

Composing monads*

Mark P. Jones
Yale University
New Haven, CT, U.S.A.
jones-mark@cs.yale.edu

Luc Duponcheel
Alcatel Bell Telephone
Antwerp, Belgium
ldup@ra.alcbel.be

Research Report YALEU/DCS/RR-1004, December 1993

Abstract

Monads are becoming an increasingly important tool for functional programming. Different monads can be used to model a wide range of programming language features. However, real programs typically require a combination of different features, so it is important to have techniques for combining several features in a single monad. In practice, it is usually possible to construct a monad that supports some specific combination of features. However, the techniques used are typically ad-hoc and it is very difficult to find general techniques for combining arbitrary monads.

This report gives three general constructions for the composition of monads, each of which depends on the existence of an auxiliary function linking the monad structures of the components. In each case, we establish a set of laws that the auxiliary function must satisfy to ensure that the composition is itself a monad.

Using the notation of constructor classes, we describe some specific applications of these constructions. These results are used in the development of a simple expression evaluator that combines exceptions, output and an environment of variable bindings using a composition of three corresponding monads.

1 Introduction

In recent years, the concept of a monad – an idea that was originally motivated by high-level abstract algebra – has become an important and practical tool for functional programmers. The reason for this is that monads provide a uniform framework for describing a wide range of programming language features including, for example, state, I/O, continuations, exceptions, parsing and non-determinism, without leaving the framework of a purely functional language. Many of these techniques were already familiar to functional programmers, but there are many new insights when they are reinterpreted as specific instances of a more general concept.

*A Gopher script containing executable versions of the programs described in this report is currently available by anonymous ftp from `nebula.cs.yale.edu` in the file `pub/yale-fp/reports/RR-1004.gs`.

Much of the initial interest in monads has been motivated by the work of Wadler [14, 15] who, in turn, drew inspiration from the work of Moggi [9] and Spivey [12]. Monads are already widely used in both small and large programs (for example, the Glasgow Haskell compiler, the largest Haskell program known to us at the time of writing, makes substantial use of monads [2]). New approaches to old problems have been proposed, relying heavily on the use of monads. For example, the I/O monad proposed in [11] is already widely used and may soon be included as part of the definition of Haskell. Monads are even influencing the design of programming languages. For example, monads provide an important motivating example for the system of constructor classes presented in [5].

With the exception of [8], questions about how monads can be combined have, so far, received surprisingly little attention. This is an important topic because many real programs require a combination of features, for example, state, I/O and exceptions. In practice, it is usually possible to construct a suitable monad that supports the desired combination of features, but the methods used are typically ad-hoc and monolithic.

The goal of this report is to investigate some techniques for combining monads by a process of *composition* and to illustrate how these constructions can be used in practice. The conditions required to build a composite monad are quite complex, but we hope that this work will provide a step towards a more modular approach to the use of monads.

One of the nicest features about monads is that, with a little practice (and perhaps, some carefully chosen syntax), it is actually quite easy to write programs in a monadic style without any knowledge of the abstract theoretical underpinnings. In the same spirit, this report is directed at functional programmers with an interest in using monads as part of a practical programming project, rather than the underlying category theory. To this end, in a number of places, we have intentionally chosen to give definitions or results without reference to the corresponding categorical concepts. Readers with an interest in a more technical, category theoretic presentation of the ideas described in this report are referred to [1].

We will assume some familiarity with the motivation for monads and their use in structuring functional programs; Wadler [14, 15] provides an excellent introduction to these topics. Programming examples will be written using the syntax of Gofer, a small, experimental, purely functional language based closely on the definition of Haskell [3]. This enables us to use constructor classes [5] to show how our results can be expressed in a concrete programming language.

For completeness, we have included detailed proofs for many of our results. In keeping with our aim to avoid unnecessary technical details, most of the proofs are constructed from first principles using simple equational reasoning. Working in this manner has led to some surprising insights. For example, one of our results in Section 3 corresponds closely to a result stated in [8], but a careful study of the laws that are actually used in the proof has allowed us to weaken the hypotheses and state the result in a slightly more general form. However, recognizing that some readers may prefer to omit such details, most of the proofs are presented in boxed figures that are easily identified and skipped.

The remainder of this report is organized as follows. Section 2 gives a definition of the algebraic properties of a monad, and these are then used to describe the construction of compositions of monads in Section 3 and their converses in Section 4. Moving on to practical applications, Section 5 describes how these results about composition of monads can be expressed using the notation of constructor classes, with several examples in Section 6. A simple application, building on this framework, is included in Section 7. Finally, Section 8 illustrates that there are other ways of combining certain monads, setting a direction for future work.

2 Monads for functional programming

Wadler [14] defines a *monad* as a unary type constructor M together with three functions map , $unit$ and $join$ whose types are given by:

$$\begin{aligned} map &:: (a \rightarrow b) \rightarrow (M\ a \rightarrow M\ b) \\ unit &:: a \rightarrow M\ a \\ join &:: M\ (M\ a) \rightarrow M\ a \end{aligned}$$

In addition, these functions are required to satisfy a collection of algebraic laws:

$$map\ id = id \quad (1)$$

$$map\ f . map\ g = map\ (f . g) \quad (2)$$

$$unit . f = map\ f . unit \quad (3)$$

$$join . map\ (map\ f) = map\ f . join \quad (4)$$

$$join . unit = id \quad (5)$$

$$join . map\ unit = id \quad (6)$$

$$join . map\ join = join . join \quad (7)$$

(We use the standard infix period notation for function composition, $(f . g)\ x = f\ (g\ x)$, and the symbol id to represent the identity function $id\ x = x$. It is well known that composition of functions is associative with id as both a left and right identity; in symbols, $f . (g . h) = (f . g) . h$ and $f . id = f = id . f$ for all f , g and h of appropriate types.)

For the purposes of this report, it will be convenient to break this down into stages; if M is a unary type constructor, then we will say that:

- M is a *functor* if there is a function $map :: (a \rightarrow b) \rightarrow (M\ a \rightarrow M\ b)$ satisfying laws (1) and (2) above.
- M is a *premonad* if it is a functor with a function $unit :: a \rightarrow M\ a$ satisfying law (3) above.
- M is a *monad* if it is a premonad with a function $join :: M\ (M\ a) \rightarrow M\ a$ satisfying laws (4), (5), (6) and (7) above.

It should be mentioned that there are several other (equivalent) ways to define the concept of a monad, each of which comes with its own collection of monad operators and equational laws. For example, Wadler [15] describes how a monad can be characterized using just the *unit* function together with an operator

$$\text{bind} \quad :: \quad M \ a \rightarrow (a \rightarrow M \ b) \rightarrow M \ b$$

Both the *map* and *join* operators that we have described can be defined using *bind* and *unit*. The *bind* operator is particularly useful in practical work with monads and as a means of translating the notation of monad comprehensions [14, 15, 5]. Furthermore, only three laws are needed to specify the properties of a monad in this case. Nevertheless, we have chosen to work with the *map*, *unit*, *join* formulation of monads outlined above. One reason for this decision is that, in our opinion, many of the proofs are easier to express in this framework. In addition, we will see that the ability to distinguish between monads and premonads will also be quite useful in the following work.

3 Conditions for composition

Suppose that M and N are functors. To avoid confusion, we will write map_M and map_N for the corresponding *map* functions in each case. How is it possible to compose these functors in some sensible way? Certainly, we can think of a composition of M and N as a type constructor that takes any type a to the type $M (N \ a)$, but we also need to be able to give a definition for

$$\text{map} \quad :: \quad (a \rightarrow b) \rightarrow (M (N \ a) \rightarrow M (N \ b))$$

satisfying the functor laws (1) and (2). Fortunately, this is quite easy! If $f :: a \rightarrow b$, then we can apply map_N to obtain $\text{map}_N \ f :: N \ a \rightarrow N \ b$ and then apply map_M to obtain $\text{map}_M (\text{map}_N \ f) :: M (N \ a) \rightarrow M (N \ b)$. This gives the definition

$$\text{map} \quad = \quad \text{map}_M \cdot \text{map}_N$$

and we can use the proofs in Figure 1 to show that this does satisfy the necessary laws.

These proofs use standard techniques of equational reasoning, writing the justification for each step in the right hand column. In many cases, this is just a reference to the law that has been used, with the convention that, for example, (1M) is the version of law (1) for the constructor M . Where a step follows directly from the definition of a function (either by folding or unfolding the definition), we will simply write the name of the function involved as justification.

Composition of premonads is similar. Most of the work (the definition of the composed type constructor) has already been dealt with in the composition of functors. A suitable *unit* function for the composition is:

$$\begin{aligned} \text{unit} & \quad :: \quad a \rightarrow M (N \ a) \\ \text{unit} & \quad = \quad \text{unit}_M \cdot \text{unit}_N \end{aligned}$$

$map\ id$	= $map_M (map_N id)$	map
	= $map_M id$	(1N)
	= id	(1M)
$map\ f . map\ g$	= $map_M (map_N f) . map_M (map_N g)$	map, map
	= $map_M (map_N f . map_N g)$	(2M)
	= $map_M (map_N (f . g))$	(2N)
	= $map (f . g)$	map

Figure 1: Composition of functors, laws (1)–(2)

($unit_M$ and $unit_N$ being the unit functions for M and N respectively) and the proof in Figure 2 demonstrates that this does indeed satisfy law (3).

$unit . f$	= $unit_M . unit_N . f$	$unit$
	= $unit_M . map_N f . unit_N$	(3N)
	= $map_M (map_N f) . unit_M . unit_N$	(3M)
	= $map f . unit$	$map, unit$

Figure 2: Composition of premonads, law (3)

Unfortunately, our real goal, composition of monads, is rather more difficult. Recall that, to define the composition of two monads M and N , we need to find a function:

$$join :: M (N (M (N a))) \rightarrow M (N a)$$

that satisfies the monad laws (4)–(7). As a first guess, and following the pattern of the previous examples, we might consider the function $join_M . join_N$ where $join_M$ and $join_N$ are the $join$ functions in the component monads. But this is, in general, not even type correct! If $x :: N (N a)$ (for some a), then the result of $join_N x$ will have type $N a$ while $join_M$ expects an argument with a type of the form $M (M b)$.

In fact, we can actually prove that, in a certain sense, there is no way to construct a $join$ function with the type above using only the operations of the two monads (see the appendix for an outline of the proof). It follows that the only way that we might hope to form a composition is if there are some additional constructions linking the two components. In this report we will concentrate on four methods for constructing a composite monad, described in the following sections.

3.1 The trivial case, by definition

The composite of two premonads M and N is a monad if there is a polymorphic function:

$$join :: M (N (M (N a))) \rightarrow M (N a)$$

satisfying the laws (4), (5), (6) and (7). This, of course, follows directly from the definitions above.

3.2 The *prod* construction

The composition of a monad M with a premonad N is itself a monad if there is a polymorphic function:

$$prod :: N (M (N a)) \rightarrow M (N a)$$

with the *join* function defined by:

$$join = join_M . map_M prod$$

satisfying the following four laws:

$$\begin{aligned} prod . map_N (map f) &= map f . prod && P(1) \\ prod . unit_N &= id && P(2) \\ prod . map_N unit &= unit_M && P(3) \\ prod . map_N join &= join . prod && P(4) \end{aligned}$$

The proofs for this are given in Figure 3. One interesting point is that, when we started this work, we had assumed that it would be necessary to require that both M and N be monads. In fact, from the proofs, we see that we did not actually need any of the monad laws for N at all! As a result, we can construct a ‘composition of monads’ under somewhat weaker conditions than originally expected.

