant is returned. If the term is a quotient, its subterms are evaluated and the quotient computed. We use ‘ \div ’ to denote integer division.

The following will provide running examples.

$$\begin{aligned} answer, error & :: Term \\ answer & = (Div (Div (Con\ 1972) (Con\ 2)) (Con\ 23)) \\ error & = (Div (Con\ 1) (Con\ 0)) \end{aligned}$$

Computing $eval\ answer$ yields the value of $((1972 \div 2) \div 23)$, which is 42. The basic evaluator does not incorporate error handling, so the result of $eval\ error$ is undefined.

2.2 Variation one: Exceptions

Say it is desired to add error checking, so that the second example above returns a sensible error message. In an impure language, this is easily achieved with the use of exceptions.

In a pure language, exception handling may be mimicked by introducing a type to represent computations that may raise an exception.

$$\begin{aligned} \mathbf{data}\ M\ a & = Raise\ Exception \mid Return\ a \\ \mathbf{type}\ Exception & = String \end{aligned}$$

A value of type $M\ a$ either has the form $Raise\ e$, where e is an exception, or $Return\ a$, where a is a value of type a . By convention, a will be used both as a type variable, as in $M\ a$, and as a variable ranging over values of that type, as in $Return\ a$.

(A word on the difference between ‘data’ and ‘type’ declarations. A ‘data’ declaration introduces a new data type, in this case M , and new constructors for values of that type, in this case $Raise$ and $Return$. A ‘type’ declaration introduces a new name for an existing type, in this case $Exception$ becomes another name for $String$.)

It is straightforward, but tedious, to adapt the evaluator to this representation.

```
eval          :: Term → M Int
eval (Con a) = Return a
eval (Div t u) = case eval t of
    Raise e → Raise e
    Return a →
        case eval u of
            Raise e → Raise e
            Return b →
                if b = 0
                then Raise "divide by zero"
                else Return (a ÷ b)
```

At each call of the evaluator, the form of the result must be checked: if an exception was raised it is re-raised, and if a value was returned it is processed. Applying the new evaluator to *answer* yields (*Return 42*), while applying it to *error* yields (*Raise "divide by zero"*).

2.3 Variation two: State

Forgetting errors for the moment, say it is desired to count the number of divisions performed during evaluation. In an impure language, this is easily achieved by the use of state. Set a given variable to zero initially, and increment it by one each time a division occurs.

In a pure language, state may be mimicked by introducing a type to represent computations that act on state.

```
type M a = State → (a, State)
type State = Int
```

Now a value of type *M a* is a function that accepts the initial state, and returns the computed value paired with the final state.

Again, it is straightforward but tedious to adapt the evaluator to this representation.

```
eval          :: Term → M Int
eval (Con a) x = (a, x)
eval (Div t u) x = let (a, y) = eval t x in
                    let (b, z) = eval u y in
                    (a ÷ b, z + 1)
```

At each call of the evaluator, the old state must be passed in, and the new state extracted from the result and passed on appropriately. Computing *eval answer 0* yields $(42, 2)$, so with initial state 0 the answer is 42 and the final state is 2 , indicating that two divisions were performed.

2.4 Variation three: Output

Finally, say it is desired to display a trace of execution. In an impure language, this is easily done by inserting output commands at appropriate points.

In a pure language, output may be mimicked by introducing a type to represent computations that generate output.

```

type  $M\ a$     = (Output,  $a$ )
type Output = String

```

Now a value of type $M\ a$ consists of the output generated paired with the value computed.

Yet again, it is straightforward but tedious to adapt the evaluator to this representation.

```

eval           :: Term →  $M\ Int$ 
eval (Con  $a$ ) = (line (Con  $a$ )  $a$ ,  $a$ )
eval (Div  $t\ u$ ) = let ( $x$ ,  $a$ ) = eval  $t$  in
                   let ( $y$ ,  $b$ ) = eval  $u$  in
                   ( $x$  ++  $y$  ++ line (Div  $t\ u$ ) ( $a \div b$ ),  $a \div b$ )

line           :: Term → Int → Output
line  $t\ a$      = "eval (" ++ showterm  $t$  ++ ") ← " ++ showint  $a$  ++ "\n"

```

At each call of the evaluator, the outputs must be collected and assembled to form the output of the enclosing call. The function *line* generates one line of the output. Here *showterm* and *showint* convert terms and integers to strings, ++ concatenates strings, and "\n" represents the string consisting of a newline.

Computing *eval answer* returns the pair $(x, 42)$, where x is the string

```

eval (Con 1972) ← 1972
eval (Con 2)   ← 2
eval (Div (Con 1972) (Con 2)) ← 986
eval (Con 23)  ← 23

```

$eval (Div (Div (Con 1972) (Con 2)) (Con 23)) \Leftarrow 42$

which represents a trace of the computation.

From the discussion so far, it may appear that programs in impure languages are easier to modify than those in pure languages. But sometimes the reverse is true. Say that it was desired to modify the previous program to display the execution trace in the reverse order:

$eval (Div (Div (Con 1972) (Con 2)) (Con 23)) \Leftarrow 42$

$eval (Con 23) \Leftarrow 23$

$eval (Div (Con 1972) (Con 2)) \Leftarrow 986$

$eval (Con 2) \Leftarrow 2$

$eval (Con 1972) \Leftarrow 1972$

This is simplicity itself to achieve with the pure program: just replace the term

$x \# y \# line (Div t u) (a \div b)$

with the term

$line (Div t u) (a \div b) \# y \# x.$

It is not so easy to modify the impure program to achieve this effect. The problem is that output occurs as a side-effect of computation, but one now desires to display the result of computing a term *before* displaying the output generated by that computation. This can be achieved in a variety of ways, but all require substantial modification to the impure program.

2.5 A monadic evaluator

Each of the variations on the interpreter has a similar structure, which may be abstracted to yield the notion of a *monad*.

In each variation, we introduced a type of computations. Respectively, M represented computations that could raise exceptions, act on state, and generate output. By now the reader will have guessed that M stands for *monad*.

The original evaluator has the type $Term \rightarrow Int$, while in each variation its type took the form $Term \rightarrow M Int$. In general, a function of type $a \rightarrow b$ is replaced by a function of type $a \rightarrow M b$. This can be read as a function that accepts an argument of type a and returns a result of type b , with a possible additional effect captured by M . This effect may be to act on state, generate output, raise an exception, or what have you.

What sort of operations are required on the type M ? Examination of the examples reveals two. First, we need a way to turn a value into the computation

that returns that value and does nothing else.

$$\text{unit} :: a \rightarrow M a$$

Second, we need a way to apply a function of type $a \rightarrow M b$ to a computation of type $M a$. It is convenient to write these in an order where the argument comes before the function.

$$(\star) :: M a \rightarrow (a \rightarrow M b) \rightarrow M b$$

A *monad* is a triple (M, unit, \star) consisting of a type constructor M and two operations of the given polymorphic types. These operations must satisfy three laws given in Section 3.

We will often write expressions in the form

$$m \star \lambda a. n$$

where m and n are expressions, and a is a variable. The form $\lambda a. n$ is a lambda expression, with the scope of the bound variable a being the expression n . The above can be read as follows: perform computation m , bind a to the resulting value, and then perform computation n . Types provide a useful guide. From the type of (\star) , we can see that expression m has type $M a$, variable a has type a , expression n has type $M b$, lambda expression $\lambda a. n$ has type $a \rightarrow M b$, and the whole expression has type $M b$.

The above is analogous to the expression

$$\text{let } a = m \text{ in } n.$$

In an impure language, this has the same reading: perform computation m , bind a to the resulting value, then perform computation n and return its value. Here the types say nothing to distinguish values from computations: expression m has type a , variable a has type a , expression n has type b , and the whole expression has type b . The analogy with 'let' explains the choice of the order of arguments to \star . It is convenient for argument m to appear before the function $\lambda a. n$, since computation m is performed before computation n .

