

# Generic Programming

## – An Introduction –

Roland Backhouse<sup>1</sup>, Patrik Jansson<sup>2</sup>, Johan Jeuring<sup>3</sup>, and Lambert Meertens<sup>4</sup>

<sup>1</sup> Department of Mathematics and Computing Science  
Eindhoven University of Technology  
P.O. Box 513  
5600 MB Eindhoven  
The Netherlands  
email: [rolandb@win.tue.nl](mailto:rolandb@win.tue.nl)  
url: <http://www.win.tue.nl/~rolandb/>

<sup>2</sup> Department of Computing Science  
Chalmers University of Technology  
S-412 96 Göteborg  
Sweden  
email: [patrikj@cs.chalmers.se](mailto:patrikj@cs.chalmers.se)  
url: <http://www.cs.chalmers.se/~patrikj/>

<sup>3</sup> Department of Computer Science  
Utrecht University  
P.O. Box 80.089  
3508 TB Utrecht  
The Netherlands  
email: [johanj@cs.uu.nl](mailto:johanj@cs.uu.nl)  
url: <http://www.cs.uu.nl/~johanj/>

<sup>4</sup> CWI & Utrecht University  
P.O. Box 80.089  
3508 TB Utrecht  
The Netherlands  
email: [lambert@cw.nl](mailto:lambert@cw.nl)  
url: <http://www.cwi.nl/~lambert/>

# 1 Introduction

## 1.1 The Abstraction-Specialisation Cycle

The development of science proceeds in a cycle of activities, the so-called abstraction-specialisation cycle. *Abstraction* is the process of seeking patterns or commonalities, which are then classified, often in a formal mathematical framework. In the process of abstraction, we gain greater understanding by eliminating irrelevant detail in order to identify what is essential. The result is a collection of general laws which are then put to use in the second phase of the cycle, the *specialisation* phase. In the specialisation phase the general laws are instantiated to specific cases which, if the abstraction is a good one, leads to novel applications, yet greater understanding, and input for another round of abstraction followed by specialisation.

The abstraction-specialisation cycle is particularly relevant to the development of the science of computing because the modern digital computer is, above all else, a *general-purpose* device that is used for a dazzling range of tasks. Harnessing this versatility is the core task of software design.

Good, commercially viable, software products evolve in a cycle of abstraction and *customisation*. Abstraction, in this context, is the process of identifying a single, general-purpose product out of a number of independently arising requirements. Customisation is the process of optimizing a general-purpose product to meet the special requirements of particular customers. Software manufacturers are involved in a continuous process of abstraction followed by customisation.

## 1.2 Genericity in Programming Languages

The abstraction-specialisation/customisation cycle occurs at all levels of software design. Programming languages play an important role in facilitating its implementation. Indeed, the desire to be able to name and reuse “programming patterns” —capturing them in the form of parametrisable abstractions— has been a driving force in the evolution of high-level programming languages to the extent that the level of “genericity” of a programming language has become a vital criterion for usability.

To determine the level of genericity there are three questions we can ask:

- Which entities can be named in a definition and then referred to by that given name?
- Which entities can be supplied as parameters?
- Which entities can be used “anonymously”, in the form of an expression, as parameters? (For example, in  $y = \sin(2 \times x)$ , the number resulting from  $2 \times x$  is not given a name. In a language allowing numeric parameters, but not anonymously, we would have to write something like  $y = \sin(z)$  **where**  $z = 2 \times x$ .)

An entity for which all three are possible is called a *first-class citizen* of that language.

In one of the first high-level programming languages, FORTRAN (1957), procedures could be named, but not used as parameters. In ALGOL 60 procedures (including functions)

were made almost first-class citizens: they were allowed as parameters, but only by name. In neither language could types be named, nor passed as parameters. In ALGOL 68 procedures were made true first-class citizens, making higher-order functions possible (but not practical, because of an awkward scope restriction<sup>1</sup>). Further, types could be named, but not used as parameters.

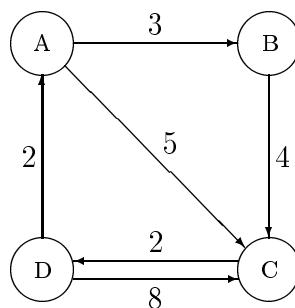
Functional programming languages stand out in the evolution of programming languages because of the high-level of abstraction that is achieved by the combination of higher-order functions and parametric polymorphism. In, for example, Haskell higher-order functions are possible *and practical*. But the level of genericity still has its limitations. Types can be defined *and* used as parameters, but ... types can only be given as parameters in “type expressions”. They cannot be passed to functions. The recent Haskell-like language Cayenne [2] which extends Haskell with dependent types does allow types as arguments and results of functions.

In these lecture notes we introduce another dimension to the level of abstraction in programming languages, namely parameterisation with respect to classes of algebras of variable signature. This first chapter is intended to introduce the key elements of the lectures in broad terms and to motivate what is to follow. We begin by giving a concrete example of a generic algorithm. (The genericity of this algorithm is at a level that can be implemented in conventional functional programming languages, since the parameter is a class of algebras with a fixed signature.) This is followed by a plan of the later chapters.

### 1.3 Path Problems

A good example of parameterising programs by a class of algebras is provided by the problem of finding “extremal” paths in a graph.

Extremal path problems have as input a finite, labelled graph such as the one shown below.



Formally, a *directed graph* consists of a finite set of *nodes*,  $V$ , a finite set of *edges*,  $E$ , and two functions **source** and **target**, each with domain  $E$  and range  $V$ . If **source**  $e$  is the node  $x$  and **target**  $e$  is the node  $y$ , we say that  $e$  is *from*  $x$  *to*  $y$ . (In the figure the nodes

---

<sup>1</sup> Anything that —in implementation terms— would have required what is now known as a “closure”, was forbidden.

are circled and an edge  $e$  is depicted by an arrow that begins at **source**  $e$  and points to **target**  $e$ .) A *path* through the graph from node  $s$  to node  $t$  of *edge length*  $n$  is a finite list of edges  $[e_1, e_2, \dots, e_n]$  such that  $s = \text{source } e_1$  and  $t = \text{target } e_n$  and, for each  $i$ ,  $0 < i < n$ , **target**  $e_i = \text{source } e_{i+1}$ . A graph is *labelled* if it is supplied with a function *label* whose domain is the set of edges,  $E$ .

In an extremal path problem the edge labels are used to weight paths, and the problem is to find the extreme (i.e. best or least, in some sense) weight of paths between given pairs of nodes. We discuss three examples: the reachability problem, the least-cost path problem and the bottleneck problem.

*Reachability* The *reachability* problem is the problem of determining for each pair of nodes  $x$  and  $y$  whether there is a path in the graph from  $x$  to  $y$ . It is solved by a very elegant (and now well-known) algorithm discovered by Roy [42] and Warshall [46]. The algorithm assumes that the nodes are numbered from 1 to  $N$  (say) and that the existence of edges in the graph is given by an  $N \times N$  matrix  $a$  where  $a_{ij}$  is **true** if there is an edge from node numbered  $i$  to the node numbered  $j$ , and **false** otherwise. The matrix is updated by the following code. On termination  $a_{ij}$  is **true** if there is a path from node  $i$  to node  $j$  of edge length at least one; otherwise  $a_{ij}$  is **false**.

```

for each  $k$ ,  $1 \leq k \leq N$ 
do for each pair  $(i, j)$ ,  $1 \leq i, j \leq N$ 
    do  $a_{ij} := a_{ij} \vee (a_{ik} \wedge a_{kj})$ 
    end_for
end_for

```

(The order in which the nodes are numbered, and the order in which pairs of nodes  $(i, j)$  are chosen in the inner loop, is immaterial.)

The reachability problem is an extremal path problem in which all edges have the same label and all paths have the same weight, namely **true**.

*Least-Cost Paths* About the same time as Warshall's discovery of the reachability algorithm, Floyd [17] discovered a very similar algorithm that computes the cost of a least cost path between each pair of nodes in the graph. The algorithm assumes that the matrix  $a$  is a matrix of numbers such that  $a_{ij}$  represents the least *cost* of traversing an edge from node  $i$  to node  $j$ . If there is no edge from  $i$  to  $j$  then  $a_{ij}$  is  $\infty$ . The *cost* of a path is the sum of the costs of the individual edges on the path. Floyd's algorithm for computing the cost of a least cost path from each node to each other node is identical to the Roy-Warshall algorithm above except for the assignment statement which instead is:

$$a_{ij} := a_{ij} \downarrow (a_{ik} + a_{kj})$$

where  $x \downarrow y$  denotes the minimum of  $x$  and  $y$ .

*Bottleneck Problem* A third problem that can be solved with an algorithm of identical shape to the Roy-Warshall algorithm is called the *bottleneck* problem. It is most easily explained as determining the best route to negotiate a high load under a series of low bridges. Suppose an edge in a graph represents a road between two cities and the label is the height of the lowest underpass on the road. The *height* of a path between two nodes is defined to be the minimum of the heights of the individual edges that make up the path. The problem is to determine, for each pair of nodes  $i$  and  $j$ , the maximum of the heights of the paths from node  $i$  to node  $j$  (thus the maximum of the minimum height underpass on a path from  $i$  to  $j$ ).

The bridge height problem is solved by an algorithm identical to the Roy-Warshall algorithm above except for the assignment statement which in this case is:

$$a_{ij} := a_{ij} \uparrow (a_{ik} \downarrow a_{kj})$$

where  $x \downarrow y$  denotes the minimum of  $x$  and  $y$  and  $x \uparrow y$  denotes their maximum. (In the case that there is no edge from  $i$  to  $j$  then the initial value of  $a_{ij}$  is 0.)

*A Generic Path Algorithm* If we abstract from the general shape of these three algorithms we obtain a single algorithm of the form

```

for each  $k$ ,  $1 \leq k \leq N$ 
do for each pair  $(i,j)$ ,  $1 \leq i,j \leq N$ 
  do  $a_{ij} := a_{ij} \oplus (a_{ik} \otimes a_{kj})$ 
  end_for
end_for

```

where  $\oplus$  and  $\otimes$  are binary operators. The initial value of  $a_{ij}$  is the label of the edge from  $i$  to  $j$  if such an edge exists, and is a constant  $\mathbf{0}$  otherwise. (For the purposes of exposition we assume that there is at most one edge from  $i$  to  $j$  for each pair of nodes  $i$  and  $j$ .) The algorithm is thus parameterised by an algebra. In the case of the Roy-Warshall algorithm the carrier of the algebra is the two-element set containing `true` and `false`, the constant  $\mathbf{0}$  is `false`, the operator  $\oplus$  is disjunction and the operator  $\otimes$  is conjunction. In the case of the least-cost path problem the carrier is the set of positive real numbers, the constant  $\mathbf{0}$  is  $\infty$ , the operator  $\oplus$  is the binary minimum operator, and the operator  $\otimes$  is addition. Finally, in the case of the bridge height problem the carrier is also the set of positive real numbers, the operator  $\oplus$  is the binary maximum operator, and the operator  $\otimes$  is minimum.