3.3 The *dorp* construction

The composition of a premonad M with a monad N is itself a monad if there is a polymorphic function:

$$dorp :: M (N (M a)) \rightarrow M (N a)$$

with the *join* function defined by:

$$join = map_M join_N . dorp$$

such that the following laws hold:

$$\begin{aligned} dorp . map (map_M f) &= map f . dorp && D(1) \\ dorp . unit &= map_M unit_N && D(2) \\ dorp . map unit_M &= id && D(3) \\ dorp . join &= join . map dorp && D(4) \end{aligned}$$

Full proofs for this are given in Figure 4. In this case, we use the fact that N is a monad, although M can be an arbitrary premonad.

$join . map (map f)$	
$= join_M . map_M prod . map_M (map_N (map f))$	$join, map$
$= join_M . map_M (prod . map_N (map f))$	(2M)
$= join_M . map_M (map f . prod)$	P(1)
$= join_M . map_M (map f) . map_M prod$	(2M)
$= join_M . map_M (map_M (map_N f)) . map_M prod$	map
$= map_M (map_N f) . join_M . map_M prod$	(4M)
$= map f . join$	$map, join$
$join . unit$	
$= join_M . map_M prod . unit_M . unit_N$	$join, unit$
$= join_M . unit_M . prod . unit_N$	(3M)
$= prod . unit_N$	(5M)
$= id$	P(2)
$join . map unit$	
$= join_M . map_M prod . map_M (map_N unit)$	$join, map$
$= join_M . map_M (prod . map_M unit)$	(2M)
$= join_M . map_M unit_M$	P(3)
$= id$	(6M)
$join . map join$	
$= join_M . map_M prod . map_M (map_N join)$	$join, map$
$= join_M . map_M (prod . map_N join)$	(2M)
$= join_M . map_M (join . prod)$	P(4)
$= join_M . map_M (join_M . map_M prod . prod)$	$join$
$= join_M . map_M join_M . map_M (map_M prod) . map_M prod$	(2M)
$= join_M . join_M . map_M (map_M prod) . map_M prod$	(7M)
$= join_M . map_M prod . join_M . map_M prod$	(4M)
$= join . join$	$join, join$

Figure 3: Composition of monads, laws (4)–(7), with $join = join_M . map_M prod$

$join . map (map f)$	
$= map_M join_N . dorp . map (map_M (map_N f))$	$join, map$
$= map_M join_N . map (map_N f) . dorp$	$D(1)$
$= map_M join_N . map_M (map_N (map_N f)) . dorp$	map
$= map_M (join_N . map_N (map_N f)) . dorp$	$(2M)$
$= map_M (map_N f . join_N) . dorp$	$(4N)$
$= map_M (map_N f) . map_M join_N . dorp$	$(2M)$
$= map f . join$	$map, join$
$join . unit$	
$= map_M join_N . dorp . unit$	$join$
$= map_M join_N . map_M unit_N$	$D(2)$
$= map_M (join_N . unit_N)$	$(2M)$
$= map_M id$	$(5N)$
$= id$	$(1M)$
$join . map unit$	
$= map_M join_N . dorp . map (unit_M . unit_N)$	$join, unit$
$= map_M join_N . dorp . map unit_M . map unit_N$	(2)
$= map_M join_N . map unit_N$	$D(3)$
$= map_M join_N . map_M (map_N unit_N)$	map
$= map_M (join_N . map_N unit_N)$	$(2M)$
$= map_M id$	$(6N)$
$= id$	$(1M)$
$join . map join$	
$= map_M join_N . dorp . map (map_M join_N . dorp)$	$join, join$
$= map_M join_N . dorp . map (map_M join_N) . map dorp$	(2)
$= map_M join_N . map join_N . dorp . map dorp$	$D(1)$
$= map_M join_N . map_M (map_N join_N) . dorp . map dorp$	map
$= map_M (join_N . map_N join_N) . dorp . map dorp$	$(2M)$
$= map_M (join_N . join_N) . dorp . map dorp$	$(7N)$
$= map_M join_N . map_M join_N . dorp . map dorp$	$(2M)$
$= map_M join_N . join . map dorp$	$join$
$= map_M join_N . dorp . join$	$D(4)$
$= join . join$	$join$

Figure 4: Composition of monads, laws (4)–(7), with $join = map_M join_N . dorp$

3.4 The *swap* construction

We can also define the composition of two monads M and N in terms of a polymorphic function:

$$\text{swap} :: N (M a) \rightarrow M (N a)$$

with the *join* function defined by:

$$\text{join} = \text{map}_M \text{join}_N . \text{join}_M . \text{map}_M \text{swap}$$

satisfying a collection of laws given below. Note that this is equivalent to:

$$\text{join} = \text{join}_M . \text{map}_M (\text{map}_M \text{join}_N . \text{swap})$$

since:

$$\begin{aligned} & \text{map}_M \text{join}_N . \text{join}_M . \text{map}_M \text{swap} \\ &= \text{join}_M . \text{map}_M (\text{map}_M \text{join}_N) . \text{map}_M \text{swap} \quad (4M) \\ &= \text{join}_M . \text{map}_M (\text{map}_M \text{join}_N . \text{swap}). \quad (2M) \end{aligned}$$

In order to state the laws for *swap*, we define two (familarly named) functions as abbreviations for expressions involving *swap*:

$$\begin{aligned} \text{prod} &= \text{map}_M \text{join}_N . \text{swap} \\ \text{dorp} &= \text{join}_M . \text{map}_M \text{swap} \end{aligned}$$

We will justify the use of these names below, but first we state the laws that the *swap* function must satisfy:

$$\begin{aligned} \text{swap} . \text{map}_N (\text{map}_M f) &= \text{map}_M (\text{map}_N f) . \text{swap} & \text{S(1)} \\ \text{swap} . \text{unit}_N &= \text{map}_M \text{unit}_N & \text{S(2)} \\ \text{swap} . \text{map}_N \text{unit}_M &= \text{unit}_M & \text{S(3)} \\ \text{prod} . \text{map}_N \text{dorp} &= \text{dorp} . \text{prod} & \text{S(4)} \end{aligned}$$

It is possible to prove the existence of the composite monad using these definitions alone. However, it is easier, and perhaps more instructive, to do this indirectly by exploring the relationship between *swap*, *prod* and *dorp*.

First of all, note that the definition of *join* coincides with the *join* that would be obtained using the *prod* construction (with the above definition of *prod*):

$$\begin{aligned} \text{join}_M . \text{map}_M \text{prod} &= \text{join}_M . \text{map}_M (\text{map}_M \text{join}_N . \text{swap}) & \text{prod} \\ &= \text{join} & \text{join} \end{aligned}$$

Furthermore, the definition of *join* also coincides with the *join* function that we would obtain using the *dorp* construction (with the above definition of *dorp*):

$$\begin{aligned} \text{map}_M \text{join}_N . \text{dorp} &= \text{map}_M \text{join}_N . (\text{join}_M . \text{map}_M \text{swap}) & \text{dorp} \\ &= \text{join} & \text{join} \end{aligned}$$

Assuming that laws S(1)–S(4) are satisfied, the proofs in Figures 5 and 6 show that *prod* and *dorp* satisfy laws P(1)–P(4) and D(1)–D(4), respectively. Note that, in each case, both M and N must be monads if the composite is also to be a monad. For example, the proof of P(1)–P(4) requires the monad laws for N , while the *prod* construction requires that M should also be a monad.

$prod . map_N (map f)$	
$= map_M join_N . swap . map_N (map_M (map_N f))$	$prod, map$
$= map_M join_N . map_M (map_N (map_N f)) . swap$	S(1)
$= map_M (join_N . map_N (map_N f)) . swap$	(2M)
$= map_M (map_N f . join_N) . swap$	(4N)
$= map_M (map_N f) . map_M join_N . swap$	(2M)
$= map f . prod$	$map, prod$
$prod . unit_N$	
$= map_M join_N . swap . unit_N$	$prod$
$= map_M join_N . map_M unit_N$	S(2)
$= map_M (join_N . unit_N)$	(2M)
$= map_M id$	(5N)
$= id$	(1M)
$prod . map_N unit$	
$= map_M join_N . swap . map_N (unit_M . unit_N)$	$prod, unit$
$\doteq map_M join_N . swap . map_N unit_M . map_N unit_N$	(2N)
$= map_M join_N . unit_M . map_N unit_N$	S(3)
$= unit_M . join_N . map_N unit_N$	(3M)
$= unit_M$	(6N)
$prod . map_N join$	
$= prod . map_N (map_M join_N . dorp)$	$join$
$= map_M join_N . swap . map_N (map_M join_N) . map_N dorp$	$prod, (2N)$
$= map_M join_N . map_M (map_N join_N) . swap . map_N dorp$	S(1)
$= map_M (join_N . map_N join_N) . swap . map_N dorp$	(2M)
$= map_M (join_N . join_N) . swap . map_N dorp$	(7N)
$= map_M join_N . map_M join_N . swap . map_N dorp$	(2M)
$= map_M join_N . prod . map_N dorp$	$prod$
$= map_M join_N . dorp . prod$	S(4)
$= join . prod$	$join$

Figure 5: Proof of P(1)–P(4) from S(1)–S(4) with $prod = map_M join_N . swap$

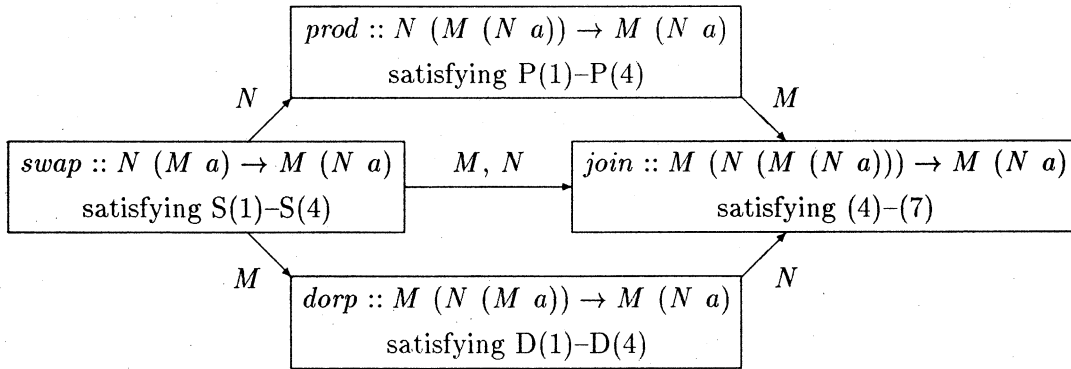
$dorp . map (map_M f)$	
$= join_M . map_M swap . map_M (map_N (map_M f))$	$dorp, map$
$= join_M . map_M (swap . map_N (map_M f))$	(2M)
$= join_M . map_M (map_M (map_N f) . swap)$	S(1)
$= join_M . map_M (map_M (map_N f)) . map_M swap$	(2M)
$= map_M (map_N f) . join_M . map_M swap$	(4M)
$= map f . dorp$	$map, dorp$
$dorp . unit$	
$= join_M . map_M swap . unit_M . unit_N$	$dorp, unit$
$= join_M . unit_M . swap . unit_N$	(3M)
$= swap . unit_N$	(5M)
$= map_M unit_N$	S(2)
$dorp . map unit_M$	
$= join_M . map_M swap . map_M (map_N unit_M)$	$dorp, map$
$= join_M . map_M (swap . map_N unit_M)$	(2M)
$= join_M . map_M unit_M$	S(3)
$= id$	(6M)
$dorp . join$	
$= join_M . map_M swap . join_M . map_M prod$	$dorp, join$
$= join_M . join_M . map_M (map_M swap) . map_M prod$	(4M)
$= join_M . map_M join_M . map_M (map_M swap) . map_M prod$	(7M)
$= join_M . map_M (join_M . map_M swap . prod)$	(2M)
$= join_M . map_M (dorp . prod)$	$dorp$
$= join_M . map_M (prod . map_N dorp)$	S(4)
$= join_M . map_M prod . map dorp$	(2M), map
$= join . map dorp$	$join$

Figure 6: Proof of D(1)–D(4) from S(1)–S(4) with $dorp = join_M . map_M swap$

3.5 Summary

At first glance, the constructions in the previous sections may seem rather mysterious; in each case, we gave a type for some polymorphic function, stated some laws that it should satisfy ... and ‘presto!’ we have another way of composing monads. In fact, these constructions were discovered largely by experimentation, ‘guessing’ a definition of *join* in some particular form, for example $join = join_M \cdot map_M \cdot prod$, then attempting to prove the monad laws, for example, to determine what properties *prod* should satisfy. Types played an essential part in this process, helping to suggest ways to define *join* and ensuring that the laws we used are well-typed.