The evaluator is easily rewritten in terms of these abstractions.

$$\begin{aligned} \text{eval} & \quad :: \text{Term} \rightarrow M \text{Int} \\ \text{eval} (\text{Con } a) & = \text{unit } a \\ \text{eval} (\text{Div } t \ u) & = \text{eval } t \star \lambda a. \text{eval } u \star \lambda b. \text{unit } (a \div b) \end{aligned}$$

A word on precedence: lambda abstraction binds least tightly and application binds most tightly, so the last equation is equivalent to the following.

$$\text{eval}(\text{Div } t \ u) = (((\text{eval } t) \star (\lambda a. ((\text{eval } u) \star (\lambda b. (\text{unit } (a \div b)))))))$$

The type $\text{Term} \rightarrow M \text{Int}$ indicates that the evaluator takes a term and performs a computation yielding an integer. To compute $(\text{Con } a)$, just return a . To compute $(\text{Div } t \ u)$, first compute t , bind a to the result, then compute u , bind b to the result, and then return $a \div b$.

The new evaluator is a little more complex than the original basic evaluator, but it is much more flexible. Each of the variations given above may be achieved by simply changing the definitions of M , unit , and \star , and by making one or two local modifications. It is no longer necessary to re-write the entire evaluator to achieve these simple changes.

2.6 Variation zero, revisited: The basic evaluator

In the simplest monad, a computation is no different from a value.

```
type M a = a
unit      :: a -> I a
unit a    = a
(★)      :: M a -> (a -> M b) -> M b
a ★ k    = k a
```

This is called the *identity* monad: M is the identity function on types, unit is the identity function, and \star is just application.

Taking M , unit , and \star as above in the monadic evaluator of Section 2.5 and simplifying yields the basic evaluator of Section 2.1

2.7 Variation one, revisited: Exceptions

In the exception monad, a computation may either raise an exception or return a value.

```
data M a      = Raise Exception | Return a
type Exception = String
unit         :: a -> M a
```

<i>unit a</i>	= <i>Return a</i>
(\star)	:: $M a \rightarrow (a \rightarrow M b) \rightarrow M b$
$m \star k$	= case <i>m</i> of <i>Raise e</i> \rightarrow <i>Raise e</i> <i>Return a</i> \rightarrow <i>k a</i>
<i>raise</i>	:: <i>Exception</i> $\rightarrow M a$
<i>raise e</i>	= <i>Raise e</i>

The call *unit a* simply returns the value *a*. The call $m \star k$ examines the result of the computation *m*: if it is an exception it is re-raised, otherwise the function *k* is applied to the value returned. Just as \star in the identity monad is function application, \star in the exception monad may be considered as a form of *strict* function application. Finally, the call *raise e* raises the exception *e*.

To add error handling to the monadic evaluator, take the monad as above. Then just replace *unit (a ÷ b)* by

```

if b = 0
  then raise "divide by zero"
  else unit (a ÷ b)

```

This is commensurate with change required in an impure language.

As one might expect, this evaluator is equivalent to the evaluator with exceptions of Section 2.2. In particular, unfolding the definitions of *unit* and \star in this section and simplifying yields the evaluator of that section.

2.8 Variation two, revisited: State

In the state monad, a computation accepts an initial state and returns a value paired with the final state.

```

type M a = State  $\rightarrow$  (a, State)
type State = Int

unit      :: a  $\rightarrow M a$ 
unit a    =  $\lambda x. (a, x)$ 

```

$$\begin{array}{ll}
(\star) & :: M a \rightarrow (a \rightarrow M b) \rightarrow M b \\
m \star k & = \lambda x. \mathbf{let} (a, y) = m x \mathbf{in} \\
& \quad \mathbf{let} (b, z) = k a y \mathbf{in} \\
& \quad (b, z) \\
tick & :: M () \\
tick & = \lambda x. ((), x + 1)
\end{array}$$

The call *unit a* returns the computation that accept initial state *x* and returns value *a* and final state *x*; that is, it returns *a* and leaves the state unchanged. The call *m * k* performs computation *m* in the initial state *x*, yielding value *a* and intermediate state *y*; then performs computation *k a* in state *y*, yielding value *b* and final state *z*. The call *tick* increments the state, and returns the empty value *()*, whose type is also written *()*.

In an impure language, an operation like *tick* would be represented by a function of type $() \rightarrow ()$. The spurious argument *()* is required to delay the effect until the function is applied, and since the output type is *()* one may guess that the function's purpose lies in a side effect. In contrast, here *tick* has type $M ()$: no spurious argument is needed, and the appearance of *M* explicitly indicates what sort of effect may occur.

To add execution counts to the monadic evaluator, take the monad as above. Then just replace *unit (a ÷ b)* by

$$tick \star \lambda(). unit (a \div b)$$

(Here $\lambda_. e$ is equivalent to $\lambda x. e$ where $x :: ()$ is some fresh variable that does not appear in *e*; it indicates that the value bound by the lambda expression must be *()*.) Again, this is commensurate with change required in an impure language. Simplifying yields the evaluator with state of Section 2.3.

2.9 Variation three, revisited: Output

In the output monad, a computation consists of the output generated paired with the value returned.

$$\mathbf{type} M a = (Output, a)$$

type *Output* = *String*

unit :: $a \rightarrow M a$

unit a = ("", *a*)

(\star) :: $M a \rightarrow (a \rightarrow M b) \rightarrow M b$

$m \star k$ = **let** (x, a) = m **in**
 let (y, b) = $k a$ **in**
 ($x \# y, b$)

out :: *Output* $\rightarrow M ()$

out x = ($x, ()$)

The call *unit a* returns no output paired with *a*. The call $m \star k$ extracts output *x* and value *a* from computation *m*, then extracts output *y* and value *b* from computation *k a*, and returns the output formed by concatenating *x* and *y* paired with the value *b*. The call *out x* returns the computation with output *x* and empty value ().

To add execution traces to the monadic evaluator, take the monad as above. Then in the clause for *Con a* replace *unit a* by

out (line (Con a) a) \star \lambda(). unit a

and in the clause for *Div t u* replace *unit (a \div b)* by

out (line (Div t u) (a \div b)) \star \lambda(). unit (a \div b)

Yet again, this is commensurate with change required in an impure language. Simplifying yields the evaluator with output of Section 2.4.

To get the output in the reverse order, all that is required is to change the definition of \star , replacing $x \# y$ by $y \# x$. This is commensurate with the change required in the pure program, and rather simpler than the change required in an impure program.

You might think that one difference between the pure and impure versions is that the impure version displays output as it computes, while the pure version will display nothing until the entire computation completes. In fact, if the pure language is lazy then output will be displayed in an incremental fashion as the computation occurs. Furthermore, this will also happen if the order of output is reversed, which is much more difficult to arrange in an impure language. Indeed, the easiest way to arrange it may be to simulate lazy evaluation.

3 Monad laws

The operations of a monad satisfy three laws.

- *Left unit.* Compute the value a , bind b to the result, and compute n . The result is the same as n with value a substituted for variable b .

$$\text{unit } a \star \lambda b. n = n[a/b].$$

- *Right unit.* Compute m , bind the result to a , and return a . The result is the same as m .

$$m \star \lambda a. \text{unit } a = m.$$

- *Associative.* Compute m , bind the result to a , compute n , bind the result to b , compute o . The order of parentheses in such a computation is irrelevant.

$$m \star (\lambda a. n \star \lambda b. o) = (m \star \lambda a. n) \star \lambda b. o.$$

The scope of the variable a includes o on the left but excludes o on the right, so this law is valid only when a does not appear free in o .

A binary operation with left and right unit that is associative is called a *monoid*. A monad differs from a monoid in that the right operand involves a binding operation.