*Correctness* The above generic algorithm will compute “something” whatever actual parameters we supply for the formal parameters  $\oplus$ ,  $\otimes$  and  $\mathbf{0}$ , the only proviso being that the parameters have compatible types. But, that “something” is only guaranteed to be meaningful if the operators obey certain algebraic properties. The more general *transitive closure* algorithm shown below

```

for each  $k$ ,  $1 \leq k \leq N$ 
do for each pair  $(i,j)$ ,  $1 \leq i,j \leq N$ 
  do  $a_{ij} := a_{ij} \oplus (a_{ik} \otimes (a_{kk})^* \otimes a_{kj})$ 
  end_for
end_for

```

is guaranteed to be correct if the algebra is *regular* [6, 8]<sup>2</sup>. By correctness is meant that if initially

$$a_{ij} = \Sigma \langle e: e \text{ is an edge from } i \text{ to } j: \text{label } e \rangle ,$$

where  $\Sigma$  is the generalisation of the binary operator  $\oplus$  to arbitrary bags, then on termination

$$a_{ij} = \Sigma \langle p: p \text{ is a path of positive edge length from } i \text{ to } j: \text{weight } p \rangle$$

where  $\text{weight } p$  is defined recursively by

$$\text{weight } [] = \mathbf{1}$$

for the empty path  $[]$ , and for paths  $e : p$  (the edge  $e$  followed by path  $p$ )

$$\text{weight } (e : p) = (\text{label } e) \otimes (\text{weight } p) .$$

**Exercise 1.1** Suppose that the edges of a graph are coloured. (So there are blue edges, red edges, etc.) We say that a path has *colour*  $c$  if all the edges on the path have colour  $c$ . Suggest how to use the above algorithm to determine for each pair of nodes  $x$  and  $y$  the set of colours  $c$  such that there is a path of colour  $c$  from  $x$  to  $y$  in the graph.

□

## 1.4 The Plan

The difference between executability and correctness is an important one that shows up time and time again, and it is important to stress it once more. The transitive closure algorithm presented above can be *executed* provided only that instantiations are given for the two constants  $\mathbf{0}$  and  $\mathbf{1}$ , the two binary operators  $\oplus$  and  $\otimes$ , the unary operator  $*$ , the number  $N$  and the matrix  $a$ . An implementation of the algorithm thus requires just the specification of these seven parameters. Moreover, if we bundle the first five parameters together into an algebra, all that is required for the implementation is the *signature* of the algebra: the knowledge that there are two binary operators (with units) and one unary operator. For the *correctness* of the algorithm much, much more is needed. We have to

<sup>2</sup> Without going into complete details, an algebra is regular if it has two constants  $\mathbf{0}$  and  $\mathbf{1}$ , two binary operators  $\oplus$  and  $\otimes$ , and one unary operator  $*$ . The constants  $\mathbf{0}$  and  $\mathbf{1}$  and operators  $\oplus$  and  $\otimes$  should behave like  $0$ ,  $1$ ,  $+$  and  $\times$  in real arithmetic except that  $\times$  is not required to be commutative, and  $+$  is required to be idempotent. The  $*$  operator is a least fixed point operator. The three algebras mentioned above are all regular, after suitably defining the constant  $\mathbf{1}$  and defining  $a^*$  to be  $\mathbf{1}$  for all  $a$ .

supply a specification relative to which correctness is asserted, and establishing correctness demands that we require the algebra to be in a certain *class* of algebras (in this case the class of regular algebras).

As for conventional programs, the specification is absent from a generic program’s implementation. Nevertheless, it is the complete process of *program construction* —from program specification to a systematic derivation of the final implementation— that will dominate the discussion in the coming pages. Our aim is not to show how to derive functional programs but to show how to derive functional programs that are correct by construction. To this end we borrow a number of concepts from category theory, emphasising the calculational properties that these concepts entail.

*Algebras, Functors and Datatypes* The emphasis on calculational properties begins right at the outset in chapter 2 where we introduce the notion of a *functor* and an *initial algebra* and relate these notions to *datatypes*.

An *algebra* (in its simplest form) is a set, called the *carrier* of the algebra, together with a number of operations on that set. A Boolean algebra, for example, has as carrier a set with two elements, commonly named *true* and *false* and binary operations  $\wedge$  (conjunction) and  $\vee$  (disjunction) and unary operation  $\neg$  (negation). The *signature* of the algebra specifies the types of the basic operations in the algebra.

In order to implement a generic algorithm we need to provide the compiler with information on the signature of the operators in the algebra on which the algorithm is parameterised. In order to calculate and reason about generic algorithms we also need a *compact* mechanism for defining signatures. The use of functors provides such a mechanism, compactness being achieved by avoiding naming the operators of the algebra. The use of functors entails much more however than just defining the signature of an algebra. As we shall see, a datatype is a functor and inductively defined datatypes are (the carriers of) *initial* algebras. The concepts of functor, datatype and algebra are thus inextricably intertwined.

*PolyP* Following the discussion of algebras and datatypes, we introduce *PolyP*, an extension of the Haskell programming language in which generic functions can be implemented.

The name of PolyP is derived from “polytypic programming”, polytypic programs being generic programs defined on a particular class of datatypes, the so-called regular datatypes. Writing programs in PolyP means that one can get hands-on experience of generic programming thus reinforcing one’s understanding and, hopefully, leading to further insights.

*A Unification Algorithm* Chapter 4 presents a more substantial example of generic programming — a generic unification algorithm. The basis for the algorithm is a generic construction of a type representing terms with variables, and substitution of terms for variables. The algorithm is implemented using type classes in a style similar to object-oriented programming.

*Relations* The discussion in chapters 2 and 3 is on *functional* programs. In chapter 5 we outline how the concepts introduced in chapter 2 are extended to *relations*, and we show how the extension is used in establishing one element of the correctness of the generic unification algorithm.

There are several reasons for wanting to take the step from functions to relations. The most pressing is that specifications are *relations* between the input and the output, and our concern is with both specifications and implementations. Related to this is that termination properties of programs are typically established by appeal to a well-founded *relation* on the state space. We will not go into termination properties in these lecture notes but the use of well-founded relations will play an integral part in our discussion of one element of the correctness of a generic unification algorithm in chapter 4.

Another reason for wanting to extend the discussion to relations lies in the theoretical basis of generic programming. In chapter 5 we demonstrate how every parametrically polymorphic function satisfies a so-called *logical relation*.

The final reason is *why not?* As we shall see, extending the theory to relations does not significantly increase the complexity whilst the benefits are substantial.

## 1.5 Why Generic Programming?

The form of genericity that we present in the coming pages is novel and has not yet proved its worth. Our goal is to stimulate your interest in exploring it further, and to provide evidence of its potential value.

Generic programming has indeed, potentially, major advantages over “one-shot” programming, since genericity makes it possible to write programs that solve a class of problems once and for all, instead of writing new code over and over again for each different instance. The two advantages that we stress here are the greater potential for reuse, since generic programs are natural candidates for incorporation in library form, and the increased reliability, due to the fact that generic programs are stripped of irrelevant detail which often makes them easier to construct. But what we want to stress most of all is that generic programming is fun. Finding the right generic formulation that captures a class of related problems can be a significant challenge, whose achievement is very satisfying.

*Acknowledgements* The work that is presented here emerged out of the Dutch STOP (Specification and Transformation of Programs) project which ran formally from 1988 to 1992. The project was particularly successful because of the real spirit of cooperation among those participating. Project members (both official and unofficial) included, in alphabetical order, Roland Backhouse, Richard Bird, Henk Doornbos, Maarten Fokkinga, Paul Hoogendijk, Johan Jeuring, Grant Malcolm, Lambert Meertens, Erik Meijer, Oege de Moor, Frans Rietman, Doaitse Swierstra, Jaap van der Woude, Nico Verwer, Ed Voermans. Our thanks go to all who made participation in the project such an enjoyable and stimulating experience.



Development of both practical applications of generic programming and the underlying theory is continuing: see the bibliography for a selection of recent (formally-published and web-published) papers.

## 2 Algebras, Functors and Datatypes

This chapter introduces the concepts fundamental to generic programming. The first section (section 2.1) introduces algebras and homomorphisms between algebras. In this section we see that datatypes (like the natural numbers) are also algebras, but of a special kind. The presentation in section 2.1 is informal. In section 2.4 we make precise in what way datatypes are special: we introduce the all-important notion of an “initial” algebra and the notion of a “catamorphism” (a special sort of homomorphism). The link between the two sections is provided by the intermediate sections on functors. The first of these (section 2.2) provides the formal definition of a functor, motivating it by examples from functional programming. Then section 2.3 introduces further examples of functors forming a class called the “polynomial functors”. Section 2.4 augments the class of polynomial functors with so-called type functors; the resulting class is called the class of “regular functors”, and generic programs defined over the regular functors are called “polytypic” programs. The final section (section 2.5) presents an elementary example of a polytypic program.

### 2.1 Algebras and Homomorphisms

In this section we review the notion of an algebra. The main purpose is to introduce several examples that we can refer to later. The examples central to the discussion are datatypes. At the end of the section we consider how we might formalise the notion of an algebra. We recall a formulation typical of ones in texts on Universal Algebra and remark why this is inadequate for our purposes. We then present the definition of an algebra in category theory based on the notion of a “functor” and outline how the latter expresses the content of the traditional definitions much more succinctly and in a much more structured way.

*Algebras* An algebra is a set, together with a number of operations (functions) that return values in that set. The set is called the *carrier* of the algebra. Here are some concrete examples of algebras:

$$\begin{array}{l} (\mathbb{N}, 0, (+)), \text{ with } 0 :: 1 \rightarrow \mathbb{N}, \quad (+) :: \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N} \\ (\mathbb{N}, 0, (\uparrow)), \text{ with } 0 :: 1 \rightarrow \mathbb{N}, \quad (\uparrow) :: \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N} \\ (\mathbb{R}, 1, (\times)), \text{ with } 1 :: 1 \rightarrow \mathbb{R}, \quad (\times) :: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R} \\ (\mathbb{B}, \text{true}, (\equiv)), \text{ with } \text{true} :: 1 \rightarrow \mathbb{B}, \quad (\equiv) :: \mathbb{B} \times \mathbb{B} \rightarrow \mathbb{B} \\ (\mathbb{B}, \text{false}, (\vee)), \text{ with } \text{false} :: 1 \rightarrow \mathbb{B}, \quad (\vee) :: \mathbb{B} \times \mathbb{B} \rightarrow \mathbb{B} \\ (\mathbb{B}, \text{true}, (\wedge)), \text{ with } \text{true} :: 1 \rightarrow \mathbb{B}, \quad (\wedge) :: \mathbb{B} \times \mathbb{B} \rightarrow \mathbb{B} \\ (A^*, \varepsilon, (+)), \text{ with } \varepsilon :: 1 \rightarrow A^*, \quad (+) :: A^* \times A^* \rightarrow A^* \end{array}$$

In the last line  $A^*$  stands for the words over some alphabet  $A$ , with “+” denoting word concatenation, and “ $\varepsilon$ ” the empty word. This is, of course, basically the same algebra as  $(List\ A, [], (+))$ , the (finite) lists of  $A$ -elements with list concatenation. Note that in the typing of the operations we use the notation “source-type  $\rightarrow$  target-type”. In an algebra all operations have the same target type<sup>3</sup>: its carrier. Note further that we use the “uncurried”

---

<sup>3</sup> We freely identify types and sets whenever convenient.

view in which a binary operation takes a pair (2-tuple) of arguments and so has some type like  $A \times B \rightarrow C$ . To make fixed elements, like  $0 \in \mathbb{N}$ , fit in, they are treated here as nullary operations: operations with a 0-tuple of arguments. This is indicated by the source type 1, which in Haskell would be denoted as “()”. Sometimes we will instantiate a generic program to a specific Haskell program, and in doing so we will switch back to the curried view for binary operations, having some type  $A \rightarrow (B \rightarrow C)$ , and to the view of nullary operations as plain elements, having type  $A$  rather than  $1 \rightarrow A$ . Conversely, going from a Haskell program to an algebraic view, we will uncurry  $n$ -ary functions,  $n \geq 2$ , and treat constants as nullary functions.