The following diagram summarizes these results and the relationship between the different constructions for compositions of monads:



In each case, the arrows between different constructions represent implications, with the labels indicating which of the constructors *M* and *N* is required to be a monad.

4 Converse results

We now have a number of different ways of composing two monads, but we have not made any attempt to see how general each approach might be. In particular, it is natural to ask what kinds of monads can be obtained using the constructions described above. Thinking of the diagram at the end of the previous section, we know by definition that all compositions of *M* and *N* must have a *join* function as specified by the rightmost box. The goal of this section is to establish what conditions are necessary to move back, in the opposite direction of the arrows, to each of the other three constructions.

More formally, suppose that *M* and *N* are monads and that there is a composition of *M* and *N* with operators *map*, *unit* and *join* satisfying the laws (1)–(7) and such that:

$$\begin{aligned} map &= map_M \cdot map_N \\ unit &= unit_M \cdot unit_N \end{aligned}$$

The problem now is to determine which composite monads defined in this way can be obtained using each of the *prod*, *dorp* and *swap* constructions.

4.1 The *prod* construction

It is actually fairly easy to give a definition for a *prod* function (of the required type) in terms of the various monad operators available:

$$\mathit{prod} = \mathit{join} . \mathit{unit}_M$$

Showing that this definition of *prod* satisfies the laws P(1)–P(4) is also straightforward, as detailed by the proofs in Figure 7. In fact, the only difficulty arises when we try to

$\mathit{map} f . \mathit{prod}$	$= \mathit{map} f . \mathit{join} . \mathit{unit}_M$	prod
	$= \mathit{join} . \mathit{map} (\mathit{map} f) . \mathit{unit}_M$	(4)
	$= \mathit{join} . \mathit{map}_M (\mathit{map}_N (\mathit{map} f)) . \mathit{unit}_M$	map
	$= \mathit{join} . \mathit{unit}_M . \mathit{map}_N (\mathit{map} f)$	(3M)
	$= \mathit{prod} . \mathit{map}_N (\mathit{map} f)$	prod
$\mathit{prod} . \mathit{unit}_N$	$= \mathit{join} . \mathit{unit}_M . \mathit{unit}_N$	prod
	$= \mathit{join} . \mathit{unit}$	unit
	$= \mathit{id}$	(5)
$\mathit{prod} . \mathit{map}_N \mathit{unit}$	$= \mathit{join} . \mathit{unit}_M . \mathit{map}_N \mathit{unit}$	prod
	$= \mathit{join} . \mathit{map}_M (\mathit{map}_N \mathit{unit}) . \mathit{unit}_M$	(3M)
	$= \mathit{join} . \mathit{map} \mathit{unit} . \mathit{unit}_M$	map
	$= \mathit{unit}_M$	(6)
$\mathit{prod} . \mathit{map}_N \mathit{join}$	$= \mathit{join} . \mathit{unit}_M . \mathit{map}_N \mathit{join}$	prod
	$= \mathit{join} . \mathit{map}_M (\mathit{map}_N \mathit{join}) . \mathit{unit}_M$	(3M)
	$= \mathit{join} . \mathit{map} \mathit{join} . \mathit{unit}_M$	map
	$= \mathit{join} . \mathit{join} . \mathit{unit}_M$	(7)
	$= \mathit{join} . \mathit{prod}$	prod

Figure 7: Proof of P(1)–P(4) from the existence of a composition

show that the *join* function we obtain from the *prod* construction is the same as the *join* function that we started with. One way to do this is as follows:

$$\begin{aligned}
 \mathit{join}_M . \mathit{map}_M \mathit{prod} &= \mathit{join}_M . \mathit{map}_M (\mathit{join} . \mathit{unit}_M) && \mathit{prod} \\
 &= \mathit{join}_M . \mathit{map}_M \mathit{join} . \mathit{map}_M \mathit{unit}_M && (2M) \\
 &= \mathit{join} . \mathit{join}_M . \mathit{map}_M \mathit{unit}_M && J(1) \\
 &= \mathit{join} && (6M)
 \end{aligned}$$

Note that this depends on an assumption about the way that *join* and *join_M* ‘commute’ with one another:

$$join_M . map_M join = join . join_M \quad J(1)$$

This condition may seem a little arbitrary, but it turns out that every composite monad obtained by the *prod* construction has this property:

$$\begin{aligned} join_M . map_M join &= join_M . map_M (join_M . map_M prod) && join \\ &= join_M . map_M join_M . map_M (map_M prod) && (2M) \\ &= join_M . join_M . map_M (map_M prod) && (7M) \\ &= join_M . map_M prod . join_M && (4M) \\ &= join . join_M && join \end{aligned}$$

It follows that the set of composite monads that can be obtained using the *prod* construction are precisely those satisfying J(1).

4.2 The *dorp* construction

As in the previous case, it is easy to find a suitably typed definition of the *dorp* function using the operations of the composite monad and its components:

$$dorp = join . map (map_M unit_N)$$

The proofs in Figure 8 show that this definition satisfies the laws D(1)–D(4) as we would hope. Once again, the most difficult task is to show that the *join* function obtained from this *dorp* function using the earlier construction is equal to the *join* operator in the composite monad. One way to prove this is as follows:

$$\begin{aligned} map_M join_N . dorp &= map_M join_N . join . map (map_M unit_N) && dorp \\ &= join . map (map_M join_N) . map (map_M unit_N) && J(2) \\ &= join . map (map_M (join_N . unit_N)) && (2), (2M) \\ &= join . map (map_M id) && (5N) \\ &= join && (1M), (1) \end{aligned}$$

This also depends on an assumed law, this time linking the behaviour of *join* and *join_N*:

$$join . map (map_M join_N) = map_M join_N . join \quad J(2)$$

In fact, law J(2) holds for all composite monads obtained using the *dorp* construction, as demonstrated by the following:

$$\begin{aligned} join . map (map_M join_N) &= map_M join_N . dorp . map (map_M join_N) && join \\ &= map_M join_N . map join_N . dorp && D(1) \\ &= map_M (join_N . map_N join_N) . dorp && map, (2M) \\ &= map_M (join_N . join_N) . dorp && (7N) \\ &= map_M join_N . map_M join_N . dorp && (2M) \\ &= map_M join_N . join && join \end{aligned}$$

$dorp . map (map_M f)$	
$= join . map (map_M unit_N) . map (map_M f)$	$dorp$
$= join . map (map_M (unit_N . f))$	(2), (2M)
$= join . map (map_M (map_N f . unit_N))$	(3N)
$= join . map (map f) . map (map_M unit_N)$	(2), (2M), map
$= map f . join . map (map_M unit_N)$	(4)
$= map f . dorp$	$dorp$
$dorp . unit$	
$= join . map (map_M unit_N) . unit$	$dorp$
$= join . unit . map_M unit_N$	(3)
$= map_M unit_N$	(5)
$dorp . map unit_M$	
$= join . map (map_M unit_N) . map unit_M$	$dorp$
$= join . map (map_M unit_N . unit_M)$	(2)
$= join . map (unit_M . unit_N)$	(3M)
$= join . map unit$	$unit$
$= id$	(6)
$dorp . join$	
$= join . map (map_M unit_N) . join$	$dorp$
$= join . join . map (map (map_M unit_N))$	(4)
$= join . map join . map (map (map_M unit_N))$	(7)
$= join . map (join . map (map_M unit_N))$	(2)
$= join . map dorp$	$dorp$

Figure 8: Proof of D(1)–D(4) from the existence of a composition

4.3 The *swap* construction

Our goal in this section is to determine the class of composite monads that can be constructed using the *swap* construction presented in Section 3.4. From the results given there, any composite obtained using *swap* can also be obtained from a *prod* or a *dorp* construction, so it follows (from the results in the last two sections) that any monad obtained from *swap* must satisfy at least J(1) and J(2). In fact, we will now show that these two properties are not only necessary, but also sufficient.

Following the pattern in the previous cases, we start with a definition for the *swap* function in terms of the *join* of the composite monad:

$$\text{swap} = \text{join} . \text{unit}_M . \text{map}_N (\text{map}_M \text{unit}_N)$$

or, equivalently:

$$\text{swap} = \text{join} . \text{map} (\text{map}_M \text{unit}_N) . \text{unit}_M$$

since:

$$\begin{aligned} & \text{join} . \text{unit}_M . \text{map}_N (\text{map}_M \text{unit}_N) \\ &= \text{join} . \text{map}_M (\text{map}_N (\text{map}_M \text{unit}_N)) . \text{unit}_M && (3M) \\ &= \text{join} . \text{map} (\text{map}_M \text{unit}_N) . \text{unit}_M && \text{map} \end{aligned}$$

Note that these definitions of *swap* can also be expressed in terms of the *prod* and *dorp* functions used in the previous two sections:

$$\begin{aligned} \text{swap} &= \text{join} . \text{unit}_M . \text{map}_N (\text{map}_M \text{unit}_N) && \text{swap} \\ &= \text{prod} . \text{map}_N (\text{map}_M \text{unit}_N) && \text{prod} \\ \\ \text{swap} &= \text{join} . \text{map} (\text{map}_M \text{unit}_N) . \text{unit}_M && \text{swap} \\ &= \text{dorp} . \text{unit}_M && \text{dorp} \end{aligned}$$

Assuming J(2), we can show that the definition of *prod* in terms of *swap* given in Section 3.4 coincides with the *prod* function specified in terms of *join* in Section 4.1:

$$\begin{aligned} & \text{map}_M \text{join}_N . \text{swap} \\ &= \text{map}_M \text{join}_N . \text{join} . \text{unit}_M . \text{map}_N (\text{map}_M \text{unit}_N) && \text{swap} \\ &= \text{join} . \text{map} (\text{map}_M \text{join}_N) . \text{unit}_M . \text{map}_N (\text{map}_M \text{unit}_N) && \text{J(2)} \\ &= \text{join} . \text{unit}_M . \text{map}_N (\text{map}_M \text{join}_N) . \text{map}_N (\text{map}_M \text{unit}_N) && \text{map, (3M)} \\ &= \text{join} . \text{unit}_M . \text{map}_N (\text{map}_M (\text{join}_N . \text{unit}_N)) && (2N), (2M) \\ &= \text{join} . \text{unit}_M && (5N), (1M), (1N) \end{aligned}$$

In a similar way, assuming J(1), we can show that the definition of *dorp* from *swap* in Section 3.4 gives the same function as the definition of *dorp* in Section 4.2:

$$\begin{aligned} & \text{join}_M . \text{map}_M \text{swap} \\ &= \text{join}_M . \text{map}_M (\text{join} . \text{map} (\text{map}_M \text{unit}_N) . \text{unit}_M) && \text{swap} \\ &= \text{join}_M . \text{map}_M \text{join} . \text{map}_M (\text{map} (\text{map}_M \text{unit}_N) . \text{unit}_M) && (2M) \\ &= \text{join} . \text{join}_M . \text{map}_M (\text{map} (\text{map}_M \text{unit}_N) . \text{unit}_M) && \text{J(1)} \\ &= \text{join} . \text{map} (\text{map}_M \text{unit}_N) . \text{join}_M . \text{map}_M \text{unit}_M && \text{map, (4M), (2M)} \\ &= \text{join} . \text{map} (\text{map}_M \text{unit}_N) && (6M) \end{aligned}$$