To demonstrate the utility of these laws, we prove that addition is associative. Consider a variant of the evaluator based on addition rather than division.

```
data Term = Con Int | Add Term Term
eval      :: Term → M Int
eval (Con a) = unit a
eval (Add t u) = eval t ★ λa. eval u ★ λb. unit (a ÷ b)
```

We show that evaluation of

$$\text{Add } t (\text{Add } u \ v) \quad \text{and} \quad \text{Add } (\text{Add } t \ u) \ v,$$

both compute the same result. Simplify the left term:

$$\begin{aligned} & \text{eval } (\text{Add } t (\text{Add } u \ v)) \\ &= \{ \text{def'n eval} \} \\ & \text{eval } t \star \lambda a. \text{eval } (\text{Add } u \ v) \star \lambda x. \text{unit } (a + x) \end{aligned}$$

$$\begin{aligned}
&= \{ \text{def'n } \textit{eval} \} \\
&\quad \textit{eval } t \star \lambda a. (\textit{eval } u \star \lambda b. \textit{eval } v \star \lambda c. \textit{unit } (b + c)) \star \lambda x. \textit{unit } (a + x) \\
&= \{ \text{associative} \} \\
&\quad \textit{eval } t \star \lambda a. \textit{eval } u \star \lambda b. \textit{eval } v \star \lambda c. \textit{unit } (b + c) \star \lambda x. \textit{unit } (a + x) \\
&= \{ \text{left unit} \} \\
&\quad \textit{eval } t \star \lambda a. \textit{eval } u \star \lambda b. \textit{eval } v \star \lambda c. \textit{unit } (a + (b + c))
\end{aligned}$$

Simplify the right term similarly:

$$\begin{aligned}
&\quad \textit{eval } (\textit{Add } (\textit{Add } t \ u) \ v) \\
&= \{ \text{as before} \} \\
&\quad \textit{eval } t \star \lambda a. \textit{eval } u \star \lambda b. \textit{eval } v \star \lambda c. \textit{unit } ((a + b) + c)
\end{aligned}$$

The result follows by the associativity of addition. This proof is trivial; without the monad laws, it would be impossible.

The proof works in any monad: exception, state, output. This assumes that the code is as above: if it is modified then the proof also must be modified. Section 2.3 modified the program by adding calls to *tick*. In this case, associativity still holds, as can be demonstrated using the law

$$\textit{tick} \star \lambda(). m = m \star \lambda(). \textit{tick}$$

which holds so long as *tick* is the only action on state within *m*. Section 2.4 modified the program by adding calls to *line*. In this case, the addition is no longer associative, in the sense that changing parentheses will change the trace, though the computations will still yield the same value.

As another example, note that for each monad we can define the following operations.

$$\begin{aligned}
\textit{map} &\quad :: (a \rightarrow b) \rightarrow (M \ a \rightarrow M \ b) \\
\textit{map } f \ m &= m \star \lambda a. \textit{unit } (f \ a) \\
\textit{join} &\quad :: M \ (M \ a) \rightarrow M \ a \\
\textit{join } z &= z \star \lambda m. m
\end{aligned}$$

The *map* operation simply applies a function to the result yielded by a computation. To compute *map f m*, first compute *m*, bind *a* to the result, and then return *f a*. The *join* operation is trickier. Let *z* be a computation that *itself* yields a computation. To compute *join z*, first compute *z*, binds *m* to the result, and then behaves as computation *m*. Thus, *join* flattens a mind-boggling double layer of computation into a run-of-the-mill single layer of computation. As we will see in Section 5.1, lists form a monad, and for this monad *map* applies a function to each element of a list, and *join* concatenates a list of lists.

Using *id* for the identity function ($id\ x = x$), and (\cdot) for function composition ($(f \cdot g)\ x = f(g\ x)$), one can then formulate a number of laws.

$$\begin{aligned}map\ id &= id \\map\ (f \cdot g) &= map\ f \cdot map\ g \\map\ f \cdot unit &= unit \cdot f \\map\ f \cdot join &= join \cdot map\ (map\ f) \\join \cdot unit &= id \\join \cdot map\ unit &= id \\join \cdot map\ join &= join \cdot join \\m \star k &= join\ (map\ k\ m)\end{aligned}$$

The proof of each is a simple consequence of the definitions of *map* and *join* and the three monad laws.

Often, monads are defined not in terms of *unit* and \star , but rather in terms of *unit*, *join*, and *map* [10, 13]. The three monad laws are replaced by the first seven of the eight laws above. If one defines \star by the eighth law, then the three monad laws follow. Hence the two definitions are equivalent.

4 State

Arrays play a central role in computing because they closely match current architectures. Programs are littered with array lookups such as $x[i]$ and array updates such as $x[i] := v$. These operations are popular because array lookup is implemented by a single indexed fetch instruction, and array update by a single indexed store.

It is easy to add arrays to a functional language, and easy to provide efficient array lookup. How to provide efficient array update, on the other hand, is a question with a long history. Monads provide a new answer to this old question.

Another question with a long history is whether it is *desirable* to base programs on array update. Since so much effort has gone into developing algorithms and architectures based on arrays, we will sidestep this debate and simply assume the answer is yes.

There is an important difference between the way monads are used in the previous section and the way monads are used here. The previous section showed monads help to use existing language features more effectively; this section shows how monads can help define new language features. No change to the program-

ming language is required, but the implementation must provide a new abstract data type, perhaps as part of the standard prelude.

Here monads are used to manipulate state internal to the program, but the same techniques can be use to manipulate extenal state: to perform input/output, or to communicate with other programming languages. The Glasgow implementation of Haskell uses a design based on monads to provide input/output and interlanguage working with the imperative programming language C [15]. This design has been adopted for version 1.3 of the Haskell standard.

4.1 Arrays

Let Arr be the type of arrays taking indexes of type Ix and yielding values of type Val . The key operations on this type are

$$\begin{aligned}newarray &:: Val \rightarrow Arr, \\index &:: Ix \rightarrow Arr \rightarrow Val, \\update &:: Ix \rightarrow Val \rightarrow Arr \rightarrow Arr.\end{aligned}$$

The call $newarray\ v$ returns an array with all entries set to v ; the call $index\ i\ x$ returns the value at index i in array x ; and the call $update\ i\ v\ x$ returns an array where index i has value v and the remainder is identical to x . The behaviour of these operations is specified by the laws

$$\begin{aligned}index\ i\ (newarray\ v) &= v, \\index\ i\ (update\ i\ v\ x) &= v, \\index\ i\ (update\ j\ v\ x) &= index\ i\ x, \text{ if } i \neq j.\end{aligned}$$

In practice, these operations would be more complex; one needs a way to specify the index bounds, for instance. But the above suffices to explicate the main points.

The efficient way to implement the update operation is to overwrite the specified entry of the array, but in a pure functional language this is only safe if there are no other pointers to the array extant when the update operation is performed. An array satisfying this property is called *single threaded*, following Schmidt [18].

Consider building an interpreter for a simple imperative language. The abstract syntax for this language is represented by the following data types.

data *Term* = *Var Id* | *Con Int* | *Add Term Term*

data *Comm* = *Asgn Id Term* | *Seq Comm Comm* | *If Term Comm Comm*

data *Prog* = *Prog Comm Term*

Here *Id* : *baastad.tex, v1.1.1.11996/02/2915 : 17 : 01wadlerExp* is an unspecified type of identifiers. A term is a variable, a constant, or the sum of two terms; a command is an assignment, a sequence of two commands, or a conditional; and a program consists of a command followed by a term.