The concrete algebras above were chosen in such a way that they all have the same number of operations with the same typing pattern. They can be unified generically into the following abstract algebra:

$$(A, e, \oplus), \text{ with } e :: 1 \rightarrow A, \oplus :: A \times A \rightarrow A$$

So they all belong to the same *class* of algebras. An example of another class of algebras is:

$$\begin{aligned} (\mathbb{N}, (+), (+1)), & \text{ with } (+) :: \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}, (+1) :: \mathbb{N} \rightarrow \mathbb{N} \\ (\mathbb{R}, (\times), (\times 2)), & \text{ with } (\times) :: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}, (\times 2) :: \mathbb{R} \rightarrow \mathbb{R} \\ (A, \oplus, f), & \text{ with } \oplus :: A \times A \rightarrow A, f :: A \rightarrow A \end{aligned}$$

Here, the first two are concrete, while the last is the generic algebra.

By just looking at an algebra, it is not possible (in general) to tell what class of algebras it belongs to: a given algebra can belong to several different classes. So the class information has to be supplied additionally. Take for example the following class:

$$\begin{aligned} (\mathbb{N}, 0, (+)), & \text{ with } 0 :: 1 \rightarrow \mathbb{N}, (+) :: \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N} \\ (\mathbb{N}, 0, (\uparrow)), & \text{ with } 0 :: 1 \rightarrow \mathbb{N}, (\uparrow) :: \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N} \\ (List\ \mathbb{N}, [], (:\)), & \text{ with } [] :: 1 \rightarrow List\ \mathbb{N}, (:) :: \mathbb{N} \times List\ \mathbb{N} \rightarrow List\ \mathbb{N} \\ (A, e, \oplus), & \text{ with } e :: 1 \rightarrow A, \oplus :: \mathbb{N} \times A \rightarrow A \end{aligned}$$

The first two concrete algebras also occur in the first class treated above, but the generic algebra reveals that this is a different class.

To give a concluding example of an algebra class:

$$\begin{aligned} (\mathbb{N}, 0, (+1)), & \text{ with } 0 :: 1 \rightarrow \mathbb{N}, (+1) :: \mathbb{N} \rightarrow \mathbb{N} \\ (\mathbb{R}, 1, (\times 2)), & \text{ with } 1 :: 1 \rightarrow \mathbb{R}, (\times 2) :: \mathbb{R} \rightarrow \mathbb{R} \\ (\mathbb{B}, true, (\neg)), & \text{ with } true :: 1 \rightarrow \mathbb{B}, (\neg) :: \mathbb{B} \rightarrow \mathbb{B} \\ (\mathbb{B}, false, (\neg)), & \text{ with } false :: 1 \rightarrow \mathbb{B}, (\neg) :: \mathbb{B} \rightarrow \mathbb{B} \\ (A, e, f), & \text{ with } e :: 1 \rightarrow A, f :: A \rightarrow A \end{aligned}$$

A recursively defined datatype determines, in a natural way, an algebra. A simple example is the datatype *Nat* defined by<sup>4</sup>:

```
data Nat = zero | succ Nat
```

<sup>4</sup> We use Haskell syntax for defining datatypes, except that we write constructors using a sans serif font where Haskell would capitalize the first letter. The Haskell definition of *Nat* would be `data Nat = Zero | Succ Nat`.

The corresponding algebra is:

$$(Nat, zero, succ), \text{ with } zero :: 1 \rightarrow Nat, \text{ succ} :: Nat \rightarrow Nat$$

This belongs to the last class mentioned; in fact, if we ignore the possibility of infinite data structures —made possible by lazy evaluation— this is essentially the same algebra as  $(\mathbb{N}, 0, (+1))$ . Another example is:

$$\mathbf{data} \text{ Natlist} = \mathbf{nil} \mid \mathbf{cons} \ \mathbb{N} \ \text{Natlist}$$

The corresponding algebra is:

$$(\text{Natlist}, \mathbf{nil}, \mathbf{cons}), \text{ with } \mathbf{nil} :: 1 \rightarrow \text{Natlist}, \text{ cons} :: \mathbb{N} \times \text{Natlist} \rightarrow \text{Natlist}$$

This is basically the same as  $(List \ \mathbb{N}, [], (:))$ . Both of these examples illustrate the general phenomenon that a recursively defined datatype determines an algebra in which the carrier of the algebra is the datatype itself, and the constructors of the datatype are the operations of the algebra.

*Homomorphisms* A homomorphism between two algebras, which must be *from the same class*, is a function between their carrier sets that “respects the structure” of the class. For example, the function  $\mathbf{exp} :: \mathbb{N} \rightarrow \mathbb{R}$  is a homomorphism with as *source algebra*  $(\mathbb{N}, 0, (+))$  and as *target algebra*  $(\mathbb{R}, 1, (\times))$ . In this case, respecting the structure of this algebra class means that it satisfies the following two properties:

$$\begin{aligned} \mathbf{exp} \ 0 &= 1 \\ \mathbf{exp}(x + y) &= (\mathbf{exp} \ x) \times (\mathbf{exp} \ y) \end{aligned}$$

Another example in the same class is  $\mathbf{length} :: (A^*, \varepsilon, (+)) \rightarrow (\mathbb{N}, 0, (+))$ . (This notation is shorthand for the statement that the function  $\mathbf{length} :: A^* \rightarrow \mathbb{N}$  is a homomorphism from source algebra  $(A^*, \varepsilon, (+))$  to target algebra  $(\mathbb{N}, 0, (+))$ ). In this case, respecting the structure means:

$$\begin{aligned} \mathbf{length} \ \varepsilon &= 0 \\ \mathbf{length}(x + y) &= (\mathbf{length} \ x) + (\mathbf{length} \ y) \end{aligned}$$

In general (for this class of algebras),  $h :: (A, u, \otimes) \rightarrow (B, e, \oplus)$  means:

$$\begin{aligned} h &:: A \rightarrow B \\ h \ u &= e \\ h(x \otimes y) &= (h \ x) \oplus (h \ y) \end{aligned}$$

So to apply  $h$  to a value in  $A$  that resulted from a  $u$ -operation (and there is only one such value), we may equally apply  $h$  to the operands (of which there are none) and apply  $e$  to the resulting 0-tuple. Similarly, to apply  $h$  to a value in  $A$  that resulted from a  $\otimes$ -operation, we may equally well apply  $h$  to the operands (which gives two  $B$ -values) and combine these with the operation  $\oplus$ . Here are some more examples of homomorphisms in this class:

$$\begin{aligned} (\downarrow 1) &:: (\mathbb{N}, 0, (+)) \rightarrow (\mathbb{N}, 0, (\uparrow)) \\ \mathbf{even} &:: (\mathbb{N}, 0, (+)) \rightarrow (\mathbb{B}, \mathbf{true}, (\equiv)) \\ (> 0) &:: (\mathbb{N}, 0, (\uparrow)) \rightarrow (\mathbb{B}, \mathbf{false}, (\vee)) \\ (\neg) &:: (\mathbb{B}, \mathbf{false}, (\vee)) \rightarrow (\mathbb{B}, \mathbf{true}, (\wedge)) \\ (\neg) &:: (\mathbb{B}, \mathbf{true}, (\wedge)) \rightarrow (\mathbb{B}, \mathbf{false}, (\vee)) \end{aligned}$$

If we have two homomorphisms  $h :: (A, e, \oplus) \rightarrow (B, u, \otimes)$  and  $k :: (B, u, \otimes) \rightarrow (C, z, \odot)$ , in which the target algebra of the first is the source algebra of the second, then their composition is also a homomorphism  $k \cdot h :: (A, e, \oplus) \rightarrow (C, z, \odot)$ . For example,

$$\begin{aligned} (> 0) \cdot (\downarrow 1) &:: (\mathbb{N}, 0, (+)) \rightarrow (\mathbb{B}, \text{false}, (\vee)) \\ (\neg) \cdot (\neg) &:: (\mathbb{B}, \text{false}, (\vee)) \rightarrow (\mathbb{B}, \text{false}, (\vee)) \end{aligned}$$

Now  $(> 0) \cdot (\downarrow 1) = (> 0)$  on  $\mathbb{N}$ , and  $(\neg) \cdot (\neg) = \text{id}_{\mathbb{B}}$  (the identity function on  $\mathbb{B}$ ), so we have

$$\begin{aligned} (> 0) &:: (\mathbb{N}, 0, (+)) \rightarrow (\mathbb{B}, \text{false}, (\vee)) \\ \text{id} &:: (\mathbb{B}, \text{false}, (\vee)) \rightarrow (\mathbb{B}, \text{false}, (\vee)) \end{aligned}$$

The identity function  $\text{id}_A$  is of course a homomorphism between *any* algebra with carrier  $A$  and itself.

For the class of algebras whose generic algebra is

$$(A, e, \oplus), \quad \text{with } e :: 1 \rightarrow A, \quad \oplus :: \mathbb{N} \times A \rightarrow A$$

we have that  $h :: (A, e, \oplus) \rightarrow (B, u, \otimes)$  means:

$$\begin{aligned} h &:: A \rightarrow B \\ h e &= u \\ h(x \oplus y) &= x \otimes (h y) \end{aligned}$$

So why is  $h$  for this class not applied to the occurrence of  $x$  in the righthand side of the second equality? The answer is that that would not make sense, since  $h$  has source type  $A$ , but  $x$  is of type  $\mathbb{N}$ . (Later, after we have introduced functors, we shall see how to define the notion of homomorphism generically, independent of the specific algebra class.) We have:

$$\begin{aligned} \text{sum} &:: (\text{List } \mathbb{N}, [], (:)) \rightarrow (\mathbb{N}, 0, (+)) \\ \text{foldr } \oplus e &:: (\text{List } \mathbb{N}, [], (:)) \rightarrow (A, e, \oplus) \end{aligned}$$

In fact,  $\text{sum} = \text{foldr } (+) 0$ .