With these results, it is straightforward to show that, given a composite monad satisfying J(1) and J(2), the *swap* function defined above satisfies the laws S(1)–S(4) required for the *swap* construction; see Figure 9 for the proofs. For convenience, we have used

$swap . map_N (map_M f)$		
$= dorp . unit_M . map_N (map_M f)$		<i>swap</i>
$= dorp . map (map_M f) . unit_M$		(3M), <i>map</i>
$= map f . dorp . unit_M$		D(1)
$= map_M (map_N f) . swap$		<i>map</i> , <i>swap</i>
$swap . unit_N$		
$= dorp . unit_M . unit_N$		<i>swap</i>
$= dorp . unit$		<i>unit</i>
$= map_M unit_N$		D(2)
$swap . map_N unit_M$		
$= dorp . unit_M . map_N unit_M$		<i>swap</i>
$= dorp . map unit_M . unit_M$		(3M), <i>map</i>
$= unit_M$		D(3)
$dorp . prod$		
$= dorp . join . unit_M$		<i>prod</i>
$= join . map dorp . unit_M$		D(4)
$= join . unit_M . map_N dorp$		<i>map</i> , (3M)
$= prod . map_N dorp$		<i>prod</i>

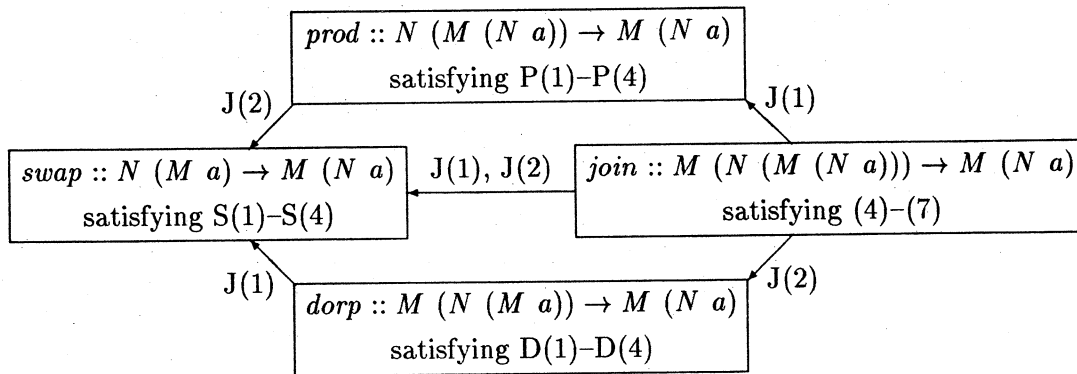
Figure 9: Proof of S(1)–S(4) from the existence of a composition

the properties of *dorp* described by the laws D(1)–D(4) to establish these properties of *swap*. Similar derivations are also possible using the laws for *prod*, or directly using the definition of *swap* in terms of *join* although the proofs are a little longer, particularly in the second case.

Given the comments above, this completes the proof that the set of monads which can be constructed from a *swap* function are precisely those satisfying both J(1) and J(2).

4.4 Summary

Corresponding to the diagram in Section 3.5, we can summarize the results about the converses for the *prod*, *dorp* and *swap* constructions as follows:



This time, we have labeled the arrows between the different constructions with the additional properties required to establish the converse.

The properties J(1) and J(2) that we have used in these results are actually fairly similar, each given by a law of the form:

$$join . map f = f . join$$

For J(1), the *join* and *map* functions here should be read as the corresponding operators for the monad *M* with *f* the *join* function in the composite monad. For J(2), the *join* and *map* operators are taken from the composite monad with $f = map_M join_N$.

Incidentally, functions satisfying a law of the form shown above are actually quite widely studied in functional programming. For example, in the special case of the list monad where *join* is just the concatenation of a list of lists (the *concat* function in Haskell), functions satisfying the law above are often referred to as *list homomorphisms*.

5 Programming monad composition

Our goal in this and remaining sections is to show how the different constructions for monad composition presented above can be used in a practical programming language. For convenience, we will use the notation of *constructor classes* [5] implemented as part of the Gofer system, although the same ideas can also be applied to a much wider range of languages.

In this section, we present a general framework for working with different forms of monad composition. Later, we will describe a number of concrete examples of monad composition, and use some of these in a simple application.

5.1 Representing functors

The Haskell programming language [3] makes use of a system of *type classes* to provide a flexible treatment of ad-hoc polymorphism. A type class is a set of types, often referred to as the *instances* of the class, together with a family of operations that are defined for each instance. Constructor classes are a natural extension, allowing classes of

type constructors. Since types are just nullary constructors (i.e. taking no arguments), this includes Haskell style type classes as a special case. For example, the following declaration introduces a class *Functor* and specifies that, for each instance *f*, there is a *map* function that maps functions of type $(a \rightarrow b)$ to functions of type $(f\ a \rightarrow f\ b)$:

```
class Functor f where
  map :: (a → b) → (f a → f b)
```

This corresponds fairly closely to the definition of a functor in Section 2, except that it does not include the functor laws (1) and (2). Equality of functions is not computable so there is no way for a compiler to ensure that these laws are satisfied. However, it is useful to be able to include these as part of the program to check that they are, at the very least, type correct. One way to do this is to define an operator representing the desired equality:

```
(===)    :: a → a → Law a
x === y  =  error "uncomputable equality"
```

The type signature here is particularly important, specifying that the two values being compared must be of the same type. The exact definition of the *Law a* type isn't important¹, but the argument type is used to record the type of the values that are asserted as being equal. Note however that any attempt to use this operator to compare two values will produce a runtime error.

Using this technique, we can represent the functor laws by the function definitions:

```
law1      :: Functor f ⇒ () → Law (f a → f a)
law1 ()   =  map id === id

law2      :: Functor f ⇒ (b → c) → (a → b) → Law (f a → f c)
law2 f g  =  map f . map g === map (f . g)
```

Notice the use of function arguments to model free variables in the definition of *law2*. The types of these variables, as well as the type of the expressions being compared, can be read directly from the type signature. Note that, although we have written this explicitly as part of the program, the same type could also have been obtained automatically by the Gofer type inference mechanism.

In the first law we have used the unit value $()$ to emphasize the fact that the law has no free variables².

¹The definition we used was: `data Law a = Unspecified.`

²The addition of a dummy argument also helps to avoid the monomorphism restriction in Haskell/Gofer so that the type of *law1* can be calculated by the type inference mechanism. Without the dummy argument, an explicit type signature would be mandatory, not optional as it is here.

5.2 Representing premonads

With the framework established above, the representation of premonads is straightforward. The class of premonads can be described by the class:

```
class Functor m ⇒ Premonad m where
  unit :: a → m a
```

Note the first line in this declaration that captures the requirement that every premonad is also a functor; another way of saying this is that *Premonad* is a subclass of *Functor*. The premonad law, referred to as (3) in the first part of this report, can now be represented by:

```
law3    :: Premonad m ⇒ (a → b) → Law (a → m b)
law3 f  = map f . unit === unit . f
```

5.3 Representing monads

The representation of monads also follows directly from our earlier definitions, captured by the class *Monad*, which is a subclass of *Premonad*:

```
class Premonad m ⇒ Monad m where
  join :: m (m a) → m a
```

The monad laws (4)–(7) are represented by the following:

```
law4    :: Monad m ⇒ (a → b) → Law (m (m a) → m b)
law4 f  = join . map (map f) === map f . join

law5    :: Monad m ⇒ () → Law (m a → m a)
law5 () = join . unit === id

law6    :: Monad m ⇒ () → Law (m a → m a)
law6 () = join . map unit === id

law7    :: Monad m ⇒ () → Law (m (m (m a)) → m a)
law7 () = join . map join === join . join
```

5.4 A general framework for composition constructions

It is easy to describe the composition of functors and premonads using the definitions:

```
mapC    :: (Functor f, Functor g) ⇒ (a → b) → (f (g a) → f (g b))
mapC    = map . map

unitC   :: (Premonad m, Premonad n) ⇒ a → m (n a)
unitC   = unit . unit
```

However, neither of these functions has a type of the right form to be able to define an instance of the *Functor* or *Premonad* classes. Furthermore, when we consider the different constructions for monad composition, there is nothing in a type expression of the form $f (g x)$ to indicate which construction is intended. To avoid this problem, we will define a different constructor c for each of the composition constructions with the intention that $c f g x$ is isomorphic to the composition $f (g x)$, identifying the construction used.

To simplify the task of converting between values of type $c f g x$ and those of type $f (g x)$, we introduce a constructor class to describe the required isomorphisms³:

```
class Composer  $c$  where
  open  ::  $c f g x \rightarrow f (g x)$ 
  close ::  $f (g x) \rightarrow c f g x$ 
```

Using these functions, we can package up the *mapC* operator defined above to give an instance of the *Functor* class:

```
instance (Composer  $c$ , Functor  $f$ , Functor  $g$ )  $\Rightarrow$  Functor ( $c f g$ ) where
  map  $f$  = close . mapC  $f$  . open
```

Note that this definition can be used with any suitable c ; we do not need to repeat the definition of the *map* function for each different construction.

The definition of the composition of premonads can also be dealt with in a similar manner.

```
instance (Composer  $c$ , Premonad  $m$ , Premonad  $n$ )  $\Rightarrow$  Premonad ( $c m n$ ) where
  unit = close . unitC
```

In each of these instance declarations, we used compositions with *open* and *close* to modify the type of a value so that it could be used to define an instance of a particular class. The following function will be used a number of times in subsequent sections to wrap up the definition of a *join* function in an instance of the *Monad* class:

```
wrap    :: (Composer  $c$ , Functor  $m$ , Functor  $n$ )  $\Rightarrow$ 
           ( $m (n (m (n a))) \rightarrow m (n a) \rightarrow$ 
            ( $c m n (c m n a) \rightarrow c m n a$ )
  wrap  $j$  = close .  $j$  . mapC open . open
```

The type of this function may seem a little daunting. However, all it really does is convert a function with a type suitable for the *join* function of a composition to an equivalent form using an instance c of the *Composer* class.

³As another, technical aside, the implementations that we give in later sections are not strictly isomorphisms – the representations we use for each $c f g x$ are actually isomorphic to the lifted representation of $f (g x)$, not to $f (g x)$ itself. This could be fixed using strict constructors or irrefutable pattern matching; for the purposes of this work, we will assume that *open* and *close* are genuine isomorphisms and we will not consider the details any further here.

We will also need a way to embed computations in one component into the composite monad. This can be accomplished using the following functions:

```

right :: (Composer c, Premonad f) => g a -> c f g a
right = close . unit

left  :: (Composer c, Functor f, Premonad g) => f a -> c f g a
left  = close . map unit

```

See Section 7 for a concrete example of this.