The current state of execution will be modelled by an array where the indexes are identifiers and the corresponding values are integers.

type *State* = *Arr*

type *Ix* = *Id*

type *Val* = *Int*

Here is the interpreter.

eval :: *Term* → *State* → *Int*

eval (*Var i*) *x* = *index i x*

eval (*Con a*) *x* = *a*

eval (*Add t u*) *x* = *eval t x* + *eval u x*

exec :: *Comm* → *State* → *State*

exec (*Asgn i t*) *x* = *update i (eval t x) x*

exec (*Seq c d*) *x* = *exec d (exec c x)*

exec (*If t c d*) *x* = **if** *eval t x* = 0 **then** *exec c x* **else** *exec d x*

elab :: *Prog* → *Int*

elab (*Prog c t*) = *eval t (exec c (newarray 0))*

This closely resembles a denotational semantics. The evaluator for terms takes a term and a state and returns an integer. Evaluation of a variable is implemented by indexing the state. The executor for commands takes a command and the initial state and returns the final state. Assignment is implemented by updating the state. The elaborator for programs takes a program and returns an integer. It executes the command in an initial state where all identifiers map to 0, then evaluates the given expression in the resulting state and returns its value.

The state in this interpreter is single threaded: at any moment of execution there is only one pointer to the state, so it is safe to update the state in place. In order for this to work, the update operation must evaluate the new value

before placing it in the array. Otherwise, the array may contain a closure that itself contains a pointer to the array, violating the single threading property. In semantic terms, one says that *update* is strict in all three of its arguments.

A number of researchers have proposed analyses that determine whether a given functional program uses an array in a single threaded manner, with the intent of incorporating such an analysis into an optimising compiler. Most of these analyses are based on abstract interpretation [1]. Although there has been some success in this area, the analyses tend to be so expensive as to be intractable [2, 7].

Even if such analyses were practicable, their use may be unwise. Optimising update can affect a program's time and space usage by an order of magnitude or more. The programmer must be assured that such an optimisation will occur in order to know that the program will run adequately fast and within the available space. It may be better for the programmer to indicate explicitly that an array should be single threaded, rather than leave it to the vagaries of an optimising compiler.

Again, a number of researchers have proposed techniques for indicating that an array is single threaded. Most of these techniques are based on type systems [6, 19, 22]. This area seems promising, although the complexities of these type systems remain formidable.

The following section presents another way of indicating explicitly the intention that an array be single threaded. Naturally, it is based on monads. The advantage of this method is that it works with existing type systems, using only the idea of an abstract data type.

4.2 Array transformers

The monad of array transformers is simply the monad of state transformers, with the state taken to be an array. The definitions of M , *unit*, \star are as before.

```
type  $M\ a = State \rightarrow (a, State)$ 
type  $State = Arr$ 

unit       $:: a \rightarrow M\ a$ 
unit  $a$     $= \lambda x. (a, x)$ 

 $(\star)$      $:: M\ a \rightarrow (a \rightarrow M\ b) \rightarrow M\ b$ 
 $m\ \star\ k$   $= \lambda x. \text{let } (a, y) = m\ x \text{ in}$ 
```

$$\text{let } (b, z) = k a y \text{ in } \\ (b, z)$$

Previously, our state was an integer and we had an additional operation *tick* acting upon the state. Now our state is an array, and we have additional operations corresponding to array creation, indexing, and update.

$$\begin{aligned} \text{block} &:: \text{Val} \rightarrow M a \rightarrow a \\ \text{block } v m &= \text{let } (a, x) = m (\text{newarray } v) \text{ in } a \\ \text{fetch} &:: Ix \rightarrow M \text{Val} \\ \text{fetch } i &= \lambda x. (\text{index } i x, x) \\ \text{assign} &:: Ix \rightarrow \text{Val} \rightarrow M () \\ \text{assign } i v &= \lambda x. ((), \text{update } i v x) \end{aligned}$$

The call *block v m* creates a new array with all locations initialised to *v*, applies monad *m* to this initial state to yield value *a* and final state *x*, deallocates the array, and returns *a*. The call *fetch i* returns the value at index *i* in the current state, and leaves the state unchanged. The call *assign i v* returns the empty value *()*, and updates the state so that index *i* contains value *v*.

A little thought shows that these operations are indeed single threaded. The only operation that could duplicate the array is *fetch*, but this may be implemented as follows: first fetch the entry at the given index in the array, and then return the pair consisting of this value and the pointer to the array. In semantic terms, *fetch* is strict in the array and the index, but not in the value located at the index, and *assign* is strict in the array and the index, but not the value assigned.

(This differs from the previous section, where in order for the interpreter to be single threaded it was necessary for *update* to be strict in the given value. In this section, as we shall see, this strictness is removed but a spurious sequencing is introduced for evaluation of terms. In the following section, the spurious sequencing is removed, but the strictness will be reintroduced.)

We may now make *M* into an *abstract data type* supporting the five operations

The rewritten interpreter is slightly longer than the previous version, but perhaps slightly easier to read. For instance, execution of *(Seq c d)* can be read: compute the execution of *c*, then compute the execution of *d*, then return nothing. Compare this with the previous version, which has the unnerving property that *exec d* appears to the left of *exec c*.

One drawback of this program is that it introduces too much sequencing. Evaluation of *(Add t u)* can be read: compute the evaluation of *t*, bind *a* to the result, then compute the evaluation of *u*, bind *b* to the result, then return *a + b*. This is unfortunate, in that it imposes a spurious ordering on the evaluation of *t* and *u* that was not present in the original program. The order does not matter because although *eval* depends on the state, it does not change it. To remedy this we will augment the monad of state transformers *M* with a second monad *M'* of state readers.

4.3 Array readers

Recall that the monad of array transformers takes an initial array and returns a value and a final array.

```

type M a = State → (a, State)
type State = Arr
  
```

The corresponding monad of array readers takes an array and returns a value. No array is returned because it is assumed identical to the original array.

```

type M' a = State → a
unit'      :: a → M' a
unit' a    = λx. a

(★')      :: M' a → (a → M' b) → M' b
m ★' k    = λx. let a = m x in k a x

fetch'    :: Ix → M' Val
fetch' i  = λx. index i x
  
```

The call *unit' a* ignores the given state *x* and returns *a*. The call *m ★' k* performs computation *m* in the given state *x*, yielding value *a*, then performs computation *k a* in the same state *x*. Thus, *unit'* discards the state and *★'* duplicates it. The call *fetch' i* returns the value in the given state *x* at index *i*.

Clearly, computations that only read the state are a subset of the computa-

tions that may read and write the state. Hence there should be a way to coerce a computation in monad M' into one in monad M .

$$\begin{aligned} \text{coerce} & \quad :: M' a \rightarrow M a \\ \text{coerce } m & = \lambda x. \mathbf{let} \ a = m x \ \mathbf{in} \ (a, x) \end{aligned}$$

The call $\text{coerce } m$ performs computation m in the initial state x , yielding a , and returns a paired with state x . The function coerce enjoys a number of mathematical properties to be discussed shortly.

Again, these operations maintain single threading if suitably implemented. The definitions of \star' and coerce must both be strict in the intermediate value a . This guarantees that when $\text{coerce } m$ is performed in state x , the computation of $m x$ will reduce to a form a that contains no extant pointers to the state x before the pair (a, x) is returned. Hence there will be only one pointer extant to the state whenever it is updated.

A monad is *commutative* if it satisfies the law

$$m \star \lambda a. n \star \lambda b. o = n \star \lambda b. m \star \lambda a. o.$$

The scope of a includes n on the right and not on the left, so this law is valid only when a does not appear free in n . Similarly, b must not appear free in m . In a commutative monad the order of computation does not matter.