*Uniqueness* We have given several examples of algebra classes and their homomorphisms. The first class had generic algebra

$$(A, e, \oplus) \quad \text{with } e :: 1 \rightarrow A, \quad \oplus :: A \times A \rightarrow A .$$

Note that the fact that a function is a homomorphism of algebras in this class does not uniquely define the function. For example, we observed above that `length` is a homomorphism with source  $(A^*, \varepsilon, (++)$ ) and target  $(\mathbb{N}, 0, (+))$ . But the function that is constantly 0 for all lists is also a homomorphism with exactly the same source and target algebras. Indeed, in the case of all the examples we gave of homomorphisms between algebras in this class the constant function returning the value  $e$  of the target algebra has the same homomorphism type as the given function.

Contrast this with the third class of algebras. The generic algebra has the form

$$(A, e, \oplus) \quad \text{with } e :: 1 \rightarrow A, \quad \oplus :: \mathbb{N} \times A \rightarrow A$$

Again, the fact that a function is a homomorphism of algebras in this class does not uniquely define the function. But there is something rather special about the algebra  $(List\ \mathbb{N}, [], (:))$  in this class of algebras. Specifically,  $foldr\ \oplus\ e$  is the *unique* homomorphism with source algebra  $(List\ \mathbb{N}, [], (:))$  and target algebra  $(A, e, \oplus)$ . For example, `sum` is the unique homomorphism with source  $(List\ \mathbb{N}, [], (:))$  and target  $(\mathbb{N}, 0, (+))$ . That is, function  $h$  satisfies the equations

$$\begin{aligned} h &:: List\ \mathbb{N} \rightarrow \mathbb{N} \\ h\ [] &= 0 \\ h(x:xs) &= x + (h\ xs) \end{aligned}$$

if and only if  $h = \text{sum}$ .

This uniqueness is an important property that will be a focus of later discussion.

*Isomorphisms* Above, we said several times that two algebras were “basically” or “essentially” the same. We want to make this notion precise. The technical term for this is that these algebras are *isomorphic*. In set theory, two sets  $A$  and  $B$  are called isomorphic whenever there exists a bijection between  $A$  and  $B$ . Equivalently,  $A$  and  $B$  are isomorphic whenever there exist functions  $f :: A \rightarrow B$  and  $g :: B \rightarrow A$  that cancel each other, that is:

$$\begin{aligned} f \circ g &= \text{id}_B \\ g \circ f &= \text{id}_A \end{aligned}$$

The generalisation for algebras is now that we require these functions to be *homomorphisms* between the algebras involved. A homomorphism that has a cancelling homomorphism is called an *isomorphism*. From the examples above we see that the algebras  $(\mathbb{B}, \text{true}, (\wedge))$  and  $(\mathbb{B}, \text{false}, (\vee))$  are isomorphic.

*Algebras with laws* Although we will hardly use this, no account of the notion of algebra is complete without mentioning the following. A class of algebras can be further determined by a set of laws. In a “lawful” class of algebras, all algebras satisfy the same set of (possibly conditional) equational laws. Monoids form the best-known example of a lawful algebra class. The generic monoid is  $(A, e, \oplus)$ , with  $e :: 1 \rightarrow A$ ,  $\oplus :: A \times A \rightarrow A$ , and the monoid laws are the following two:

$$\begin{aligned} \oplus \text{ is associative: } & (x \oplus y) \oplus z = x \oplus (y \oplus z) \\ e \text{ is neutral for } \oplus: & e \oplus x = x = x \oplus e \end{aligned}$$

If an operation  $\oplus$  has a neutral element, it is unique, and we denote it as  $\nu_{\oplus}$ . For example,  $\nu_{+} = 0$  and  $\nu_{\times} = 1$ . The examples of concrete algebras from the first class treated in this chapter are actually all monoids. For lawful algebras the definition of homomorphism is the same as before.

*Graphs* The notion of homomorphism is more general than that of a “structure-respecting” function between algebras. Homomorphisms can generally be defined for anything having

structure. As an example, we consider homomorphisms between directed graphs. Recall that a directed graph is a structure

$$(V, E, \text{source}, \text{target}), \quad \text{with } \text{source} :: E \rightarrow V, \text{target} :: E \rightarrow V$$

in which the elements of  $V$  are called “vertices” or “nodes”, and the elements of  $E$  are called “edges” or “arcs”. If edge  $e$  is an edge from node  $m$  to node  $n$ , we have:  $\text{source } e = m$  and  $\text{target } e = n$ . Directed graphs are just like an algebra class, except that we have two “carrier sets”:  $V$  and  $E$ . (There is a term for algebras with more carrier sets: *heterogeneous* or *multi-sorted* algebras.) A homomorphism from graph  $(V_0, E_0, \text{source}_0, \text{target}_0)$  to graph  $(V_1, E_1, \text{source}_1, \text{target}_1)$  is a *pair* of functions, one with the typing  $V_0 \rightarrow V_1$  and one with the typing  $E_0 \rightarrow E_1$ , and if we overload the identifier  $h$  to denote both functions, they satisfy:

$$\begin{aligned} h(\text{source } a) &= \text{source}(h a) \\ h(\text{target } a) &= \text{target}(h a) \end{aligned}$$

As before for algebras, two graphs are isomorphic whenever there are cancelling homomorphisms between them. Informally, this means that one graph can be obtained from the other by systematic renaming. In standard Graph Theory, for unlabelled graphs like the ones we are considering here, two isomorphic graphs are usually considered *identical*. Still, there can be non-trivial *automorphisms*, that is, isomorphisms between a graph and itself that are not the identity isomorphism.

*Summarising and looking ahead* In this section we have introduced the notion of a class of algebras and homomorphisms between algebras in the same class. We have observed that datatype definitions in a functional programming language define an algebra, the carrier of the algebra being the datatype itself and the operations being the constructors of the datatype. We have also made the important observation that in some cases a function is uniquely characterised by its homomorphism type (the fact that it is a homomorphism combined with knowledge about its source and target algebras).

In the remaining sections of this chapter our goal is to formalise all these ideas in a way that facilitates the calculational construction of programs. Let us give an outline of what is in store.

The notion of an algebra is formalised in many textbooks on Universal Algebra. Here is an example of such a definition. This is *not* the definition we intend to use so you don’t need to understand it in detail.

*$\Sigma$ -algebra* A  $\Sigma$ -algebra with respect to a signature with operators  $\Sigma = (S, F)$  is a pair  $(V, F)$  such that

- $V$  is an  $S$ -sorted set, and

–  $F = \{\gamma: \gamma \in \cup \Gamma: f_\gamma\}$  is a set of functions such that

$$\gamma \in \Gamma_{\langle\langle s_0, \dots, s_{n-1} \rangle, r \rangle} \Rightarrow f_\gamma \in V_{s_0} \times \dots \times V_{s_{n-1}} \rightarrow V_r$$

$$\gamma \in \Gamma_{\langle s, r \rangle} \Rightarrow f_\gamma \in V_s \rightarrow V_r$$

$V$  is called the *carrier set* of the  $\Sigma$ -algebra and set  $F$  is its *operator set*.

Contrast this with the definition we are going to explain in the coming sections.

*F-algebra* Suppose  $F$  is a functor. Then an  $F$ -algebra is a pair  $(A, \alpha)$  such that  $\alpha \in FA \rightarrow A$ .

Neither definition is complete since in the first definition the notion of a signature has not been defined, and in the second the notion of a functor hasn't been defined. In the first definition, however, it's possible to guess what the definition of a signature is and, after struggling some time with the subscripts of subscripts, it is possible to conclude that the definition corresponds to the "intuitive" notion of an algebra. The disadvantage is that the definition is grossly unwieldy. If the definitions of one's basic concepts are as complicated as this then one should give up altogether any hope that one can calculate with them.

The second definition is very compact and, as we shall see, gives an excellent basis for program construction. Its disadvantage, however, is that it is impossible to guess what the definition of a functor might be, and it is difficult to see how it corresponds to the familiar notion of an algebra. How is it possible to express the idea that an algebra consists of a set of operations? On the face of it, it would appear that an  $F$ -algebra has just one operation  $\alpha$ . Also, how does one express the fact that the operations in an algebra have various arities?

The answer to these questions is hidden in the definition of a "functor". And, of course, if its definition is long and complicated then all the advantages of the compactness of the definition of an algebra are lost. We shall see, however, that the definition of a functor is also very compact. We shall also see that functors can be constructed from primitive functors in a systematic way. The "disjoint sum" of two functors enables one to express the idea that an algebra has a set of operations; the "cartesian product" of functors allows one to express the arity of the various operations; "constant functors" enable the expression of the existence of designated constants in an algebra. An additional major bonus is that the categorical notion of an "initial algebra" leads to a very compact and workable definition of inductively defined datatypes in a programming language. The remaining sections of this chapter thus provide a veritable arsenal of fundamental concepts whose mastery is tremendously worthwhile.

**Exercise 2.1** Check the claim that  $\text{even} :: (\mathbb{N}, 0, (+)) \rightarrow (\mathbb{B}, \text{true}, (\equiv))$  is a homomorphism.

**Exercise 2.2** Give the composition of the following two homomorphisms:

$$\begin{aligned} (\neg) &:: (\mathbb{B}, \text{false}, (\vee)) \rightarrow (\mathbb{B}, \text{true}, (\wedge)) \\ (> 0) &:: (\mathbb{N}, 0, (+)) \rightarrow (\mathbb{B}, \text{false}, (\vee)) \end{aligned}$$



**Exercise 2.3** An automorphism is an isomorphism with the same source and target algebra. Show that the only automorphism on the algebra  $(\mathbb{B}, \text{true}, (\equiv))$  is the trivial automorphism  $\text{id}$ .

**Exercise 2.4** Give an example of a non-trivial automorphism on the algebra  $(\mathbb{R}, 0, (\times))$ .

## 2.2 Functors

To a first approximation, datatypes are just sets. A second approximation, which we have just seen, is that a datatype is the carrier of an algebra. In this section we identify *parameterised* datatypes with the categorical notion of functor, giving us a third approximation to what it is to be a datatype. It is in this section that we take the first steps towards a generic theory of datatypes.

*Examples* The best way to introduce the notion of a functor is by abstraction from a number of examples. Here are a few datatype definitions:

```

data List a = nil | cons a (List a)

data Maybe a = none | one a

data Bin a = tip a | join (Bin a) (Bin a)

data Rose a = fork a (List(Rose a))

```

Each of these types can be viewed as a structured repository of information, the type of information being specified by the parameter  $a$  in the definition. Each of these types has its own `map` combinator. “Mapping” a function over an instance of one of these datatypes means applying the function to all the values stored in the structure without changing the structure itself. The typings of the individual `map` combinators are thus as follows.

```

mapList    :: (a → b) → (List a → List b)
mapMaybe  :: (a → b) → (Maybe a → Maybe b)
mapBin     :: (a → b) → (Bin a → Bin b)
mapRose    :: (a → b) → (Rose a → Rose b)

```

A datatype that has more than one type parameter also has a `map` combinator, but with more arguments. For instance, defining the type of trees with leaves of type  $a$  and interior nodes of type  $b$  by

```

data Tree a b = leaf a | node (Tree a b) b (Tree a b)

```

the corresponding `map` combinator has type

```

mapTree    :: (a → c) → (b → d) → (Tree a b → Tree c d)

```

Given a tree of type  $\text{Tree } a \ b$ , the combinator applies a function of type  $a \rightarrow c$  to all the leaves of the tree, and a function of type  $b \rightarrow d$  to all the nodes, thus creating a tree of type  $\text{Tree } c \ d$ .