5.5 Programming the *prod* construction

Our first construction, and our first application of the *Composer* class introduced above, is based on the *prod* construction first described in Section 3.2. We will use the following composer to identify this particular construction.

```

data PComp f g x = PC (f (g x))

instance Composer PComp where
  open (PC x) = x
  close      = PC

```

The construction itself requires a *prod* function, as described by the following class declaration:

```

class (Monad m, Premonad n) => PComposable m n where
  prod :: n (m (n a)) -> m (n a)

```

Note that *PComposable* uses two parameters and that the superclass constraints capture the requirement that *m* is a monad, while only a premonad structure is needed for *n*. The definition of the composite *join* function follows directly from our earlier results:

```

joinP :: (PComposable m n) => m (n (m (n a))) -> m (n a)
joinP = join . map prod

```

For good measure, we will also include the laws P(1)–P(4) that are required for the *prod* construction. For brevity, we omit the corresponding type signatures, all of which can in any case, be inferred automatically:

```

p1 f = prod . map (mapC f) === mapC f . prod
p2 () = prod . unit === id
p3 () = prod . map unitC === unit
p4 () = prod . map joinP === joinP . prod

```

Finally, we can package up the *joinP* function defined above to define a new instance of the *Monad* class (the corresponding superclass instances for *Functor* and *Premonad* are already covered by the definitions in the previous section):

```

instance PComposable m n => Monad (PComp m n) where
  join = wrap joinP

```

5.6 Programming the *dorp* construction

With the previous section still fresh in our minds, the treatment of the *dorp* construction in this section is unlikely to cause any big surprises. We begin with the definition of a new composer to identify compositions obtained from this construction:

```
data DComp f g x = DC (f (g x))  
instance Composer DComp where  
  open (DC x) = x  
  close      = DC
```

The construction requires a function *dorp*, specified by:

```
class (Premonad m, Monad n) => DComposable m n where  
  dorp :: m (n (m a)) -> m (n a)
```

and yields a monad structure with a *join* function given by:

```
joinD :: (DComposable m n) => m (n (m (n a))) -> m (n a)  
joinD = map join . dorp
```

Any definition of *dorp* is expected to satisfy the laws D(1)–D(4) which are represented as follows:

```
d1 f = dorp . mapC (map f) == mapC f . dorp  
d2 () = dorp . unitC == map unit  
d3 () = dorp . mapC unit == id  
d4 () = dorp . joinD == joinD . mapC dorp
```

And we can package the *dorp* construction as an instance of the *Monad* class using the following declaration:

```
instance (DComposable m n) => Monad (DComp m n) where  
  join = wrap joinD
```

5.7 Programming the *swap* construction

Following, once again, the pattern of the previous sections, we begin our implementation of the *swap* construction with the definition of a corresponding new composer:

```
data SComp f g x = SC (f (g x))  
instance Composer SComp where  
  open (SC x) = x  
  close      = SC
```

To compose two monads using this technique, we require a *swap* function given by:

```
class (Monad m, Monad n) => SComposable m n where  
  swap :: n (m a) -> m (n a)
```

and satisfying the laws S(1)–S(4), represented by:

```

s1 f = swap . mapC f === mapC f . swap
s2 () = swap . unit === map unit
s3 () = swap . map unit === unit
s4 () = prod . map dorp === dorp . prod
      where prod = map join . swap
            dorp = join . map swap

```

This allows us to define a monad structure on the composition with *join* function given by:

```

joinS :: (SComposable m n) => m (n (m (n a))) -> m (n a)
joinS = map join . join . map swap

```

which leads to the following instance of the *Monad* class:

```

instance (SComposable m n) => Monad (SComp m n) where
  join = wrap joinS

```

Finally, we can capture some aspects of the relationship between the different constructions using the following instance declarations.

```

instance (SComposable m n) => PComposable m n where
  prod = map join . swap

instance (SComposable m n) => DComposable m n where
  dorp = join . map swap

```

These definitions reflect the fact, proved in Section 3.4, that values for *prod* and *dorp* can always be derived from a suitable definition of *swap*.

6 Some specific monad constructions

Now that we have established the basic framework for our approach to monad composition, we will show how our results can be used to compose some specific monads. We have already seen that we cannot define a composition that works for two arbitrary monads. The next best option is to fix one of the components in a composition to be a particular monad and allow the other component to range over a family of different monads.

To illustrate the three different constructions, the following sections show how we can define certain compositions with the *Maybe*, reader, and list monads using the *prod*, *dorp* and *swap* constructions, respectively. We will also present some further examples using the same techniques to obtain compositions with other standard monads.

6.1 The *Maybe* datatype

The *Maybe* datatype, used in [12] to model a form of exception handling, is defined by:

```
data Maybe a = Just a | Nothing
```

There is a natural functor and monad structure corresponding to this datatype, given by the following declarations (proofs that these functions satisfy the appropriate laws are left as an exercise for the reader):

```
instance Functor Maybe where
```

```
  map f (Just x) = Just (f x)
```

```
  map f Nothing = Nothing
```

```
instance Premonad Maybe where
```

```
  unit = Just
```

```
instance Monad Maybe where
```

```
  join (Just m) = m
```

```
  join Nothing = Nothing
```

Our goal now is to show how to construct a new monad by composing the *Maybe* constructor with an arbitrary datatype. Using the *prod* construction, a composition of the form *Maybe* . *n* would require a function:

```
prod :: n (Maybe (n a)) → Maybe (n a).
```

However, using only the monad operators for *n*, there does not appear to be any reasonable way to define a suitable function with this type. (This could probably be proved formally using the same kind of techniques as in the proof in the appendix.) On the other hand, for a composition of the form *m* . *Maybe*, we require a function:

```
prod :: Maybe (m (Maybe a)) → m (Maybe a).
```

In this case, since we know something about the structure of objects constructed using *Maybe*, it is relatively easy to find a suitable definition for *prod*:

```
instance PComposable m Maybe where
```

```
  prod (Just m) = m
```

```
  prod Nothing = unit Nothing
```

Proofs that this definition satisfies the necessary laws P(1)–P(4) are given in Figure 10. In each case, except for law P(2), we split the proof into two case – one for values of the form *Nothing* and a second for values *Just* m.

<i>prod (map (mapC f) Nothing)</i>	<i>= prod Nothing</i>	<i>map</i>
	<i>= unit Nothing</i>	<i>prod</i>
	<i>= unit (map f Nothing)</i>	<i>map</i>
	<i>= map (map f) (unit Nothing)</i>	<i>law3</i>
	<i>= mapC f (unit Nothing)</i>	<i>mapC</i>
	<i>= mapC f (prod Nothing)</i>	<i>prod</i>
<i>prod (map (mapC f) (Just m))</i>	<i>= prod (Just (mapC f m))</i>	<i>map</i>
	<i>= mapC f m</i>	<i>prod</i>
	<i>= mapC f (prod (Just m))</i>	<i>prod</i>
<i>prod (unit m)</i>	<i>= prod (Just m)</i>	<i>unit</i>
	<i>= m</i>	<i>prod</i>
<i>prod (map unitC Nothing)</i>	<i>= prod Nothing</i>	<i>map</i>
	<i>= unit Nothing</i>	<i>prod</i>
<i>prod (map unitC (Just x))</i>	<i>= prod (Just (unitC x))</i>	<i>map</i>
	<i>= unitC x</i>	<i>prod</i>
	<i>= unit (unit x)</i>	<i>unitC</i>
	<i>= unit (Just x)</i>	<i>unit</i>
<i>prod (map joinP Nothing)</i>	<i>= prod Nothing</i>	<i>map</i>
	<i>= join (unit (prod Nothing))</i>	<i>law5</i>
	<i>= join (map prod (unit Nothing))</i>	<i>law3</i>
	<i>= joinP (unit Nothing)</i>	<i>joinP</i>
	<i>= joinP (prod Nothing)</i>	<i>prod</i>
<i>prod (map joinP (Just m))</i>	<i>= prod (Just (joinP m))</i>	<i>map</i>
	<i>= joinP m</i>	<i>prod</i>
	<i>= joinP (prod (Just m))</i>	<i>prod</i>

Figure 10: Proof of P(1)–P(4) for the *Maybe* monad

6.2 Monad Comprehension Syntax

For most of the monad constructions introduced in the following sections, it is convenient to work with the notation of *monad comprehensions*. Several functional programming languages, including Haskell and Gofer, provide a special syntax for *list comprehensions* which allow some list-based computations to be expressed very clearly and concisely. Lists form a monad (see Section 6.4 for details) and, noticing this, Wadler [14] showed how the comprehension notation could be generalized to an arbitrary monad. A comprehension is written using the notation $[exp \mid gs]$ where exp is an expression and gs is a list of *generators* i.e. expressions of the form $x \leftarrow e$. The meaning of a comprehension can be defined by translating it into a form using the standard monad operators:

$$\begin{aligned} [exp \mid x \leftarrow e] &= map (\backslash x \rightarrow exp) e && \text{mapComp} \\ [exp] &= unit exp && \text{unitComp} \\ [exp \mid gs, hs] &= join [[exp \mid hs] \mid gs] && \text{joinComp} \end{aligned}$$

The first equation can be considered as a way of defining *map* using the comprehension notation, with $map f [exp \mid gs] = [f exp \mid gs]$. The second equation gives the meaning of a comprehension with an empty sequence of qualifiers. In the last equation, the variables gs and hs range over (possibly empty) sequences of generators. Using the monad laws, we can show that the way in which these rules are used to find a translation of a monad comprehension does not have any effect on the meaning of the result. See [14] for further details.

For convenience, we will also use the following laws about monad comprehensions, each of which can be derived from the definitions above and the monad laws.

$$\begin{aligned} [x \mid x \leftarrow xs] &= xs && \text{compId} \\ [f x \mid x \leftarrow map g e] &= [f (g x) \mid x \leftarrow e] && \text{compMap} \\ [f x \mid x \leftarrow unit e] &= [f e] && \text{compUnit} \\ [exp \mid x \leftarrow join e] &= [exp \mid z \leftarrow e, x \leftarrow z] && \text{compJoin} \end{aligned}$$

For example, the first of these is just another way of writing the functor law (1), while the second follows directly (2) and can be used to justify equalities such as:

$$\begin{aligned} [h x y \mid x \leftarrow map f u, y \leftarrow map g v] &= [h (f x) (g y) \mid x \leftarrow u, y \leftarrow v] \\ [exp \mid m \leftarrow map f n, x \leftarrow m] &= [exp \mid m \leftarrow n, x \leftarrow f m] \end{aligned}$$

In a similar way, *compUnit* can also be used in the justification of laws like:

$$[h x y \mid x \leftarrow unit u, y \leftarrow unit v] = [h u v].$$

Finally, we mention that the *compId* and *compJoin* can be used together to give a definition of the *join* operator using the comprehension notation:

$$join e = [x \mid x \leftarrow join e] = [z \mid z \leftarrow e, x \leftarrow z].$$

6.3 Reader monads

A reader monad is described by a constructor of the form $(r \rightarrow)$ mapping each type a to the function type $r \rightarrow a$. We refer to computations in this monad as readers because they can read the value passed in the parameter of type r , but they cannot change that value. The functor, premonad and monad structures for readers are given by:

```
instance Functor (r →) where
  map f g = f . g

instance Premonad (r →) where
  unit x y = x

instance Monad (r →) where
  join f x = f x x
```

The following declaration can be used to compose an arbitrary monad n with a reader monad $(r \rightarrow)$ using the *dorp* construction:

```
instance DComposable (r →) n where
  dorp m r = [g r | g ← m r]
```

Proofs that this definition satisfies the laws D(1)–D(4) are included in Figure 11.

6.4 The *List* monad

Lists are one of the most widely used monads in functional programming. The structure of the list monad is captured by the type constructor *List*⁴, together with the following instance declarations:

```
instance Functor List where
  map f []      = []
  map f (x : xs) = f x : map f xs

instance Premonad List where
  unit x = [x]

instance Monad List where
  join = foldr (++) []
```

The *foldr* function and the list append operator $(++)$ used here are both taken from the standard prelude.

This time, we will use the *swap* construction to obtain a composition of a monad m with the *List* monad (we will see shortly that this construction only yields a composite

⁴In fact, the concrete syntax of Gofer currently requires that we write this constructor as $[]$, but we use the notation *List* here for clarity. The type *List* a is the same as the type $[a]$ in standard Haskell notation.