The state reader monad is commutative, while the state transformer monad is not. So no spurious order is imposed on computations in the state reader monad. In particular, the call $m \star' k$ may safely be implemented so that m and $k a$ are computed in parallel. However, the final result must still be strict in a . For instance, with the annotations used in the GRIP processor, \star' could be defined as follows.

$$\begin{aligned} m \star' k & = \lambda x. \mathbf{let} \ a = m x \ \mathbf{in} \\ & \quad \mathbf{let} \ b = k a x \ \mathbf{in} \\ & \quad \mathbf{par} \ a \ (\mathbf{par} \ b \ (\mathbf{seq} \ a \ b)) \end{aligned}$$

The two calls to \mathbf{par} spark parallel computations of a and b , and the call to \mathbf{seq} waits for a to reduce to a non-bottom value before returning b .

These operations may be packaged into two abstract data types, M and M' , supporting the eight operations unit , \star , unit' , \star' , block , assign , fetch' , and coerce . The abstraction guarantees single threading, so assign may be implemented by an in-place update.

The interpreter may be rewritten again.

$$\text{eval} \quad \quad :: \text{Term} \rightarrow M' \text{Int}$$

$$\begin{aligned}
eval (Var i) &= fetch' i \\
eval (Con a) &= unit' a \\
eval (Add t u) &= eval t \star' \lambda a. eval u \star' \lambda b. unit' (a + b) \\
exec &:: Comm \rightarrow M () \\
exec (Asgn i t) &= coerce (eval t) \star \lambda a. assign i a \\
exec (Seq c d) &= exec c \star \lambda(). exec d \star \lambda(). unit () \\
exec (If t c d) &= coerce (eval t) \star \lambda a. \\
&\quad \text{if } a = 0 \text{ then } exec c \text{ else } exec d \\
elab &:: Prog \rightarrow Int \\
elab (Prog c t) &= block 0 (exec c \star \lambda(). coerce (eval t) \star \lambda a. unit a)
\end{aligned}$$

This differs from the previous version in that *eval* is written in terms of M' rather than M , and calls to *coerce* surround the calls of *eval* in the other two functions. The new types make it clear that *eval* depends upon the state but does not alter it, while *exec* may both depend upon and alter the state.

The excessive sequencing of the previous version has been eliminated. In the evaluation of $(Add\ t\ u)$ the two subexpressions may be evaluated in either order or concurrently.

A *monad morphism* from a monad M' to a monad M is a function $h :: M' a \rightarrow M a$ that preserves the monad structure:

$$\begin{aligned}
h (unit' a) &= unit a, \\
h (m \star' \lambda a. n) &= (h m) \star \lambda a. (h n).
\end{aligned}$$

It often happens that one wishes to use a combination of monads to achieve a purpose, and monad morphisms play the key role of converting from one monad to another [9].

In particular, *coerce* is a monad morphism, and it follows immediately from this that the two versions of the interpreter are equivalent.

4.4 Conclusion

How a functional language may provide in-place array update is an old problem. This section has presented a new solution, consisting of two abstract data types with eight operations between them. No change to the programming language is required, other than to provide an implementation of these types, perhaps as part of the standard prelude. The discovery of such a simple solution comes as a surprise, considering the plethora of more elaborate solutions that have been

proposed.

A different way of expressing the same solution, based on continuation passing style, has subsequently been proposed by Hudak [8]. But Hudak's solution was inspired by the monad solution, and the monad solution still appears to have some small advantages [15].

Why was this solution not discovered twenty years ago? One possible reason is that the data types involve higher-order functions in an essential way. The usual axiomatisation of arrays involves only first-order functions, and so perhaps it did not occur to anyone to search for an abstract data type based on higher-order functions. That monads led to the discovery of the solution must count as a point in their favour.

5 Parsers

Parsers are the great success story of theoretical computing. The BNF formalism provides a concise and precise way to describe the syntax of a programming language. Mathematical tests can determine if a BNF grammar is ambiguous or vacuous. Transformations can produce an equivalent grammar that is easier to parse. Compiler-compilers can turn a high-level specification of a grammar into an efficient program.

This section shows how monads provide a simple framework for constructing recursive descent parsers. This is of interest in its own right, and also because the basic structures of parsing – sequencing and alternation – are fundamental to all of computing. It also provides a demonstration of how monads can model backtracking (or angelic non-determinism).

5.1 Lists

Our representation of parsers depends upon lists. Lists are ubiquitous in functional programming, and it is surprising that we have managed to get by so far while barely mentioning them. Actually, they have appeared in disguise, as strings are simply lists of characters.

We review some notation. We write $[a]$ for the type of a list with elements all of type a , and $:$ for 'cons'. Thus $[1, 2, 3] = 1 : 2 : 3 : []$, and both have type $[Int]$. Strings are lists of characters, so *String* and $[Char]$ are equivalent, and "monad" is just an abbreviation for $['m', 'o', 'n', 'a', 'd']$.

It is perhaps not surprising that lists form a monad.

$$\begin{aligned}
\mathit{unit} &:: a \rightarrow [a] \\
\mathit{unit} \ a &= [a] \\
(\star) &:: [a] \rightarrow (a \rightarrow [b]) \rightarrow [b] \\
[] \star k &= [] \\
(a : x) \star k &= k \ a \ ++ \ (x \star k)
\end{aligned}$$

The call $\mathit{unit} \ a$ simply forms the unit list containing a . The call $m \star k$ applies k to each element of the list m , and appends together the resulting lists.

If monads encapsulate effects and lists form a monad, do lists correspond to some effect? Indeed they do, and the effect they correspond to is choice. One can think of a computation of type $[a]$ as offering a choice of values, one for each element of the list. The monadic equivalent of a function of type $a \rightarrow b$ is a function of type $a \rightarrow [b]$. This offers a choice of results for each argument, and hence corresponds to a relation. The operation unit corresponds to the identity relation, which associates each argument only with itself. If $k :: a \rightarrow [b]$ and $h :: b \rightarrow [c]$, then

$$\lambda a. k \ a \star \lambda b. h \ b :: a \rightarrow [c]$$

corresponds to the relational composition of k and h .

The *list comprehension* notation provides a convenient way of manipulating lists. The behaviour is analogous to set comprehensions, except the order is significant. For example,

$$\begin{aligned}
[\mathit{sqr} \ a \mid a \leftarrow [1, 2, 3]] &= [1, 4, 9] \\
[(a, b) \mid a \leftarrow [1, 2], b \leftarrow \text{"list"}] &= [(1, 'l'), (1, 'i'), (1, 's'), (1, 't'), \\
&\quad (2, 'l'), (2, 'i'), (2, 's'), (2, 't')]
\end{aligned}$$

The list comprehension notation translates neatly into monad operations.

$$\begin{aligned}
[t \mid x \leftarrow u] &= u \star \lambda x. \mathit{unit} \ t \\
[t \mid x \leftarrow u, y \leftarrow v] &= u \star \lambda x. v \star \lambda y. \mathit{unit} \ t
\end{aligned}$$

Here t is an expression, x and y are variables (or more generally patterns), and u and v are expressions that evaluate to lists. Connections between comprehensions and monads have been described at length elsewhere [21].

5.2 Representing parsers

Parsers are represented in a way similar to state transformers.

```

type  $M a = State \rightarrow [(a, State)]$ 
type  $State = String$ 

```

That is, the parser for type a takes a state representing a string to be parsed, and returns a *list* of containing the value of type a parsed from the string, and a state representing the remaining string yet to be parsed. The list represents all the alternative possible parses of the input state: it will be empty if there is no parse, have one element if there is one parse, have two elements if there are two different possible parses, and so on.

Consider a simple parser for arithmetic expressions, which returns a tree of the type considered previously.

```

data  $Term = Con Int \mid Div Term Term$ 

```

Say we have a parser for such terms.

```

term ::  $M Term$ 

```

Here are some examples of its use.

```

term "23"           = [(Con 23, "")]
term "23 and more" = [(Con 23, " and more")]
term "not a term"  = []
term "((1972 ÷ 2) ÷ 23)" = [(Div (Div (Con 1972) (Con 2)) (Con 23)), ""]

```

A parser m is *unambiguous* if for every input x the list of possible parses $m x$ is either empty or has exactly one item. For instance, *term* is unambiguous. An ambiguous parser may return a list with two or more alternative parsings.