In general, the `map` combinator for an  $n$ -ary datatype maps  $n$  functions over the values stored in the datatype. (This also holds for the case that  $n$  is zero. Datatypes having *no* type parameter also have a `map` combinator, but with *no* functional arguments! The `map` in this case is the identity function on the elements of the datatype.)

*Functors Defined* The idea that parameterised datatypes are structured repositories of information over which arbitrary functions can be mapped is captured by the concept of a *functor*. We first explain the concept informally for unary functors. Consider the world of typed functions. Functors are the structure-respecting functions for that world. So what is the structure involved? First, that world can be viewed as a directed graph, in which the nodes are types and the arcs are functions. So, as for graphs, we require that a functor is a *pair* of mappings, one acting on types and one acting on functions, and if we overload the identifier  $F$  to denote both functions, they satisfy the typing rule:

$$\frac{f :: a \rightarrow b}{Ff :: Fa \rightarrow Fb}$$

Further, functions can be composed with the operation “ $\bullet$ ”, which is associative and has neutral element the identity function, `id`, so this world forms a monoid algebra. Functors also respect the monoid structure:

$$\begin{aligned} F(f \bullet g) &= (Ff) \bullet (Fg) \\ F \text{ id}_a &= \text{id}_{Fa} \end{aligned}$$

The first of these laws says that there is no difference between mapping the composition of two functions over an  $F$  structure in one go and mapping the functions over the structure one by one. The second law says that mapping the identity function over an  $F$  structure of  $a$ 's has no effect on the structure.

To be completely precise, the world of functions is not quite a monoid, since the algebra is *partial*: the meaning of  $f \bullet g$  is only defined when this composition is well-typed, that is, when the source type of  $f$  is the target type of  $g$ . The first equality above should therefore only be applied to cases for which  $f \bullet g$  is defined, and from now on we assume this as a tacit condition on such equations. It follows from the typing rule that then also the composition  $(Ff) \bullet (Fg)$  is well-typed, so that is not needed as a condition.

Now, in general, an  $n$ -ary *functor*  $F$  is a pair of mappings that maps an  $n$ -tuple of types  $a_0, \dots, a_{n-1}$  to a type  $F a_0 \cdots a_{n-1}$  and an  $n$ -tuple of functions  $f_0, \dots, f_{n-1}$  to a function  $F f_0 \cdots f_{n-1}$  in such a way that typing, composition and identity are respected:

$$\frac{f_i \quad :: \quad a_i \quad \rightarrow \quad b_i \quad \text{for } i = 0, \dots, n-1}{Ff_0 \cdots f_{n-1} :: Fa_0 \cdots a_{n-1} \rightarrow Fb_0 \cdots b_{n-1}}$$

$$\begin{aligned} F(f_0 \bullet g_0) \cdots (f_{n-1} \bullet g_{n-1}) &= (Ff_0 \cdots f_{n-1}) \bullet (Fg_0 \cdots g_{n-1}) \\ F \text{ id} \cdots \text{id} &= \text{id} \end{aligned}$$

*Examples Revisited* As anticipated in the introduction to this section, the pairs of mappings  $F$  (on types) and  $\text{map}_F$  (on functions) for  $F = \text{List}, \text{Maybe}$ , etcetera, are all unary functors since they satisfy the typing rule

$$\frac{f :: a \rightarrow b}{\text{map}_F f :: F a \rightarrow F b}$$

and the functional equalities

$$\begin{aligned} \text{map}_F (f \bullet g) &= (\text{map}_F f) \bullet (\text{map}_F g) \\ \text{map}_F \text{id} &= \text{id} \end{aligned}$$

An example of a binary functor is the pair of mappings  $\text{Tree}$  and  $\text{map}_{\text{Tree}}$  since the pair satisfies the typing rule

$$\frac{\begin{array}{ccc} f :: & a & \rightarrow & c \\ g :: & b & \rightarrow & d \end{array}}{\text{map}_{\text{Tree}} f g :: \text{Tree } a b \rightarrow \text{Tree } c d}$$

and the functional equalities

$$\begin{aligned} \text{map}_{\text{Tree}} (f \bullet g) (h \bullet k) &= (\text{map}_{\text{Tree}} f h) \bullet (\text{map}_{\text{Tree}} g k) \\ \text{map}_{\text{Tree}} \text{id} \text{id} &= \text{id} \end{aligned}$$

*Notational convention* Conventionally, the *same* notation is used for the type mapping and the function mapping of a functor, and we follow that convention here. Moreover, when applicable, we use the name of the type mapping. So, from here on, for function  $f$ , we write  $\text{List } f$  rather than  $\text{map}_{\text{List}} f$ .

**Exercise 2.5** Consider the following datatype declarations. Each defines a mapping from types to types. For example,  $\text{Error}$  maps the type  $a$  to the type  $\text{Error } a$ . Extend the definition of each so that it becomes a functor.

```
data Error a = error String | ok a

data Drawing a = above (Drawing a) (Drawing a)
                | beside (Drawing a) (Drawing a)
                | atom a
```

□

## 2.3 Polynomial Functors

Now that we have defined the notion of a functor and have seen some non-trivial examples it is time to consider more basic examples. Vital to the usefulness of the notion is that non-trivial functors can be constructed by composing more basic functors. In this section we consider the *polynomial* functors. As the name suggests, these are the functors that

can be obtained by “addition” and “multiplication” possibly combined with the use of a number of “constants”.

The technical terms for addition and multiplication are “disjoint sum” and “cartesian product”. The use of disjoint sum enables one to capture in a *single* functor the fact that an algebra has a *set* of operations. The use of cartesian product enables one to express the fact that an operator in an algebra has an arity greater than one. We also introduce constant functors and the identity functor; these are used to express the designated constants (functions of arity zero) and unary functions in an algebra, respectively. For technical reasons, we also introduce a couple of auxiliary functors in order to complete the class of polynomial functors. We begin with the simpler cases.

*The identity functor* The simplest example of a functor is the *identity functor* which is the trivial combination of two identity functions, the function that maps every type to itself and the function that maps every function to itself. Although trivial, this example is important and shouldn’t be forgotten. We denote the identity functor by `Id`.

*Constant functors* For the constant mapping that maps any  $n$ -tuple of arguments to the same result  $x$  we use the notation  $x^k$ . As is easily verified, the pair of mappings  $a^k$  and  $\text{id}_a^k$ , where  $a$  is some type, is also a functor. It is  $n$ -ary for all  $n$ .

Following the naming convention introduced above, we write  $a^k$  to denote both the mapping on types and the mapping on functions. That is, we write  $a^k$  where strictly we should write  $\text{id}_a^k$ . So, for functions  $f_0 \dots f_{n-1}$ ,  $a^k f_0 \dots f_{n-1} = \text{id}_a$ .

A constant functor that we will use frequently is the constant functor associated with the *unit type*, `1`. The unit type is the type that is denoted `()` in Haskell. It is a type having exactly one element (which element is also denoted `()` in Haskell). This functor will be denoted by `1` rather than  $1^k$ .

*Extraction* Each extraction combinator

$$\text{Ex}_i^n z_0 \dots z_{n-1} = z_i, \quad \text{for } i = 0, \dots, n-1$$

is an  $n$ -ary functor. The extractions that we have particular use for are the identity functor `Id`, which is the same as  $\text{Ex}_0^1$ , and the binary functors  $\text{Ex}_0^2$  and  $\text{Ex}_1^2$ , for which we use the more convenient notations `Par` and `Rec`. (The reason for this choice of identifiers will become evident in chapter 3. When defining recursive datatypes like *List*, we identify a binary “pattern functor”. The first parameter of the pattern functor is the parameter of the recursive datatype —and is thus called the `Par` parameter— and the second parameter is used as the argument for recursion —and is thus called the `Rec` parameter.)

*The sum functor* The binary sum functor `+` gives the “disjoint union” of two types. We write it as an infix operator. It is defined by:

```
data a + b = Inl a | Inr b
```

$$\begin{aligned}
f + g &= h \text{ where} \\
&h(\text{inl } u) = \text{inl}(f u) \\
&h(\text{inr } v) = \text{inr}(g v)
\end{aligned}$$

$$\begin{aligned}
f \nabla g &= h \text{ where} \\
&h(\text{inl } u) = f u \\
&h(\text{inr } v) = g v
\end{aligned}$$

The datatype definition introduces both the type  $a+b$ , called the *disjoint sum* of  $a$  and  $b$ , and the two constructor functions  $\text{inl} :: a \rightarrow a + b$  and  $\text{inr} :: b \rightarrow a + b$ . The name “disjoint sum” is used because  $a+b$  is like the set union of  $a$  and  $b$  except that each element of the sets  $a$  and  $b$  is, in effect, tagged with either the label  $\text{inl}$ , to indicate that it originated in set  $a$ , or  $\text{inr}$ , to indicate that it originated in set  $b$ . In this way  $a+a$  is different from  $a$  since it effectively contains two copies of every element in  $a$ , one with label  $\text{inl}$  and one with label  $\text{inr}$ . In particular  $1+1$  has two elements. The constructors  $\text{inl}$  and  $\text{inr}$  are called *injections* and are said to *inject* elements of  $a$  and  $b$  into the respective components of  $a+b$ .

In order to extend the sum mapping on types to a functor we have to define the sum of two functions. This is done in the definition of  $f+g$  above. Its definition is obtained by type considerations — if  $+$  is to be a functor, we require that if  $f :: a \rightarrow b$  and  $g :: c \rightarrow d$  then  $f+g :: a+c \rightarrow b+d$ . It is easily checked that the above definition of  $f+g$  meets this requirement; indeed, there is no other way to do so.

In addition to defining  $f+g$  we have defined another way of combining  $f$  and  $g$ , namely  $f \nabla g$ , which we pronounce  $f$  “junc”  $g$ . (“Junc” is short for “junction”.) As we’ll see shortly,  $f \nabla g$  is more basic than  $f+g$ . The meaning of  $f \nabla g$  is only defined when  $f$  and  $g$  have the same target type; its source type is a disjoint sum of two types. Operationally, it inspects the label on its argument to see whether the argument originates from the left or right component of the disjoint sum. Depending on which component it is, either the function  $f$  or the function  $g$  is applied to the argument after first stripping off the label. In other words,  $f \nabla g$  acts like a case statement, applying  $f$  or  $g$  depending on which component of the disjoint sum the argument comes from.

The *typing rule* for  $\nabla$  is a good way of memorising its functionality:

$$\frac{f :: a \rightarrow c \quad g :: b \rightarrow c}{f \nabla g :: a+b \rightarrow c}$$

(Haskell’s prelude contains a definition of disjoint sum:

```
data Either a b = Left a | Right b
```

with *either* playing the role of  $\nabla$ .)

Now that we have defined  $+$  on types and on functions in such a way as to fulfill the typing requirements on a (binary) functor it remains to verify that it respects identities

and composition. We do this now. In doing so, we establish a number of calculational properties that will prove to be very useful for other purposes.