$dorp (mapC (map f) m) r$	$= [g r \mid g \leftarrow map (map (map f)) m r]$ $= [g r \mid g \leftarrow map (map f) (m r)]$ $= [map f g r \mid g \leftarrow m r]$ $= [f (g r) \mid g \leftarrow m r]$ $= map f [g r \mid g \leftarrow m r]$ $= map f (dorp m r)$ $= map (map f) (dorp m) r$ $= mapC f (dorp m) r$	<i>dorp, mapC</i> <i>map</i> <i>compMap</i> <i>map</i> <i>compMap</i> <i>dorp</i> <i>map</i> <i>mapC</i>
$dorp (unitC f) r$	$= [g r \mid g \leftarrow unit (unit f) r]$ $= [g r \mid g \leftarrow unit f]$ $= [f r]$ $= unit (f r)$ $= map unit f r$	<i>dorp, unitC</i> <i>unit</i> <i>compUnit</i> <i>unitComp</i> <i>map</i>
$dorp (mapC unit f) r$	$= [g r \mid g \leftarrow map (map unit) f r]$ $= [g r \mid g \leftarrow map unit (f r)]$ $= [unit g r \mid g \leftarrow f r]$ $= [g \mid g \leftarrow f r]$ $= f r$	<i>dorp, mapC</i> <i>map</i> <i>compMap</i> <i>unit</i> <i>compId</i>
$dorp (joinD f) r$	$= [g r \mid g \leftarrow map join (dorp f) r]$ $= [g r \mid g \leftarrow join (dorp f r)]$ $= [g r \mid h \leftarrow dorp f r, g \leftarrow h]$ $= [g r \mid h \leftarrow [k r \mid k \leftarrow f r], g \leftarrow h]$ $= [g r \mid h \leftarrow f r, g \leftarrow h r]$ $= join [[g r \mid g \leftarrow h r] \mid h \leftarrow f r]$ $= join [dorp h r \mid h \leftarrow f r]$ $= join [h r \mid h \leftarrow map dorp (f r)]$ $= join [h r \mid h \leftarrow map (map dorp) f r]$ $= join [h r \mid h \leftarrow mapC dorp f r]$ $= join (dorp (mapC dorp f) r)$ $= map join (dorp (mapC dorp f)) r$ $= joinD (mapC dorp f) r$	<i>dorp, joinD</i> <i>map</i> <i>compJoin</i> <i>dorp</i> <i>mapComp</i> <i>joinComp</i> <i>dorp</i> <i>compMap</i> <i>map</i> <i>mapC</i> <i>dorp</i> <i>map</i> <i>joinD</i>

Figure 11: Proof of D(1)–D(4) for reader monads

monad if m has a certain commutativity property):

```
instance SComposable m List where
  swap []          = unit []
  swap (x : xs)   = [y : ys | y ← x, ys ← swap xs]
```

The proofs for S(1)–S(3) in Figure 12 hold for any monad m . Structural induction is required for the proofs of S(1) and S(3); we have not shown the case for \perp values which follow directly from the use of pattern matching in the definition of the monad operators.

The proof of law S(4) is rather more involved. To begin with, we need to introduce the auxiliary functions:

```
prod = map join . swap
dorp = join . map swap
```

A quick calculation (or, if you prefer, proof by induction) based on the definition of *swap* allows us to use the following definition for the *prod* function:

```
prod []          = unit []
prod (x : xs)   = [u ++ v | u ← x, v ← prod xs]
```

The proof of S(4), by structural induction, is presented in Figure 13.

There are two steps in the proof which require further explanation. The first is a lemma describing the way that *swap* distributes over $(++)$:

```
swap (xs ++ ys) = [x ++ y | x ← swap xs, y ← swap ys]
```

This law holds for any monad m (i.e. the monad in which the comprehension is interpreted) and can be proved using the simple structural induction on xs in Figure 14.

Notice that we made use of (yet another) variant of *compMap* :

```
[f x y | x ← u, t ← [h y z | y ← v, z ← w]]
= [f x (h y z) | x ← u, y ← v, z ← w]
```

The other important step is the use of a commutativity property that enables us to swap two of the generators in a comprehension. We require that the following law is satisfied :

```
[r ++ s | v ← x, r ← swap v, u ← prod xs, s ← swap u]
= [r ++ s | v ← x, u ← prod xs, r ← swap v, s ← swap u]
```

One way to ensure that this law holds is to insist that m is a commutative monad – i.e. that it satisfies the law:

```
[(x, y) | x ← u, y ← v] = [(x, y) | y ← v, x ← u].
```

<i>swap</i> (<i>mapC</i> <i>f</i> [])	
= <i>swap</i> (<i>map</i> (<i>map</i> <i>f</i>) [])	<i>mapC</i>
= <i>swap</i> []	<i>mapC</i> , <i>map</i>
= <i>unit</i> []	<i>swap</i>
= <i>unit</i> (<i>map</i> <i>f</i> [])	<i>map</i>
= <i>map</i> (<i>map</i> <i>f</i>) (<i>unit</i> [])	<i>law3</i>
= <i>mapC</i> <i>f</i> (<i>swap</i> [])	<i>mapC</i> , <i>swap</i>
<i>swap</i> (<i>mapC</i> <i>f</i> (<i>x</i> : <i>xs</i>))	
= <i>swap</i> (<i>map</i> (<i>map</i> <i>f</i>) (<i>x</i> : <i>xs</i>))	<i>mapC</i>
= <i>swap</i> (<i>map</i> <i>f</i> <i>x</i> : <i>map</i> (<i>map</i> <i>f</i>) <i>xs</i>)	<i>map</i>
= [<i>y</i> : <i>ys</i> <i>y</i> ← <i>map</i> <i>f</i> <i>x</i> , <i>ys</i> ← <i>swap</i> (<i>map</i> (<i>map</i> <i>f</i>) <i>xs</i>)]	<i>swap</i>
= [<i>y</i> : <i>ys</i> <i>y</i> ← <i>map</i> <i>f</i> <i>x</i> , <i>ys</i> ← <i>map</i> (<i>map</i> <i>f</i>) (<i>swap</i> <i>xs</i>)]	induction
= [<i>f</i> <i>y</i> : <i>map</i> <i>f</i> <i>ys</i> <i>y</i> ← <i>x</i> , <i>ys</i> ← <i>swap</i> <i>xs</i>]	<i>compMap</i>
= [<i>map</i> <i>f</i> (<i>y</i> : <i>ys</i>) <i>y</i> ← <i>x</i> , <i>ys</i> ← <i>swap</i> <i>xs</i>]	<i>compMap</i>
= <i>map</i> (<i>map</i> <i>f</i>) [<i>y</i> : <i>ys</i> <i>y</i> ← <i>x</i> , <i>ys</i> ← <i>swap</i> <i>xs</i>]	<i>map</i>
= <i>mapC</i> <i>f</i> (<i>swap</i> (<i>x</i> : <i>xs</i>))	<i>mapC</i> , <i>swap</i>
<i>swap</i> (<i>unit</i> <i>x</i>)	
= <i>swap</i> [<i>x</i>]	<i>unit</i>
= [<i>y</i> : <i>ys</i> <i>y</i> ← <i>x</i> , <i>ys</i> ← <i>swap</i> []]	<i>swap</i>
= [<i>y</i> : <i>ys</i> <i>y</i> ← <i>x</i> , <i>ys</i> ← <i>unit</i> []]	<i>swap</i>
= [[<i>y</i>] <i>y</i> ← <i>x</i>]	<i>compUnit</i>
= [<i>unit</i> <i>y</i> <i>y</i> ← <i>x</i>]	<i>unitComp</i>
= <i>map</i> <i>unit</i> <i>x</i>	<i>mapComp</i>
<i>swap</i> (<i>map</i> <i>unit</i> [])	
= <i>swap</i> []	<i>map</i>
= <i>unit</i> []	<i>swap</i>
<i>swap</i> (<i>map</i> <i>unit</i> (<i>x</i> : <i>xs</i>))	
= <i>swap</i> (<i>unit</i> <i>x</i> : <i>map</i> <i>unit</i> <i>xs</i>)	<i>map</i>
= [<i>y</i> : <i>ys</i> <i>y</i> ← <i>unit</i> <i>x</i> , <i>ys</i> ← <i>swap</i> (<i>map</i> <i>unit</i> <i>xs</i>)]	<i>swap</i>
= [<i>y</i> : <i>ys</i> <i>y</i> ← <i>unit</i> <i>x</i> , <i>ys</i> ← <i>unit</i> <i>xs</i>]	induction
= [<i>x</i> : <i>xs</i>]	<i>compUnit</i>
= <i>unit</i> (<i>x</i> : <i>xs</i>)	<i>unitComp</i>

Figure 12: Proof of S(1)–S(3) for the list monad

<i>prod (map dorp [])</i>	
<i>= prod []</i>	<i>map</i>
<i>= unit []</i>	<i>prod</i>
<i>= swap []</i>	<i>swap</i>
<i>= join (unit (swap []))</i>	<i>law5</i>
<i>= join (map swap (unit []))</i>	<i>law3</i>
<i>= dorp (unit [])</i>	<i>dorp</i>
<i>= dorp (prod [])</i>	<i>prod</i>
<i>prod (map dorp (x : xs))</i>	
<i>= prod (dorp x : map dorp xs)</i>	<i>map</i>
<i>= [r++s r ← dorp x, s ← prod (map dorp xs)]</i>	<i>prod</i>
<i>= [r++s r ← dorp x, s ← dorp (prod xs)]</i>	<i>induction</i>
<i>= [r++s r ← dorp x, s ← join (map swap (prod xs))]</i>	<i>dorp</i>
<i>= [r++s r ← dorp x, u ← map swap (prod xs), s ← u]</i>	<i>compJoin</i>
<i>= [r++s r ← dorp x, u ← prod xs, s ← swap u]</i>	<i>compMap</i>
<i>= [r++s r ← join (map swap x), u ← prod xs, s ← swap u]</i>	<i>dorp</i>
<i>= [r++s v ← map swap x, r ← v, u ← prod xs, s ← swap u]</i>	<i>compJoin</i>
<i>= [r++s v ← x, r ← swap v, u ← prod xs, s ← swap u]</i>	<i>compMap</i>
<i>= [r++s v ← x, u ← prod xs, r ← swap v, s ← swap u]</i>	<i>commute</i>
<i>= join [[r++s r ← swap v, s ← swap u]</i>	<i>joinComp</i>
<i> v ← x, u ← prod xs]</i>	
<i>= join [swap (v++u) v ← x, u ← prod xs]</i>	<i>lemma</i>
<i>= join (map swap [v++u v ← x, u ← prod xs])</i>	<i>compMap</i>
<i>= dorp [v++u v ← x, u ← prod xs]</i>	<i>dorp</i>
<i>= dorp (prod (x : xs))</i>	<i>prod</i>

Figure 13: Proof of S(4) for the list monad

$swap ([] ++ ys)$	
$= swap\ ys$	(++)
$= [vs \mid vs \leftarrow swap\ ys]$	<i>compId</i>
$= [[] ++ vs \mid vs \leftarrow swap\ ys]$	(++)
$= [us ++ vs \mid us \leftarrow unit\ [],\ vs \leftarrow swap\ ys]$	<i>compUnit</i>
$= [us ++ vs \mid us \leftarrow swap\ [],\ vs \leftarrow swap\ ys]$	<i>swap</i>
$swap ((x : xs) ++ ys)$	
$= swap\ (x : xs ++ ys)$	(++)
$= [z : zs \mid z \leftarrow x,\ zs \leftarrow swap\ (xs ++ ys)]$	<i>swap</i>
$= [z : zs \mid z \leftarrow x,\ zs \leftarrow [us ++ vs \mid us \leftarrow swap\ xs,\ vs \leftarrow swap\ ys]]$	induction
$= [z : us ++ vs \mid z \leftarrow x,\ us \leftarrow swap\ xs,\ vs \leftarrow swap\ ys]$	<i>compMap</i>
$= [us ++ vs \mid us \leftarrow swap\ (x : xs),\ vs \leftarrow swap\ ys]$	<i>swap</i>

Figure 14: Proof of $swap\ (xs ++ ys) = [x ++ y \mid x \leftarrow swap\ xs,\ y \leftarrow swap\ ys]$.

The identity monad, the set monad and reader monads all have this property. On the other hand, the list monad, *Maybe* monad and state monad are non-commutative.

In summary, if m is a commutative monad, then $SComp\ m\ List$ is a monad. Since we require only a single special case of the commutativity property, it is possible that we may be able to relax the restriction to commutative monads to some degree. However, it is not possible to remove the restrictions altogether. In particular, our constructions cannot be used to compose the *List* monad with itself.