5.3 Parsing an item

The basic parser returns the first item of the input, and fails if the input is exhausted.

```

item      ::  $M Char$ 
item []   = []
item (a : x) = [(a, x)]

```

Here are two examples.

```

item "" = []
item "monad" = [( 'm', "onad" )]

```

Clearly, *item* is unambiguous.

5.4 Sequencing

To form parsers into a monad, we require operations *unit* and \star .

$$\begin{aligned} \textit{unit} &:: a \rightarrow M a \\ \textit{unit} a x &= [(a, x)] \\ (\star) &:: M a \rightarrow (a \rightarrow M b) \rightarrow M b \\ (m \star k) x &= [(b, z) \mid (a, y) \leftarrow m x, (b, z) \leftarrow k a y] \end{aligned}$$

The parser *unit a* accepts input *x* and yields one parse with value *a* paired with remaining input *x*. The parser *m \star k* takes input *x*; parser *m* is applied to input *x* yielding for each parse a value *a* paired with remaining input *y*; then parser *k a* is applied to input *y*, yielding for each parse a value *b* paired with final remaining output *z*.

Thus, *unit* corresponds to the empty parser, which consumes no input, and \star corresponds to sequencing of parsers.

Two items may be parsed as follows.

$$\begin{aligned} \textit{twoItems} &:: M (\textit{Char}, \textit{Char}) \\ \textit{twoItems} &= \textit{item} \star \lambda a. \textit{item} \star \lambda b. \textit{unit} (a, b) \end{aligned}$$

Here are two examples.

$$\begin{aligned} \textit{twoItems} \textit{ "m"} &= [] \\ \textit{twoItems} \textit{ "monad"} &= [((\textit{ 'm' }, \textit{ 'o' }), \textit{ "nad"})] \end{aligned}$$

The parse succeeds only if there are at least two items in the list.

The three monad laws express that the empty parser is an identity for sequencing, and that sequencing is associative.

$$\begin{aligned} \textit{unit} a \star \lambda b. n &= n[a/b] \\ m \star \lambda a. \textit{unit} a &= m \\ m \star (\lambda a. n \star \lambda b. o) &= (m \star \lambda a. n) \star \lambda b. o \end{aligned}$$

If *m* is unambiguous and *k a* is unambiguous for every *a*, then *m \star k* is also unambiguous.

5.5 Alternation

Parsers may also be combined by alternation.

$$\begin{aligned} \text{zero} &:: M a \\ \text{zero } x &= [] \\ (\oplus) &:: M a \rightarrow M a \rightarrow M a \\ (m \oplus n) x &= m x \uplus n x \end{aligned}$$

The parser *zero* takes input *x* and always fails. The parser $m \oplus n$ takes input *x* and yields all parses of *m* applied to input *x* and all parses of *n* applied to the same input *x*.

Here is a parser that parses one or two items from the input.

$$\begin{aligned} \text{oneOrTwoItems} &:: M \text{String} \\ \text{oneOrTwoItems} &= (\text{item} \star \lambda a. \text{unit } [a]) \\ &\oplus (\text{item} \star \lambda a. \text{item} \star \lambda b. \text{unit } [a, b]) \end{aligned}$$

Here are three examples.

$$\begin{aligned} \text{oneOrTwoItems} \text{ ""} &= [] \\ \text{oneOrTwoItems} \text{ "m"} &= [(\text{"m"}, \text{" "})] \\ \text{oneOrTwoItems} \text{ "monad"} &= [(\text{"m"}, \text{"onad"}), (\text{"mo"}, \text{"nad"})] \end{aligned}$$

The last yields two alternative parses, showing that alternation can yield ambiguous parsers.

The parser that always fails is the identity for alternation, and alternation is associative.

$$\begin{aligned} \text{zero} \oplus n &= n \\ m \oplus \text{zero} &= m \\ m \oplus (n \oplus o) &= (m \oplus n) \oplus o \end{aligned}$$

Furthermore, *zero* is indeed a zero of \star , and \star distributes through \oplus .

$$\begin{aligned} \text{zero} \star k &= \text{zero} \\ m \star \lambda a. \text{zero} &= \text{zero} \\ (m \oplus n) \star k &= (m \star k) \oplus (n \star k) \end{aligned}$$

It is *not* the case that \star distributes rightward through \oplus only because we are representing alternative parses by an ordered list; if we used an unordered bag,

then $m \star \lambda a. (k a \oplus h a) = (m \star k) \oplus (m \star h)$ would also hold. An unambiguous parser yields a list of length at most one, so the order is irrelevant, and hence this law also holds whenever either side is unambiguous (which implies that both sides are).

5.6 Filtering

A parser may be filtered by combining it with a predicate.

$$\begin{aligned} (\triangleright) \quad & :: M a \rightarrow (a \rightarrow Bool) \rightarrow M a \\ m \triangleright p & = m \star \lambda a. \mathbf{if} \ p a \ \mathbf{then} \ \mathit{unit} \ a \ \mathbf{else} \ \mathit{zero} \end{aligned}$$

Given a parser m and a predicate on values p , the parser $m \triangleright p$ applies parser m to yield a value a ; if $p a$ holds it succeeds with value a , otherwise it fails. Note that filtering is written in terms of previously defined operators, and need not refer directly to the state.

Let $\mathit{isLetter}$ and $\mathit{isDigit}$ be the obvious predicates. Here are two parsers.

$$\begin{aligned} \mathit{letter} & :: M Char \\ \mathit{letter} & = \mathit{item} \triangleright \mathit{isLetter} \\ \\ \mathit{digit} & :: M Int \\ \mathit{digit} & = (\mathit{item} \triangleright \mathit{isDigit}) \star \lambda a. \mathit{unit} (\mathit{ord} \ a - \mathit{ord} \ '0') \end{aligned}$$

The first succeeds only if the next input item is a letter, and the second succeeds only if it is a digit. The second also converts the digit to its corresponding value, using $\mathit{ord} :: Char \rightarrow Int$ to convert a character to its ASCII code. Assuming that \triangleright has higher precedence than \star would allow some parentheses to be dropped from the second definition.

A parser for a literal recognises a single specified character.

$$\begin{aligned} \mathit{lit} & :: Char \rightarrow M Char \\ \mathit{lit} \ c & = \mathit{item} \triangleright (\lambda a. a = c) \end{aligned}$$

The parser $\mathit{lit} \ c$ succeeds if the input begins with character c , and fails otherwise.

$$\begin{aligned} \mathit{lit} \ 'm' \ \mathit{"monad"} & = [(\ 'm', \ \mathit{"onad"})] \\ \mathit{lit} \ 'm' \ \mathit{"parse"} & = [] \end{aligned}$$

From the previous laws, it follows that filtering preserves zero and distributes through alternation.

$$\begin{aligned} \mathit{zero} \triangleright p & = \mathit{zero} \\ (m \oplus n) \triangleright p & = (m \triangleright p) \oplus (n \triangleright p) \end{aligned}$$

If m is an unambiguous parser, so is $m \triangleright p$.

5.7 Iteration

A single parser may be iterated, yielding a list of parsed values.

$$\begin{aligned} \textit{iterate} &:: M\ a \rightarrow M\ [a] \\ \textit{iterate}\ m &= (m \star \lambda a. \textit{iterate}\ m \star \lambda x. \textit{unit}\ (a : x)) \\ &\oplus \textit{unit}\ [] \end{aligned}$$

Given a parser m , the parser $\textit{iterate}\ m$ applies parser m in sequence zero or more times, returning a list of all the values parsed. In the list of alternative parses, the longest parse is returned first.