Note first that the definitions of  $+$  (on functions) and of  $\nabla$  can be rewritten in point-free style as the following *characterisations*:

$$\boxed{\begin{array}{l} h = f + g \quad \equiv \quad h \bullet \text{inl} = \text{inl} \bullet f \wedge h \bullet \text{inr} = \text{inr} \bullet g \\ h = f \nabla g \quad \equiv \quad h \bullet \text{inl} = f \wedge h \bullet \text{inr} = g \end{array}}$$

This style is convenient for reasoning. For example, we can prove the *identity rule*:

$$\boxed{\text{inl} \nabla \text{inr} = \text{id}}$$

by calculating as follows:

$$\begin{aligned} \text{id} &= \alpha \nabla \beta \\ &\equiv \quad \{ \text{characterisation of } \nabla \} \\ \text{id} \bullet \text{inl} &= \alpha \wedge \text{id} \bullet \text{inr} = \beta \\ &\equiv \quad \{ \text{id is the identity of composition} \} \\ \text{inl} &= \alpha \wedge \text{inr} = \beta \quad . \end{aligned}$$

This last calculation is a simple illustration of the way we often *derive* programs. In this case the goal is to express  $\text{id}$  in terms of  $\nabla$ . We therefore introduce the unknowns  $\alpha$  and  $\beta$ , and calculate expressions for  $\alpha$  and  $\beta$  that satisfy the goal.

If we substitute  $f + g$  or  $f \nabla g$  for  $h$  in the corresponding characterisation, the left-hand sides of the equivalences become trivially true. The right-hand sides are then also true, giving the *computation rules*:

$$\boxed{\begin{array}{ll} (f + g) \bullet \text{inl} = \text{inl} \bullet f & (f + g) \bullet \text{inr} = \text{inr} \bullet g \\ (f \nabla g) \bullet \text{inl} = f & (f \nabla g) \bullet \text{inr} = g \end{array}}$$

The validity of the so-called  $\nabla$ -*fusion* rule:

$$\boxed{h \bullet (f \nabla g) = (h \bullet f) \nabla (h \bullet g)}$$

is shown by the following calculation<sup>5</sup>:

$$\begin{aligned} h \bullet f \nabla g &= \alpha \nabla \beta \\ &\equiv \quad \{ \text{characterisation of } \nabla \} \\ h \bullet f \nabla g \bullet \text{inl} &= \alpha \wedge h \bullet f \nabla g \bullet \text{inr} = \beta \\ &\equiv \quad \{ \text{computation rules for } \nabla \} \\ h \bullet f &= \alpha \wedge h \bullet g = \beta \quad . \end{aligned}$$

<sup>5</sup> We adopt the convention that composition has lower precedence than all other operators. Thus  $h \bullet f \nabla g$  should be read as  $h \bullet (f \nabla g)$ . In the statement of the basic rules, however, we always parenthesise fully.

Note once again the style of calculation in which the right side of the law is constructed rather than verified.

It is also possible to express  $+$  in terms of  $\nabla$ , namely by:

$$\boxed{f + g = (\text{inl} \bullet f) \nabla (\text{inr} \bullet g)}$$

We derive the rhs of this rule as follows:

$$\begin{aligned} f+g &= \alpha \nabla \beta \\ \equiv & \quad \{ \text{characterisation of } \nabla \} \\ f+g \bullet \text{inl} &= \alpha \wedge f+g \bullet \text{inr} = \beta \\ \equiv & \quad \{ \text{computation rules for } + \} \\ \text{inl} \bullet f &= \alpha \wedge \text{inr} \bullet g = \beta \quad . \end{aligned}$$

Another fusion rule is the  $\nabla$ - $+$  *fusion* rule:

$$\boxed{(f + g) \bullet (h \nabla k) = (f \bullet h) \nabla (g \bullet k)}$$

We leave its derivation as an exercise.

These rules are useful by themselves, but they were proved to lead to the result that  $+$  respects function composition:

$$\boxed{(f + g) \bullet (h + k) = (f \bullet h) + (g \bullet k)}$$

The proof is simple:

$$\begin{aligned} & f+g \bullet h+k \\ = & \quad \{ \text{definition of } + \} \\ & (\text{inl} \bullet f) \nabla (\text{inr} \bullet g) \bullet h+k \\ = & \quad \{ \nabla\text{-}+\text{ fusion} \} \\ & (\text{inl} \bullet f \bullet h) \nabla (\text{inr} \bullet g \bullet k) \\ = & \quad \{ \text{definition of } + \} \\ & (f \bullet h) + (g \bullet k) \quad . \end{aligned}$$

The proof that  $+$  also respects  $\text{id}$ , that is,

$$\boxed{\text{id} + \text{id} = \text{id}}$$

is also left as an exercise.

An important property that we shall use is that the mapping  $\nabla$  is injective, that is:

$$f \nabla g = h \nabla k \quad \equiv \quad f = h \wedge g = k \quad .$$

Just two simple steps are needed for the proof. Note, in particular, that there is no need for separate “if” and “only if” arguments.

$$\begin{aligned}
& f \nabla g = h \nabla k \\
\equiv & \quad \{ \text{characterisation} \} \\
& f \nabla g \bullet \text{inl} = h \quad \wedge \quad f \nabla g \bullet \text{inr} = k \\
\equiv & \quad \{ \text{computation rules} \} \\
& f = h \wedge g = k \quad .
\end{aligned}$$

Further, the mapping is surjective (within the typing constraints): if  $h :: a + b \rightarrow c$ , then there exist functions  $f :: a \rightarrow c$  and  $g :: b \rightarrow c$  such that  $h = f \nabla g$ . In fact, they can be given explicitly by  $f = h \bullet \text{inl}$  and  $g = h \bullet \text{inr}$ .

*The product functor* While sums give a choice between values of two types, products combine two values. In Haskell the product type former and the pair constructor are syntactically equal. However, we want to distinguish between the type former  $\times$  and the value constructor  $(-,_-)$ . The binary product functor  $\times$  is given by:

**data**  $a \times b = (a, b)$

$\text{exl}(u, v) = u$   
 $\text{exr}(u, v) = v$

$f \times g = h$  **where**  
 $h(u, v) = (f\ u, g\ v)$

$f \triangle g = h$  **where**  
 $h\ u = (f\ u, g\ u)$

The functions  $\text{exl} :: a \times b \rightarrow a$  and  $\text{exr} :: a \times b \rightarrow b$  are called *projections* and are said to *project* a pair onto its components.

Just as for disjoint sum, we have defined  $f \times g$  in such a way that it meet the type requirements on a functor. Specifically, if  $f :: a \rightarrow b$  and  $g :: c \rightarrow d$  then  $f \times g :: a \times c \rightarrow b \times d$ , as is easily checked. Also, we have defined a second combination of  $f$  and  $g$ , namely  $f \triangle g$ , which we pronounce  $f$  “split”  $g$ .

The operational meaning of  $f \times g$  is easy to see. Given a pair of values, it produces a pair by applying  $f$  to the first component and  $g$  to the second component. The operational meaning of  $f \triangle g$  is that it constructs a pair of values from a single value by applying both  $f$  and  $g$  to the given value. (In particular,  $\text{id} \triangle \text{id}$  constructs a pair by “splitting” a given value into two copies of itself.)

A curious fact is the following. *All* the rules for sums are *also* valid for products under the following systematic replacements: replace  $+$  by  $\times$ ,  $\nabla$  by  $\triangle$ ,  $\text{inl}$  and  $\text{inr}$  by  $\text{exl}$  and  $\text{exr}$ , and switch the components  $f$  and  $g$  of each composition  $f \bullet g$ . (In category theory this is called *dualisation*.) This gives us the *characterisations*:

$h = f \times g$	$\equiv$	$\text{exl} \bullet h = f \bullet \text{exl} \wedge \text{exr} \bullet h = g \bullet \text{exr}$
$h = f \triangle g$	$\equiv$	$\text{exl} \bullet h = f \quad \wedge \quad \text{exr} \bullet h = g$



the *identity rule*:

$$\boxed{\text{exl} \triangle \text{exr} = \text{id}}$$

the *computation rules*:

$$\boxed{\begin{array}{ll} \text{exl} \bullet (f \times g) = f \bullet \text{exl} & \text{exr} \bullet (f \times g) = g \bullet \text{exr} \\ \text{exl} \bullet (f \triangle g) = f & \text{exr} \bullet (f \triangle g) = g \end{array}}$$

the  $\triangle$ -*fusion rule*:

$$\boxed{(f \triangle g) \bullet h = (f \bullet h) \triangle (g \bullet h)}$$

$\times$  expressed in terms of  $\triangle$ :

$$\boxed{f \times g = (f \bullet \text{exl}) \triangle (g \bullet \text{exr})}$$

the  $\times$ - $\triangle$ -*fusion rule*:

$$\boxed{(f \times g) \bullet (h \triangle k) = (f \bullet h) \triangle (g \bullet k)}$$

and finally the fact that  $\times$  is a *binary functor*:

$$\boxed{\begin{array}{lll} (f \times g) \bullet (h \times k) & = & (f \bullet h) \times (g \bullet k) \\ \text{id} \times \text{id} & = & \text{id} \end{array}}$$

*Functional Composition of Functors* It is easily verified that the composition of two unary functors  $F$  and  $G$  is also a functor. By their composition we mean the pair of mappings, the first of which maps type  $a$  to  $F(Ga)$  and the second maps function  $f$  to  $F(Gf)$ . We use juxtaposition—thus  $FG$ —to denote the composition of unary functors  $F$  and  $G$ . For example, *Maybe Rose* denotes the composition of the functors *Maybe* and *Rose*. The order of composition is important, of course. The functor *Maybe Rose* is quite different from the functor *Rose Maybe*.

It is also possible to compose functors of different arities. For instance we may want to compose a binary functor like disjoint sum with a unary functor like *List*. A simple notational device to help define such a functor is to overload the meaning of the symbol “+” and write *List+List*, whereby we mean the functor that maps  $x$  to  $(\text{List } x) + (\text{List } x)$ . Similarly we can compose disjoint sum with two unary functors  $F$  and  $G$ : we use the notation  $F+G$  and mean the functor that maps  $x$  to  $(F x) + (G x)$ .

Two ways of reducing the arity of a functor are *specialisation* and *duplication*. An example of specialisation is when we turn the binary disjoint sum functor into a unary functor by specialising its first argument to the unit type. We write  $1+\text{id}$  and mean the functor that maps type  $a$  to the type  $1+a$ , and function  $f$  to the function  $\text{id}_1+f$ . Duplication means

that we duplicate the argument as many times as necessary. For example, the mapping  $x \mapsto x+x$  is a unary functor.

Both duplication and specialisation are forms of functional composition of functors. To formulate them precisely we need to extend the notion of functor so that the arity of the target of a functor may be more than one. (Up till now we have always said that a functor maps an  $n$ -tuple of types/functions to a single type/function.) Then a tuple of functors is also a functor, and, for each  $n$ , there is a duplication functor of arity  $n$ . In this way duplication and specialisation can be expressed as the composition of a functor with a tuple of functors. (In the case of specialisation, one of the functors is a constant functor.)