The importance of commutative monads has also been recognized in other situations. For example, in [6], monads satisfying the commutativity axiom are used to capture explicit parallel execution of programs written in a monadic style.

6.4.1 Composing a monad with itself

In Section 2 we commented that, just by looking at the types involved, it was clear that a definition $join = join_M . join_N$ could not be used to form a composition of two monads M and N . The argument was that, while the function $join_N$ produces a result with type of the form $N\ a$, the function $join_M$ expects a value with a type of the form $M\ (M\ a)$. In the general case, these types do not match. However, as the reader may have realized, this argument fails in the special case where $M = N$.

In fact, although it has the required type, we still cannot use the function $join_M . join_M$, because it does not satisfy the monad laws. For example, using the list monad, we can demonstrate a direct counter example to law (6):

$$(join . join . map\ unitC)\ [[]] = [] \neq [[]] = id\ [[]]$$

So, in general, the composition of a monad with itself cannot be treated as a special

case; we still need to use one of the constructions described above.

6.4.2 Comparison with 'Combining monads'

Our results restrict composition with *List* to commutative monads. In contrast, King and Wadler [8] give a slightly different construction which seems more general, avoiding any restrictions on the choice of monads that can be composed with *List*. Unfortunately, this promise of a more general construction turns out to be something of a mirage; although developed in a rather different manner, their approach turns out to be equivalent to our *prod* construction and does, in fact, require some form of commutativity.

The construction given in [8], is based on the definitions:

$$\begin{aligned} (\otimes) &:: \text{Monad } m \Rightarrow m [a] \rightarrow m [a] \rightarrow m [a] \\ a \otimes b &= [x ++ y \mid x \leftarrow a, y \leftarrow b] \\ \text{prod} &:: \text{Monad } m \Rightarrow [m [a]] \rightarrow m [a] \\ \text{prod} &= \text{foldr } (\otimes) (\text{unit } []) \end{aligned}$$

This *prod* function is used to define $\text{join} = \text{join} . \text{map prod}$, the proof that this defines a monad structure for *m* composed with the list monad resting on a number of properties of *prod*, including:

$$\text{prod} . \text{map } (\text{join} . \text{map prod}) = \text{join} . \text{map prod} . \text{prod}.$$

However, using the following, somewhat contrived counter example, we find that this law does not always hold:

```
? (prod . map (join . map prod)) [[[[[1], [3]]], [[4]], [[2]]]]
[[1, 4], [1, 2], [3, 4], [3, 2]]
? (join . map prod . prod) [[[[[1], [3]]], [[4]], [[2]]]]
[[1, 4], [3, 4], [1, 2], [3, 2]]
```

Of course, none of what we have said here proves that it is impossible to construct monads by composition of arbitrary monads with the list monad; all we know is that it is not possible using the constructions described in this report.

6.5 Some additional monad compositions

The following sections deal with some further examples, providing compositions with standard monads. In each case, we use the *swap* construction to obtain the composition, but detailed proofs are not included since they are very similar to those in previous sections⁵.

⁵Detailed proofs of results in this section are included, as program comments, in a Gofer script containing an executable version of the definitions in this report. See comments on the Page 1 for details of availability.

6.5.1 Writer monads

A writer monad can be used for describing programs that produce both output and a return value. It pays to take a general approach, allowing the type of values used as output to be provided as a parameter to the monad constructor, rather than committing ourselves to a particular output type at an early stage. The following declarations define a constructor *Writer* and, assuming an output type *s*, a corresponding functor and monad structure for each *Writer s* constructor:

```
data Writer s a = Result s a

instance Functor (Writer s) where
  map f (Result s a) = Result s (f a)

instance Monoid s  $\Rightarrow$  Premonad (Writer s) where
  unit = Result zero

instance Monoid s  $\Rightarrow$  Monad (Writer s) where
  join (Result s (Result t x)) = Result (add s t) x

  write           :: s  $\rightarrow$  Writer s ()
  write msg      = Result msg ()
```

The *write* function defined here is used to perform output; the argument *msg* is returned as the output and the value returned is just *()*, the unit value.

The values *zero* and *add* in the definitions above represent the null output, and the sequencing of one output after another. To establish the monad laws for *Writer s* constructors, we need to insist that *add* is associative with *zero* as both a left and right identity. In other words, we require that these values form a monoid. This can be captured by the following class declaration:

```
class Monoid s where
  zero :: s
  add  :: s  $\rightarrow$  s  $\rightarrow$  s
```

Two obvious choices for output types, each of which forms a monad, are lists and functions, as described by the following instance declarations.

```
instance Monoid [a] where
  zero = []
  add  = (++)

instance Monoid (a  $\rightarrow$  a) where
  zero = id
  add  = (.)
```

For the example in Section 7, we will use lists of strings (i.e. values of type *String*, each corresponding to a line of output) for output values. In practice, an output type

of the form $String \rightarrow String$ might be more sensible, allowing tree like structures to be printed in linear time rather than in time quadratic in the size of the tree. Both of these possibilities are permitted by the definitions above.

Incidentally, we will also mention that the monad *Writer Int*, assuming the monoid structure:

```
instance Monoid Int where
  zero = 0
  add  = (+)
```

gives us another example of a commutative monad (see Section 6.4). This might be used, for example, in a simple profiler, counting the number of times that particular tasks are carried out and using the command *write 1* to increment the counter at suitable points in the program.

Writer monads can be composed with arbitrary monads using the following definition of *swap*:

```
instance Monoid s => SComposable m (Writer s) where
  swap (Result s m) = [Result s a | a <- m]
```

6.5.2 The *Error* monad

As a simple variation on the *Maybe* monad described in Section 6.1, it is often useful to be able to return some form of error message when an exception occurs. This can be described by the *Error* monad:

```
data Error a = Ok a | Error String

instance Functor Error where
  map f (Ok x)      = Ok (f x)
  map f (Error msg) = Error msg

instance Premonad Error where
  unit = Ok

instance Monad Error where
  join (Ok x)      = x
  join (Error msg) = Error msg
```

As with *Maybe*, we can use the *prod* construction to obtain a composition of *Error* (on the left) with any monad *m*. Suitable definitions and proofs can be obtained by replacing each reference to *Just* with *Ok* and each *Nothing* with a value of the form *Error msg* in Section 6.1. Alternatively, we can use the *swap* construction, based on the definition:

```
instance SComposable m Error where
  swap (Ok m)      = map Ok m
  swap (Error msg) = unit (Error msg)
```

6.5.3 The *Tree* monad

The final example we will consider is tree monad, defined by the constructor *Tree*, and mapping a type *a* to the type *Tree a* of binary trees with leaf values of type *a*. The monad structure is given by the following declarations:

```
data Tree a = Leaf a | Tree a : $\wedge$ : Tree a

instance Functor Tree where
  map f (Leaf x)    = Leaf (f x)
  map f (lx : $\wedge$ : rx) = map f lx : $\wedge$ : map f rx

instance Premonad Tree where
  unit = Leaf

instance Monad Tree where
  join (Leaf m)    = m
  join (lm : $\wedge$ : rm) = join lm : $\wedge$ : join rm
```

Composition of the *Tree* monad with a monad *m* can be described using the *swap* construction with:

```
instance SComposable m Tree where
  swap (Leaf m)    = [Leaf m | x  $\leftarrow$  m]
  swap (lm : $\wedge$ : rm) = [lx : $\wedge$ : rx | lx  $\leftarrow$  swap lm, rx  $\leftarrow$  swap rm]
```

As in the case of list monads, the proof of S(4) depends on a commutativity property for *m*. The proofs for S(1)–S(3) however (and hence, the proofs of the monad laws (4)–(6) for *SComp m Tree*) do not require any special properties of *m*.

7 A simple example: an evaluator

In this section, we show how two of the monad constructions introduced above can be used in a small, but practical, example – an expression evaluator. To make the example more interesting we will use an environment mapping variables to values (with a careful treatment of unbound variables) and we will provide a simple trace facility.

We will base the expression evaluator on the following types representing values, variable names, environments (i.e. mappings from variable names to values) and expressions:

```
type Value = Int
type Name = String
type Env = [(Name, Value)]
data Expr = Const Value | Var Name | Expr :+ : Expr | Trace String Expr
```

To keep this example short, we have restricted ourselves to a very simple expression language, allowing only constants, variables, addition and a simple trace mechanism.

The evaluator itself will need to access variables bound in an environment, may produce output if the trace facility is used, and requires some form of error handling to deal with unbound variables. This could be captured by defining a monad structure for the type constructor:

```
type M a = Env → Writer [String] (Error a),
```

working out suitable versions of the monad operators for this particular type. If we're going to be really fussy, we should also verify the monad laws for whatever definitions we choose for these operators. Given that the definition of M is quite complex, this is likely to require a long, error-prone, and largely uninspiring calculation.

The alternative, using the tools introduced in this report, is to recognize that the definition of M can be expressed as a composition of monads. Assuming that all the components are themselves monads, our results mean that we can use this approach without requiring any further proofs:

```
type M a = DComp (Env →) (SComp (Writer [String]) Error) a
```

In all honesty, we have to admit that this looks rather ugly. However, in a language designed from scratch to support this kind of work, we might reasonably expect to be able to define M in a less clumsy manner, using an equation of the form:

$$M = (Env \rightarrow) . Writer [String] . Error$$

We will also use the following functions, defined in terms of the general *left* and *right* functions introduced in Section 5.4, to embed computations in each of the component monads into the full monad M :

```
inError  :: Error a → M a
inError  = right . right

inReader :: (Env → a) → M a
inReader = left

inWriter :: Writer [String] a → M a
inWriter = right . left
```

It is possible to package these functions in a more general way using the overloading mechanism (all have types of the form $N a \rightarrow M a$ for some constructor N), but we will not consider this any further here.

Before we can give the definition of the evaluator, we need a simple utility function to determine the value bound to a particular variable in a given environment. We will use the *Maybe* type to enable us to deal with the two cases, depending on whether the variable is bound or not in the given environment.

```
lookup      :: Name → Env → Error Value
lookup x ((n, v) : env) = if x == n then Ok v else lookup x env
lookup x []           = Error ("unbound variable " ++ x)
```

The complete evaluator is defined as follows, with a single case for each possible form of expression:

```

eval      :: Expr → M Value
eval (Const v) = [unit v]
eval (Var n)   = [ x | r ← inReader (lookup n), x ← inError r ]
eval (e :+: f) = [ x + y | x ← eval e, y ← eval f ]
eval (Trace m e) = [ x | x ← eval e,
                    () ← inWriter (write [m ++ " = " ++ show r])]

```

Finally, we provide a function that can be used to execute a computation in the M monad, using a given environment value and returning a string as its result:

```

result      :: Text a ⇒ M a → Env → String
result m env = unlines (["Output : "] ++ s ++ ["Result : " ++ val])
               where Result s x = open (open m env)
                     val      = case x of
                                   Ok x      → show x
                                   Error msg → msg

```

For example, using the programs described here in the Gofer interpreter with the expression `testExpr = Trace "sum" (Const 1 :+: Const 2) :+: Var "x"` gives the following results:

```

? result (eval testExpr) []
Output:
sum = 3
Result: unbound variable x

? result (eval testExpr) [("x",42)]
Output:
sum = 3
Result: 45

```

8 Other ways to combine monads

Having spent so much time concentrating on the composition of monads, it is important to point out that there are other methods that can be used to combine monads.