Here is an example.

$$\begin{aligned} \textit{iterate}\ \textit{digit}\ \textit{"23 and more"} &= [([2, 3], \textit{"and more"}), \\ &([2], \textit{"3 and more"}), \\ &([], \textit{"23 and more"})] \end{aligned}$$

Here is one way to parse a number.

$$\begin{aligned} \textit{number} &:: M\ Int \\ \textit{number} &= \textit{digit} \star \lambda a. \textit{iterate}\ \textit{digit} \star \lambda x. \textit{unit}\ (\textit{asNumber}\ (a : x)) \end{aligned}$$

Here $\textit{asNumber}$ takes a list of one or more digits and returns the corresponding number. Here is an example.

$$\begin{aligned} \textit{number}\ \textit{"23 and more"} &= [(23, \textit{"and more"}), \\ &(2, \textit{"3 and more"})] \end{aligned}$$

This supplies two possible parses, one which parses both digits, and one which parses only a single digit. A number is defined to contain at least one digit, so there is no parse with zero digits.

As this last example shows, often it is more natural to design an iterator to yield only the longest possible parse. The next section describes a way to achieve this.

5.8 Biased choice

Alternation, written $m \oplus n$, yields all parses yielded by m followed by all parses yielded by n . For some purposes, it is more sensible to choose one or the other: all parses by m if there are any, and all parses by n otherwise. This is called biased choice.

$$\begin{aligned} (\odot) \quad & :: M a \rightarrow M a \rightarrow M a \\ (m \odot n) x & = \text{if } m x \neq [] \text{ then } m x \text{ else } n x \end{aligned}$$

Biased choice, written $m \odot n$, yields the same parses as m , unless m fails to yield any parse, in which case it yields the same parses as n .

Here is iteration, rewritten with biased choice.

$$\begin{aligned} \text{reiterate} & :: M a \rightarrow M [a] \\ \text{reiterate } m & = (m \star \lambda a. \text{reiterate } m \star \lambda x. \text{unit } (a : x)) \\ & \quad \odot \text{unit } [] \end{aligned}$$

The only difference is to replace \oplus with \odot . Instead of yielding a list of all possible parses with the longest first, this yields only the longest possible parse.

Here is the previous example revisited.

$$\text{reiterate digit "23 and more"} = [[[2, 3], " and more"]]$$

In what follows, *number* is taken to be rewritten with *reiterate*.

$$\begin{aligned} \text{number} & :: M \text{Int} \\ \text{number} & = \text{digit} \star \lambda a. \text{reiterate digit} \star \lambda x. \text{unit } (\text{asNumber } (a : x)) \end{aligned}$$

Here is an example that reveals a little of how ambiguous parsers may be used to search a space of possibilities. We use *reiterate* to find all ways of taking one or two items from a string, zero or more times.

$$\begin{aligned} \text{reiterate oneOrTwoItems "many"} & = [(["m", "a", "n", "y", ""), \\ & \quad (["m", "a", "ny", ""), \\ & \quad (["m", "an", "y", ""), \\ & \quad (["ma", "n", "y", ""), \\ & \quad (["ma", "ny", "")] \end{aligned}$$

This combines alternation (in *oneOrTwoItems*) with biased choice (in *reiterate*).

There are several possible parses, but for each parse *oneOrTwoItems* has been applied until the entire input has been consumed. Although this example is somewhat fanciful, a similar technique could be used to find all ways of breaking a dollar into nickels, dimes, and quarters.

If m and n are unambiguous, then $m \circ n$ and *reiterate* m are also unambiguous. For unambiguous parsers, sequencing distributes right through biased choice:

$$(m \star k) \circ (m \star h) = m \star \lambda a. k a \circ h a$$

whenever m is unambiguous. Unlike with alternation, sequencing does not distribute left through biased choice, even for unambiguous parsers.

5.9 A parser for terms

It is now possible to write the parser for terms alluded to at the beginning. Here is a grammar for fully parenthesised terms, expressed in BNF.

$$\textit{term} ::= \textit{number} \mid (' \textit{term} \textit{'\div'} \textit{term} \textit{'})'$$

This translates directly into our notation as follows. Note that our notation, unlike BNF, specifies exactly how to construct the returned value.

$$\begin{aligned} \textit{term} &:: M \textit{Term} \\ \textit{term} &= (\textit{number} \quad \star \lambda a. \\ &\quad \textit{unit}(\textit{Con} a)) \\ &\oplus (\textit{lit} \textit{'}' \quad \star \lambda_. \\ &\quad \textit{term} \quad \star \lambda t. \\ &\quad \textit{lit} \textit{'\div'} \quad \star \lambda_. \\ &\quad \textit{term} \quad \star \lambda u. \\ &\quad \textit{lit} \textit{'}'}) \\ &\quad \textit{unit}(\textit{Div} t u)) \end{aligned}$$

(Here $\lambda_. e$ is equivalent to $\lambda x. e$ where x is some fresh variable that does not appear in e ; it indicates that the value bound by the lambda expression is not of interest.) Examples of the use of this parser appeared earlier.

The above parser is written with alternation, but as it is unambiguous, it could just as well have been written with biased choice. The same is true for all the parsers in the next section.

5.10 Left recursion

The above parser works only for fully parenthesised terms. If we allow unparenthesised terms, then the operator \div should associate to the left. The usual way to express such a grammar in BNF is as follows.

$$\begin{aligned} \textit{term} & ::= \textit{term} \div \textit{factor} \mid \textit{factor} \\ \textit{factor} & ::= \textit{number} \mid \textit{'term' } \end{aligned}$$

This translates into our notation as follows.

$$\begin{aligned} \textit{term} & :: M \textit{Term} \\ \textit{term} & = (\textit{term} \quad \star \lambda t. \\ & \quad \textit{lit} \div \quad \star \lambda _ \\ & \quad \textit{factor} \quad \star \lambda u. \\ & \quad \textit{unit} (\textit{Div} \textit{t} \textit{u})) \\ & \oplus \textit{factor} \\ \textit{factor} & :: M \textit{Term} \\ \textit{factor} & = (\textit{number} \quad \star \lambda a. \\ & \quad \textit{unit} (\textit{Con} \textit{a})) \\ & \oplus (\textit{lit} \textit{' } \quad \star \lambda _ \\ & \quad \textit{term} \quad \star \lambda t. \\ & \quad \textit{lit} \textit{' } \quad \star \lambda _ \\ & \quad \textit{unit} \textit{t}) \end{aligned}$$

There is no problem with *factor*, but any attempt to apply *term* results in an infinite loop. The problem is that the first step of *term* is to apply *term*, leading to an infinite regress. This is called the *left recursion problem*. It is a difficulty for all recursive descent parsers, functional or otherwise.

The solution is to rewrite the grammar for *term* in the following equivalent form.

$$\begin{aligned} \textit{term} & ::= \textit{factor} \textit{term}' \\ \textit{term}' & ::= \div \textit{factor} \textit{term}' \mid \textit{unit} \end{aligned}$$

where as usual *unit* denotes the empty parser. This then translates directly into our notation.

$$\begin{aligned} \textit{term} & :: M \textit{Term} \\ \textit{term} & = \textit{factor} \star \lambda t. \textit{term}' \textit{t} \\ \textit{term}' & :: \textit{Term} \rightarrow M \textit{Term} \end{aligned}$$

$$\begin{array}{lll}
 \text{term}' t = & (\text{lit } \div) & \star \lambda_. \\
 & \text{factor} & \star \lambda u. \\
 & \text{term}' (\text{Div } t u) & \\
 \oplus & \text{unit } t &
 \end{array}$$

Here *term'* parses the remainder of a term; it takes an argument corresponding to the term parsed so far.