For our current purposes, a complete formalisation is an unnecessary complication and the ad hoc notation introduced above will suffice. Formalisations can be found in [18, 19, 37, 22].

*Polynomial functors* A functor built only from constants, extractions, sums, products and composition is called a *polynomial* functor.

An example of a polynomial functor is *Maybe* introduced in section 2.2. Recalling its definition:

```
data Maybe a = none | one a
```

we see that, expressed in the notation introduced above,  $Maybe = 1+ld$

The remaining examples introduced in section 2.2 are not polynomial because they are defined recursively. We need one more mechanism for constructing functors. That is the topic of the next section.

**Exercise 2.6 ( $\nabla$ - $\Delta$  abide)** Prove that, for all  $f, g, h$  and  $k$ ,

$$(f \nabla g) \Delta (h \nabla k) = (f \Delta h) \nabla (g \Delta k) .$$

□

**Exercise 2.7 (Abide laws)** The law proved in exercise 2.6 is called the  $\nabla$ - $\Delta$  *abide* law because of the following two-dimensional way of writing the law in which the two operators are written either *above* or *beside* each other. (The two-dimensional way of writing is originally due to C.A.R.Hoare, the catchy name is due to Richard Bird.)

$$\begin{array}{ccc} f \nabla g & f & g \\ \Delta & = & \Delta \nabla \Delta \\ h \nabla k & h & k \end{array}$$

What other operators abide with each other in this way? (You have already seen examples in this text, but there are also other examples from simple arithmetic.)

□

**Exercise 2.8** Consider the mapping *Square* that takes a type  $a$  to  $a \times a$  and a function  $f$  to  $f \times f$ . Check that *Square* is a functor.

**Exercise 2.9** In checking that something is a functor, we must check that it respects composition *and* identity. The last part may not be omitted, as is shown by the existence of “almost-functors”. Call  $F$  an almost-functor when  $F$  is a pair of mappings on types and functions (just like true functors) that respects typing and composition, but fails to respect identity:  $Fid \neq id$ . Can you find a simple example of such an almost-functor? (Hint: Look at constant mappings.)

**Exercise 2.10** If  $inl :: a \rightarrow a + b$  and  $inr :: b \rightarrow a + b$ , what is the typing of  $id$  in the identity rule  $inl \vee inr = id$ ?

**Exercise 2.11** Complete the verification that  $+$  is a functor by proving the  $\vee$ - $+$  fusion rule and the identity rule ( $id + id = id$ ). In the calculation you may use all the other rules stated before these two rules.

□

## 2.4 Datatypes Generically

By now the notion of a functor should be becoming familiar to you. Also, it should be clear how to extend the definition of non-inductive datatypes not involving function spaces to a polynomial functor. In this section we take the step to inductively defined datatypes.

The basic idea is that an inductively defined datatype is a fixed point of a functor, which functor we call the *pattern* functor of the datatype. For the simplest examples (such as the natural numbers) the pattern functor is polynomial but for more complicated examples (like the *Rose* datatype) it is not. We therefore need to extend the class of functors we can define beyond the polynomial functors to the so-called *regular* functors by adding the *type* functors. The basic technical device to achieve this is the *catamorphism*, which is a generalisation of the *fold* function on lists.

We begin by discussing pattern functors following which we can, at long last, define the notion of an  $F$ -algebra. Catamorphisms form the next —substantial— topic, following which we introduce type functors and the class of regular functors.

*Pattern functors and recursion* We first look at a simple inductively (= recursively) defined datatype, that of the Peano naturals, which we also saw in section 2.1:

```
data Nat = zero | succ Nat
```

There is only one number zero, which we can make explicit by:

```
data Nat = zero 1 | succ Nat
```

Instead of fancy constructor function names like `succ` and `zero` we now employ boring standard ones:

**data**  $Nat = \text{inl } 1 \mid \text{inr } Nat$

The choice here is that afforded by sum, so we replace this by

**data**  $Nat = \text{in}(1 + Nat)$

in which there is one explicit constructor function left, called “in”.

Now note that  $Nat$  occurs both on the left and the right of the datatype definition (which is why it is called an inductively defined or recursive datatype). In order to view this as a fixed point definition, let us abstract from  $Nat$  on the right side replacing it by the variable  $z$ . In this way we are led to consider the unary functor  $N$  defined by

$N z = 1 + z$

(Note that, although we have only defined  $N$  explicitly on types, we understand its extension to a functor. Using the notations introduced earlier, this functor is expressed as  $N = 1^{\mathbb{K}} + \text{Id}$ .) The functor  $N$  captures the pattern of the inductive formation of the Peano naturals. The point is that we can use this to rewrite the definition of  $Nat$  to

**data**  $Nat = \text{in}(N Nat)$

Apparently, the *pattern functor*  $N$  uniquely determines the datatype  $Nat$ . Whenever  $F$  is a unary polynomial functor, as is the case here, a definition of the form **data**  $Z = \text{in}(F Z)$  uniquely determines  $Z$ .

We need a notation to denote the datatype  $Z$  that is obtained, and write  $Z = \mu F$ . So  $Nat = \mu N$ . Replacing  $Z$  by  $\mu F$  in the datatype definition, and adding a subscript to the single constructor function  $\text{in}$  in order to disambiguate it, we obtain:

**data**  $\mu F = \text{in}_F(F \mu F)$

Now  $\text{in}_F$  is a generic function, with typing

$\text{in}_F :: F \mu F \rightarrow \mu F$

We can “reconstruct” the original functions `zero` and `succ` by defining:

`zero` =  $\text{in}_N \bullet \text{inl} :: 1 \rightarrow Nat$

`succ` =  $\text{in}_N \bullet \text{inr} :: Nat \rightarrow Nat$

Conversely,  $\text{in}_N :: N Nat \rightarrow Nat$  is then of course

$\text{in}_N = \text{zero} \vee \text{succ}$

Playing the same game on the definition of  $List$  gives us:

**data**  $List a = \text{in}(1 + (a \times List a))$

Replacing the datatype being defined,  $List a$ , systematically by  $z$ , we obtain the “equation”

**data**  $z = \text{in}(1 + (a \times z))$

Thus, we see that the pattern functor here is  $(z \mapsto 1 + (a \times z))$ . It has a parameter  $a$ , which we make explicit by putting

$$L a = (z \mapsto 1 + (a \times z))$$

Now  $List a = \mu(L a)$ , or, abstracting from  $a$ :

$$List = (a \mapsto \mu(L a))$$

**Exercise 2.12** What is the pattern functor for *Bin*? Is it polynomial? What is the pattern functor for *Rose*? Is it polynomial?

□

*F-algebras* Before we traded in the names of the constructor functions for the uniform ‘in’, we saw that the algebra naturally corresponding to the datatype *Nat*, together with the generic algebra of its class, were:

$$(Nat, zero, succ), \text{ with } zero :: 1 \rightarrow Nat, \text{ succ} :: Nat \rightarrow Nat \\ (A, e, f), \text{ with } e :: 1 \rightarrow A, \text{ f} :: A \rightarrow A$$

Using ‘in’, this should be replaced by:

$$(Nat, in_N), \text{ with } in_N :: 1 + Nat \rightarrow Nat \\ (A, \varphi), \text{ with } \varphi :: 1 + A \rightarrow A$$

in which the relation between  $\varphi$  and the pair  $(e, f)$  is, of course,

$$\varphi = e \vee f \\ e = \varphi \bullet inl \\ f = \varphi \bullet inr$$

Using the pattern functor  $N$ , we can also write:

$$(Nat, in_N), \text{ with } in_N :: N Nat \rightarrow Nat \\ (A, \varphi), \text{ with } \varphi :: N A \rightarrow A$$

In general, for a functor  $F$ , an algebra  $(A, \varphi)$  with  $\varphi :: FA \rightarrow A$  is called an *F-algebra* and  $A$  is called the *carrier* of the algebra. So *Nat* is the carrier of an  $N$ -algebra, and likewise  $List a$  is the carrier of an  $(L a)$ -algebra.

*Catamorphisms* In the class of  $F$ -algebras, a homomorphism  $h :: (A, \varphi) \rightarrow (B, \psi)$  is a function  $h :: A \rightarrow B$  that satisfies:

$$h \bullet \varphi = \psi \bullet Fh$$

This can be expressed in a diagram:

$$\begin{array}{ccc} FA & \xrightarrow{\varphi} & A \\ Fh \downarrow & \text{☺} & \downarrow h \\ FB & \xrightarrow{\psi} & B \end{array}$$

The smiley face signifies that the diagram *commutes*: the two paths from  $FA$  to  $B$  are equivalent.

A specific example of such a diagram is given by the homomorphism `even` from the natural numbers to the booleans:

$$\begin{array}{ccc}
 1+Nat & \xrightarrow{\text{zero}\nabla\text{succ}} & Nat \\
 \downarrow 1+\text{even} & \text{☺} & \downarrow \text{even} \\
 1+Bool & \xrightarrow{\text{true}\nabla\text{not}} & Bool
 \end{array}$$

which expresses the equation

$$\text{even} \bullet (\text{zero} \nabla \text{succ}) = (\text{true} \nabla \text{not}) \bullet (1 + \text{even}) .$$

Rather than use such a diagram, the standard way of defining a function on an inductive datatype is by “pattern matching” on the argument, giving a clause for each constructor function. For the naturals, the typical definition has this form:

```

data Nat = zero | succ Nat

h zero    = e
h (succ n) = f (h n)
  
```

For example, the function `even` is defined by the equations:

```

even zero    = true
even (succ n) = not (even n)
  
```

(Exercise 2.13 asks you to show that these two equations are equivalent to the commuting diagram above.) For lists, the typical pattern-matching has the form:

```

data List a = nil | cons a (List a)

h nil        = e
h (cons x xs) = x ⊕ h xs
  
```

In these definitions, the function being defined,  $h$ , is “pushed down” recursively to the components to which the constructor functions are applied. The effect is to replace the constructor functions by the corresponding arguments in the definition of  $h$  — in the case of the natural numbers, `zero` is replaced by  $e$  and `succ` is replaced by  $f$ , and in the case of lists `nil` is replaced by  $e$  and `cons` is replaced by  $\oplus$ .

For the naturals, the function  $h$  defined above is determined *uniquely* by  $e$  and  $f$ . Likewise, for lists,  $h$  is *uniquely* determined by  $e$  and  $\oplus$ , and there is a standard notation for the function thus defined, namely `foldr`  $\oplus e$ . Generalizing this, we get the following:

**data**  $\mu F = \text{in}_F(F \mu F)$

$h(\text{in}_F x) = \varphi((F h) x)$

in which simple typing considerations show that  $\varphi$  has to have a typing of the form  $F A \rightarrow A$ , and then  $h$  has the typing  $\mu F \rightarrow A$ ; in other words,  $\varphi$  is the operation of some  $F$ -algebra whose carrier is the target type of  $h$ . The function  $h$  thus defined is uniquely determined by  $\varphi$ . We call such functions *catamorphisms* and use the following notation:  $h = \llbracket \varphi \rrbracket$ . So  $\llbracket \_ \rrbracket$  is defined by:

$\llbracket \varphi \rrbracket = h$  **where**  
 $h(\text{in}_F x) = \varphi((F h) x)$

In words, when catamorphism  $\llbracket \varphi \rrbracket$  is applied to a structure of type  $\mu F$ , this means it is applied recursively to the components of the structure, and the results are combined by applying its “body”  $\varphi$ . Specialised to lists, the  $\llbracket \_ \rrbracket$ -combinator becomes `foldr` restricted to finite lists. The importance of having generic catamorphisms is that they embody a closed expression for a familiar inductive definition technique and thereby allow the generic expression of important programming rules.