In some cases, a composition of monads is not suitable because it implies a certain level of independence between the components than may not be desired. For example, one common application of monads is to model state-based computations using *state transformers*; i.e. functions taking an initial state and returning a pair containing a final state and a return value of some type. This is described, for example in [5], using a type constructor:

```

data State s a = ST (s → (s, a))

```

The type s used here gives the type of values used for the state while a represents the type of return values. The standard functor and monad structure are given by the declarations:

```

instance Functor (State s) where
  map f (ST st) = ST (\s → let (s', x) = st s in (s', f x))

instance Premonad (State s) where
  unit x = ST (\s → (s, x))

instance Monad (State s) where
  join (ST m) = ST (\s → let (s', ST m') = m s in m' s')

```

Notice that, as a type constructor, *State s* is isomorphic to the composition of a reader monad with a writer monad, $DComp (s \rightarrow) (Writer s)$. However, the monad structure associated with these two constructors is very different. For example, the *unit* function for the state monad returns the initial state value unchanged, while the *unit* operator for the composition returns the *zero* of s (with the additional constraint that s is a monoid). Perhaps there is another general construction for combining two monads that could be applied to $(s \rightarrow)$ and *Writer s* to obtain a state monad?

As a more interesting example, we can combine a simple state monad with an arbitrary using the following definitions, adapted from [5]:

```

data StateM m s a = STM (s → m (s, a))

instance Monad m ⇒ Functor (StateM m s) where
  map f (STM xs) = STM (\s → [(s', f x) | (s', x) ← xs s])

instance Monad m ⇒ Premonad (StateM m s) where
  unit x = STM (\s → unit (s, x))

instance Monad m ⇒ Monad (StateM m s) where
  join (STM xss) = STM (\s → [(s'', x) | (STM xs, s') ← xss s,
                                         (s'', x) ← xs s'])

```

For example, a monad of this form is often used in the construction of a combinator parser [4, 14, 15]. In this kind of application, we use the stream of tokens to be parsed as the state. Possible choices for the monad m include:

- The *List* monad, to deal with ambiguous grammars.
- The *Maybe* monad, to support backtracking parsers.
- The *Error* monad, to provide error handling.
- A composition of *Maybe* and *Error*, to allow both error handling and backtracking.

Looking only at the types involved, we can express *StateM m s* as a composition of a reader, the monad m , and a writer:

$$StateM\ m\ s = (s \rightarrow) . m . Writer\ s.$$

However, using the definitions in this report, the monad structure for the composition on the right is very different from the monad structure for the constructor *StateM m s* on the left given by the declarations above. Perhaps this combination of a state monad with another arbitrary monad can be described as an instance of some more general construction for monad combination?

Another way to view many of the examples in this report is as a generalization of the concept of a monad to include a ‘hole’ that can be filled with another monad to obtain suitable combinations of features. For example, the constructors *PComp m Maybe*, *DComp (r →) m*, and *StateM m s* are all examples of this with a hole represented by the parameter *m*. A similar idea motivates the recent work of Steele [13] except that the holes are built into so-called *pseudomonad* operators rather than individual constructors. The pseudomonad operators are difficult to express properly in the Haskell type system and appear to require some form of existential or recursive typing. Fortunately, it is quite easy to express these operators using constructor classes:

```
class Premonad p ⇒ Pseudomonad p where
  pbind      :: Monad m ⇒ p a → (a → m (p b)) → m (p b)
  pjoin      :: Monad m ⇒ p (m (p a)) → m (p a)
  pjoin m    = m 'pbind' id
  m 'pbind' f = pjoin (map f m)
```

The *pbind* function is included for those familiar with [13] and provides a pseudomonad version of the monadic *bind* operator mentioned in Section 2. The *pjoin* function is not used by Steele, but is included because of its close relationship to the *join*-based formulation of monads used in this report. Strictly speaking, only one of the *pbind* and *pjoin* functions is actually required to define an instance of the *Pseudomonad* class; the default definitions provided by the last two lines of the class declaration show how each can be defined in terms of the other.

The types of both *pjoin* and *pbind* are very similar to the types of the corresponding monad operators, except that they also make use of an additional parameter representing an arbitrary monad *m*. Replacing this parameter with the identity monad gives the familiar monad operators. On the other hand, this extra parameter enables us to describe composition in a simple, elegant manner. Starting with the definition of a new composer:

```
data Comp f g x = CC (f (g x))
instance Composer Comp where
  open (CC x) = x
  close      = CC
```

the composition of an arbitrary pseudomonad *p* with an arbitrary monad *m*, yielding a composite monad *Comp m p*, can be described by the following instance declaration:

```
instance (Pseudomonad p, Monad m) ⇒ Monad (Comp m p) where
  join = wrap (join . map pjoin)
```

In this way, we can build up a chain of pseudomonads, p_1, \dots, p_n , with the final hole plugged by a monad m :

$$\text{Comp } (\dots (\text{Comp } (\text{Comp } m \ p_1) \ p_2) \ \dots) \ p_n,$$

or, using an infix dot instead of *Comp*, just:

$$(\dots((m \cdot p_1) \cdot p_2) \ \dots \cdot p_n).$$

This certainly seems like a promising approach, and we hope to investigate its relationship with the constructions used here more fully in future work.

9 Conclusions

We have presented three different constructions that can be used to compose monads and shown how these can be encoded and used in practical programming problems to provide a combination of the features offered by the component monads. The proofs required to establish the monad laws for each composition require only simple equational reasoning (sometimes with structural induction), although they can be a little long. On the other hand, we have already developed a small library describing compositions with certain standard monads which can be extended to include other monads as necessary. These results can already be used to construct new monads without any further proof obligations.

One surprising aspect of this work is the need to restrict our attention to commutative monads in compositions with the *List* monad. This additional property is necessary only to establish the monad law (7), sometimes referred to as the associative law for monads. In practice, there are several examples where the basic framework suggested by the types of the monad operators is useful in practical programming examples, even though the corresponding monad laws are not all satisfied. Examples of this include the strictness monad in [14], state transformers in [10] and ‘composable contexts’ in [7]. In a similar way, we expect the composition of arbitrary monads with *List* may still be useful in practical programming applications, even though the associativity law does not always hold. Perhaps functional programmers will be prepared to sacrifice the algebraic properties of a full monad, gaining wider application of the techniques of the monadic style of programming as a reward.

Acknowledgments

The work described in this report grew out of a conversation between the two authors on the `comp.lang.functional` newsgroup in December 1992; we are both grateful for access to this forum. Mark Jones’s work is supported in part by a grant from ARPA, contract number N00014-91-J-4043. Luc Duponcheel, much of whose work has been carried out as a recreation and hobby, would like to thank Maritza and Jos for their

support, and Alcatel Bell for the opportunity to spend a little time on this kind of research oriented work.

References

- [1] L. Duponcheel. A short note on monad and monad morphism composition. Alcatel Bell Telephone, Antwerp, B-2018, Belgium. Manuscript, September 1993.
- [2] C. Hall, K. Hammond, W. Partain, S.L. Peyton Jones and P. Wadler. The Glasgow Haskell compiler: a retrospective. In *Proceedings of the 1992 Glasgow Workshop on Functional Programming*, Ayr, Scotland, July 1992. Springer Verlag Workshops in computing series.
- [3] P. Hudak, S.L. Peyton Jones and P. Wadler (eds.). Report on the programming language Haskell, version 1.2. *ACM SIGPLAN notices*, 27, 5, May 1992.
- [4] G. Hutton. Higher-order functions for parsing. *Journal of Functional Programming*, volume 2, part 3, July 1992.
- [5] M.P. Jones. A system of constructor classes: overloading and implicit higher-order polymorphism. In *Proceedings of the 6th ACM conference on Functional Programming Languages and Computer Architecture*, Copenhagen, Denmark, ACM Press, June 1993.
- [6] M. Jones and P. Hudak. Implicit and Explicit Parallel Programming in Haskell. Yale University, Department of Computer Science, Research Report YALEU/DCS/RR-982, August 1993.
- [7] R.B. Kieburtz, B. Agapiev and J. Hook. Three Monads for Continuations. In *Proceedings of the ACM SIGPLAN Workshop on Continuations*, Department of Computer Science, Stanford University, Report No. STAN-CS-92-1462, June 1992.
- [8] D.J. King and P. Wadler. Combining monads. In *Proceedings of the Fifth Annual Glasgow Workshop on Functional Programming*, Ayr, Scotland, Springer Verlag Workshops in Computer Science, 1992.
- [9] E. Moggi. Computational lambda-calculus and monads. *IEEE Symposium on Logic in Computer Science*, Asilomar, California, 1989.
- [10] M. Odersky, D. Rabin and P. Hudak. Call by name, assignment, and the lambda calculus. In *20th Annual Symposium on Principles of Programming Languages*, Charleston, South Carolina, 1993.
- [11] S.L. Peyton Jones and P. Wadler. Imperative functional programming. In *20th Annual Symposium on Principles of Programming Languages*, Charleston, South Carolina, 1993.

- [12] M. Spivey. A functional theory of exceptions. *Science of Computer Programming*, 14(1), 1990.
- [13] G.L. Steele Jr. Building interpreters by composing monads. To appear in proceedings of *21st annual ACM SIGPLAN-SIGACT symposium on Principles of Programming Languages*, Portland, Oregon, January 1994
- [14] P. Wadler. Comprehending Monads. *Mathematical Structures in Computer Science*, 2, 4 (December 1992). An earlier version appears in *Conference on LISP and Functional Programming*, Nice, France, June 1990.
- [15] P. Wadler. The essence of functional programming. In *19th Annual Symposium on Principles of Programming Languages*, Santa Fe, New Mexico, January 1992.

Appendix: On the non-existence of a natural join

In Section 3, we commented that, in a certain sense, it is impossible to construct a *join* function for the composition of two monads using only the operators of the component monads. We will now sketch a proof of this claim.

A full justification of the approach that we have used goes a little beyond the scope of this paper. In particular, we need to recognize the fact that, from a categorical perspective, the type of the *map* function involves two different kinds of function – arrows between objects and arrows between arrows. These distinctions are lost in functional languages like Haskell or Gofer. Nevertheless, we believe that this result, as well as the techniques used in its proof, are likely to be of interest to some readers, and we have therefore decided to outline some of the details in this appendix.

Suppose that we are working in a category with two monads given by endofunctors M and N . In the most general case, the well-typed terms that can be constructed setting using only the functors and the natural transformations *unit* and *join* for the two monads are precisely those terms which can be obtained using the following set of typing rules:

$$\begin{array}{c}
 a \xrightarrow{id} a \\
 \\
 \frac{b \xrightarrow{f} c \quad a \xrightarrow{f} b}{a \xrightarrow{f \cdot g} c} \\
 \\
 \frac{a \xrightarrow{f} b}{M a \xrightarrow{map_M f} M b} \quad \frac{a \xrightarrow{f} b}{N a \xrightarrow{map_N f} N b} \\
 \\
 a \xrightarrow{unit_M} M a \quad a \xrightarrow{unit_N} N a \\
 \\
 M (M a) \xrightarrow{join_M} M a \quad N (N a) \xrightarrow{join_N} N a
 \end{array}$$

(Of course, this assumes only the very basic properties of a category. Richer categories may also have operations for forming products, sums, exponentials etc. In addition, we have no guarantee of full abstraction; there may well be arrows in the underlying category that cannot be obtained by these rules. After all, our main aim is to use exactly this kind of arrow (i.e. a *prod*, *dorp* or *swap* function) to construct the composition!)

Using this characterization, we will show that there is no way to construct a term with type $M (N (M (N a))) \rightarrow M (N a)$, and hence there is no natural *join* function for the composition of M with N . We can regard the types in the rules above as purely formal expressions and, for convenience, will write types like $M (N (M (N X)))$ as strings $MNMNX$. We will also use the notation $rd X$ for the string obtained by removing all adjacent duplicates from X . For example, $rd MMNMNX = MNMNX$. The result that we want follows from the following lemma, proved by a simple structural induction.

LEMMA 1 *If $\vdash X \rightarrow Y$ then $rd X$ is a suffix of $rd Y$.*

Since $MNMNX$ is not a suffix of MNX , it follows that there is no way to define a *join* function with the required type. Similar arguments can be used to show that there is also no way to define natural *prod*, *dorp* or *swap* functions in this framework.