This has the desired effect.

$$\begin{array}{l}
 \text{term } "1972 \div 2 \div 23" = [((\text{Div } (\text{Div } (\text{Con } 1972) (\text{Con } 2)) (\text{Con } 23)), " ")] \\
 \text{term } "1972 \div (2 \div 23)" = [((\text{Div } (\text{Con } 1972) (\text{Div } (\text{Con } 2) (\text{Con } 23))), " ")]
 \end{array}$$

In general, the left-recursive definition

$$m = (m \star k) \oplus n$$

can be rewritten as

$$m = n \star (\text{closure } k)$$

where

$$\begin{array}{ll}
 \text{closure} & :: (a \rightarrow M a) \rightarrow (a \rightarrow M a) \\
 \text{closure } k a & = (k a \star \text{closure } k) \oplus \text{unit } a
 \end{array}$$

Here $m :: M a$, $n :: M a$, and $k :: a \rightarrow M a$.

5.11 Improving laziness

Typically, a program might be represented as a function from a list of characters – the input – to another list of characters – the output. Under lazy evaluation, usually only some of the input need be read before the first part of the output list is produced. This ‘on line’ behavior is essential for some purposes.

In general, it is unreasonable to expect such behaviour from a parser, since in general it cannot be known that the input will be successfully parsed until all of it is read. However, in certain special cases one may hope to do better.

Consider applying *reiterate* m to a string beginning with an instance of m . In this case, the parse cannot fail: regardless of the remainder of the string, one would expect the parse yielded to be a list beginning with the parsed value. Under lazy evaluation, one might expect to be able to generate output corresponding to the first digit before the remaining input has been read.

But this is not what happens: the parser reads the entire input before any output is generated. What is necessary is some way to encode that the parser

reiterate m always succeeds. (Even if the beginning of the input does not match *m*, it will yield as a value the empty list.) This is provided by the function *guarantee*.

$$\begin{aligned} \textit{guarantee} & \quad :: M\ a \rightarrow M\ a \\ \textit{guarantee}\ m\ x & = \textit{let}\ u = m\ x\ \textit{in}\ (\textit{fst}\ (\textit{head}\ u), \textit{snd}\ (\textit{head}\ u)) : \textit{tail}\ u \end{aligned}$$

Here $\textit{fst}\ (a, b) = a$, $\textit{snd}\ (a, b) = b$, $\textit{head}\ (a : x) = a$, and $\textit{tail}\ (a : x) = x$.

Here is *reiterate* with the guarantee added.

$$\begin{aligned} \textit{reiterate} & \quad :: M\ a \rightarrow M\ [a] \\ \textit{reiterate}\ m & = \textit{guarantee}\ (\ (m \star \lambda a. \textit{reiterate}\ m \star \lambda x. \textit{unit}\ (a : x)) \\ & \quad \oslash \textit{unit}\ []) \end{aligned}$$

This ensures that *reiterate m* and all of its recursive calls return a list with at least one answer. As a result, the behaviour under lazy evaluation is much improved.

The preceding explanation is highly operational, and it is worth noting that denotational semantics provides a useful alternative approach. Let \perp denote a program that does not terminate. One can verify that with the old definition

$$\textit{reiterate}\ \textit{digit}\ ('1' : \perp) = \perp$$

while with the new definition

$$\textit{reiterate}\ \textit{digit}\ ('1' : \perp) = (('1' : \perp), \perp) : \perp$$

Thus, given that the input begins with the character '1' but that the remainder of the input is unknown, with the old definition nothing is known about the output, while with the new definition it is known that the output yields at least one parse, the value of which is a list which begins with the character '1'.

Other parsers can also benefit from a judicious use of *guarantee*, and in particular *iterate* can be modified like *reiterate*.

5.12 Conclusion

We have seen that monads provide a useful framework for structuring recursive descent parsers. The empty parser and sequencing correspond directly to *unit* and \star , and the monads laws reflect that sequencing is associative and has the empty parser as a unit. The failing parser and alternation correspond to *zero*

and \oplus , which satisfy laws reflecting that alternation is associative and has the failing parser as a unit, and that sequencing distributes through alternation.

Sequencing and alternation are fundamental not just to parsers but to much of computing. If monads capture sequencing, then it is reasonable to ask: what captures both sequencing and alternation? It may be that *unit*, \star , *zero*, and \oplus , together with the laws above, provide such a structure. Further experiments are needed. One hopeful indication is that a slight variation of the parser monad yields a plausible model of Dijkstra's guarded command language.

References

1. S. Abramsky and C. Hankin, *Abstract Interpretation of Declarative Languages*, Ellis Horwood, 1987.
2. A. Bloss, Update analysis and the efficient implementation of functional aggregates. In *4'th Symposium on Functional Programming Languages and Computer Architecture*, ACM, London, September 1989.
3. R. Bird and P. Wadler, *Introduction to Functional Programming*. Prentice Hall, 1987.
4. P. Hudak, S. Peyton Jones and P. Wadler, editors, *Report on the Programming Language Haskell: Version 1.1*. Technical report, Yale University and Glasgow University, August 1991.
5. J.-Y. Girard, Linear logic. *Theoretical Computer Science*, 50:1–102, 1987.
6. J. Guzmán and P. Hudak, Single-threaded polymorphic lambda calculus. In *IEEE Symposium on Logic in Computer Science*, Philadelphia, June 1990.
7. P. Hudak, A semantic model of reference counting and its abstraction (detailed summary). In *ACM Conference on Lisp and Functional Programming*, pp. 351–363, Cambridge, Massachusetts, August 1986.
8. P. Hudak, Continuation-based mutable abstract data types, or how to have your state and munge it too. Technical report YALEU/DCS/RR-914, Department of Computer Science, Yale University, July 1992.
9. D. King and P. Wadler, Combining monads. In *Glasgow Workshop on Functional Programming*, Ayr, July 1992. Workshops in Computing Series, Springer Verlag, to appear.
10. S. Mac Lane, *Categories for the Working Mathematician*, Springer-Verlag, 1971.
11. R. Milner, M. Tofte, and R. Harper, *The definition of Standard ML*. MIT Press, 1990.
12. L. C. Paulson, *ML for the Working Programmer*. Cambridge University Press, 1991.
13. E. Moggi, Computational lambda-calculus and monads. In *Symposium on Logic in Computer Science*, Asilomar, California; IEEE, June 1989. (A longer version is

- available as a technical report from the University of Edinburgh.)
14. E. Moggi, An abstract view of programming languages. Course notes, University of Edinburgh.
 15. S. L. Peyton Jones and P. Wadler, Imperative functional programming. In *20th Symposium on Principles of Programming Languages*, Charleston, South Carolina; ACM, January 1993.
 16. G. Plotkin, Call-by-name, call-by-value, and the λ -calculus. *Theoretical Computer Science*, 1:125–159, 1975.
 17. J. Rees and W. Clinger (eds.), The revised⁹ report on the algorithmic language Scheme. *ACM SIGPLAN Notices*, 21(12):37–79, 1986.
 18. D. Schmidt, Detecting global variables in denotational specifications. *ACM Trans. on Programming Languages and Systems*, 7:299–310, 1985.
 19. V. Swarup, U. S. Reddy, and E. Ireland, Assignments for applicative languages. In *Conference on Functional Programming Languages and Computer Architecture*, Cambridge, Massachusetts; LNCS 523, Springer Verlag, August 1991.
 20. D. A. Turner, An overview of Miranda. In D. A. Turner, editor, *Research Topics in Functional Programming*. Addison Wesley, 1990.
 21. P. Wadler, Comprehending monads. In *Conference on Lisp and Functional Programming*, Nice, France; ACM, June 1990.
 22. Is there a use for linear logic? *Conference on Partial Evaluation and Semantics-Based Program Manipulation (PEPM)*, New Haven, Connecticut; ACM, June 1991.
 23. P. Wadler, The essence of functional programming (invited talk). In *19th Symposium on Principles of Programming Languages*, Albuquerque, New Mexico; ACM, January 1992.