**Exercise 2.13** Show that the single equation

$\text{even} \bullet \text{zero} \nabla \text{succ} = \text{true} \nabla \text{not} \bullet 1 + \text{even}$

is equivalent to the two equations

$\text{even zero} = \text{true}$   
 $\text{even}(\text{succ } n) = \text{not}(\text{even } n)$  .

□

*Initial Algebras* Catamorphisms enjoy a number of attractive calculational properties which we now discuss.

We start with giving the *typing rule* for  $\llbracket \_ \rrbracket$ :

$$\boxed{\frac{\varphi :: F a \rightarrow a}{\llbracket \varphi \rrbracket :: \mu F \rightarrow a}}$$

Taking the definition

$h(\text{in}_F x) = \varphi((F h) x)$

we can rewrite this equivalently as:

$(h \bullet \text{in}_F) x = (\varphi \bullet F h) x$

or, abstracting from  $x$ :

$h \bullet \text{in}_F = \varphi \bullet F h$

This functional equation in  $h$  has a unique solution, so we conclude that  $([\varphi])$  is *characterised* by

$$\boxed{h = ([\varphi]) \equiv h \bullet \text{in}_F = \varphi \bullet Fh}$$

The right-hand side of this equivalence states that  $h$  is a homomorphism, and if  $A$  is the carrier of  $\varphi$ , we can also express this characterisation as:

$$h = ([\varphi]) \equiv h :: (\mu F, \text{in}_F) \rightarrow (A, \varphi)$$

In words, every  $F$ -algebra is the target algebra of a unique homomorphism with  $(\mu F, \text{in}_F)$  as its source algebra, and the catamorphisms consist of these unique homomorphisms. Source algebras that have the property that there is a unique homomorphism to any target algebra are known as *initial algebras*. So  $(\mu F, \text{in}_F)$  is an initial algebra. It is easy to prove that all initial algebras in a given algebra class are isomorphic.

The following diagram expresses the fact that  $([\varphi]) :: (\mu F, \text{in}) \rightarrow (A, \varphi)$  (but not the uniqueness):

$$\begin{array}{ccc}
 F \mu F & \xrightarrow{\text{in}} & \mu F \\
 \downarrow F([\varphi]) & \text{☺} & \downarrow ([\varphi]) \\
 FA & \xrightarrow{\varphi} & A
 \end{array}$$

In formula form we get the *computation rule* for catamorphisms:

$$\boxed{([\varphi]) \bullet \text{in} = \varphi \bullet F([\varphi])}$$

The function  $\text{in}$  is itself an  $F$ -algebra, so  $([\text{in}])$  is defined. What is it? By substituting  $(A, \varphi) := (\mu F, \text{in})$  in the last equivalence above, we obtain:

$$h = ([\text{in}]) \equiv h :: (\mu F, \text{in}) \rightarrow (\mu F, \text{in})$$

But we know that  $\text{id} :: (\mu F, \text{in}) \rightarrow (\mu F, \text{in})$ ! The conclusion is the *identity rule* for catamorphisms:

$$\boxed{([\text{in}]) = \text{id}_{\mu F}}$$

This generalises the equality for lists:  $\text{foldr cons nil} = \text{id}$ .



*Further properties of catamorphisms* The identity rule is easy to remember if one thinks of a catamorphism as a function that replaces the constructor functions of the datatype by the supplied arguments. Thus `foldr cons nil` is the identity function on lists because `cons` is replaced by `cons` and `nil` is replaced by `nil`. In general, `([in])` replaces all occurrences of `in` by itself in an element of the datatype  $\mu F$ .

The identity rule is surprisingly important. As an illustration of its importance, we prove that `in` is a bijection between  $\mu F$  and  $F\mu F$ . That is, we use the rule to construct a function `out` of type  $\mu F \rightarrow F\mu F$  such that `in•out = idμF` and `out•in = idFμF`. Our calculation starts with the first requirement and derives a candidate for `out` in a systematic way:

$$\begin{aligned}
& \text{in} \bullet \text{out} = \text{id}_{\mu F} \\
\equiv & \quad \{ \text{identity rule} \} \\
& \text{in} \bullet \text{out} = ([\text{in}]) \\
\equiv & \quad \{ \text{catamorphism characterisation} \} \\
& \text{in} \bullet \text{out} \bullet \text{in} = \text{in} \bullet F(\text{in} \bullet \text{out}) \\
\Leftarrow & \quad \{ \text{cancel in} \bullet \text{ from both sides} \} \\
& \text{out} \bullet \text{in} = F(\text{in} \bullet \text{out}) \\
\equiv & \quad \{ F \text{ respects composition} \} \\
& \text{out} \bullet \text{in} = F \text{in} \bullet F \text{out} \\
\equiv & \quad \{ \text{catamorphism characterisation} \} \\
& \text{out} = ([F \text{in}]) \quad .
\end{aligned}$$

This completes the first step in the calculation: we have derived the candidate `([F in])` for `out`.

Note that the identity rule is not used to *simplify* `([in])` to `idμF` in this calculation; rather, it is used in quite the opposite way to *complicate* `idμF` to `([in])`. There is a tendency to view algebraic properties as left-to-right rewrite rules, where the left side is the complicated side and the right side is its simplified form. Calculations that use the rules in this way are straightforward and do not require insight. On the other hand, calculations (such as the one above) which include at least one complication step are relatively difficult and do require insight. The importance of the identity rule for catamorphisms is its use in introducing a catamorphism into a calculation (see also the MAG system [38], in which identity catamorphisms are introduced in calculations in order to be able to apply fusion). It can require ingenuity to use because it involves replacing an identity function which is not visible. That is, a step in a calculation may involve replacing some composition  $f \bullet g$  by  $f \bullet ([\text{in}]) \bullet g$ , the invisible intermediate step being to replace  $f \bullet g$  by  $f \bullet \text{id}_{\mu F} \bullet g$ . This is valid if  $f$  has source  $\mu F$  (equivalently,  $g$  has target  $\mu F$ ) so it is important to be aware of the types of the quantities involved.

To complete the calculation we have to check that the candidate `([F in])` we have derived for `out` satisfies the second requirement on `out`. That is, we have to verify that `([F in])•in = idFμF`.

This is an exercise in the use of the computation rule which we leave to the reader (specifically, exercise 2.14).

As another illustration of the use of the properties of catamorphisms we derive a condition under which it is possible to fuse a post-composed function with a catamorphism. The goal of the calculation is to eliminate the catamorphism brackets from the equation.

$$\begin{aligned}
& h \bullet ([\varphi]) = ([\psi]) \\
\equiv & \quad \{ \text{characterisation of } ([\psi]) \} \\
& h \bullet ([\varphi]) \bullet \text{in} = \psi \bullet F(h \bullet ([\varphi])) \\
\equiv & \quad \{ \text{computation rule for } ([\varphi]) \} \\
& h \bullet \varphi \bullet F([\varphi]) = \psi \bullet F(h \bullet ([\varphi])) \\
\equiv & \quad \{ F \text{ respects composition} \} \\
& h \bullet \varphi \bullet F([\varphi]) = \psi \bullet F h \bullet F([\varphi]) \\
\Leftarrow & \quad \{ \text{cancel } \bullet F([\varphi]) \text{ from both sides} \} \\
& h \bullet \varphi = \psi \bullet F h \quad .
\end{aligned}$$

So we have derived the  $([\_])$ -fusion rule:

$$\boxed{h \bullet ([\varphi]) = ([\psi]) \Leftarrow h \bullet \varphi = \psi \bullet F h}$$

Note that the condition states that  $h$  is a homomorphism. So the rule states that composing a homomorphism after a catamorphism is a catamorphism.

The way this rule is typically used is that we want to fuse a given function  $h$  into a given catamorphism  $([\varphi])$ , for example to improve efficiency. In order to do so, we try to solve the equation  $h \bullet \varphi = \psi \bullet F h$  for the unknown  $\psi$ . If we find a solution, we know that the answer is  $([\psi])$ .

*An example* We show this in action on a simple example:  $\text{sum} \bullet \text{concat}$  on lists of lists of numbers. Recall that the pattern functor of *List Nat* is

$$L \text{ Nat} = (z \mapsto 1 + (\text{Nat} \times z)) \quad .$$

By definition,  $\text{concat} = ([\text{nil} \nabla (+)])$ , so we try to fuse  $\text{sum}$  and  $\text{concat}$  into a catamorphism. Applying the fusion rule we have:

$$\begin{aligned}
& \text{sum} \bullet \text{concat} = ([\psi]) \\
\Leftarrow & \quad \{ \text{concat} = ([\text{nil} \nabla (+)]), \text{ fusion} \} \\
& \text{sum} \bullet \text{nil} \nabla (+) = \psi \bullet (L \text{ Nat}) \text{sum} \quad .
\end{aligned}$$

Now, the pattern functor  $(L \text{ Nat})$  is a disjoint sum of two functors. Also, the composition on the left side can be fused together:

$$\begin{aligned}
& \text{sum} \bullet \text{nil} \nabla (+) \\
= & \quad \{ \nabla \text{ fusion} \} \\
& (\text{sum} \bullet \text{nil}) \nabla (\text{sum} \bullet (+)) \quad .
\end{aligned}$$

This suggests that we should try instantiating  $\psi$  to  $\alpha \nabla \beta$  for some  $\alpha$  and  $\beta$ . In this way, we get:

$$\begin{aligned}
& \text{sum} \bullet \text{concat} = (\psi) \\
\Leftarrow & \quad \{ \text{two steps above, definition of } (L \text{ Nat}) \} \\
& (\text{sum} \bullet \text{nil}) \nabla (\text{sum} \bullet (+)) = \psi \bullet (\text{id} + (\text{id} \times \text{sum})) \\
\equiv & \quad \{ \text{postulate } \psi = \alpha \nabla \beta, \text{ fusion} \} \\
& (\text{sum} \bullet \text{nil}) \nabla (\text{sum} \bullet (+)) = (\alpha \bullet \text{id}) \nabla (\beta \bullet \text{id} \times \text{sum}) \\
\equiv & \quad \{ \nabla \text{ is injective, simplification} \} \\
& \text{sum} \bullet \text{nil} = \alpha \wedge \text{sum} \bullet (+) = \beta \bullet \text{id} \times \text{sum} .
\end{aligned}$$

We now continue with each conjunct in turn. The first conjunct is easy:  $\text{sum} \bullet \text{nil} = \text{zero}$ . For the second conjunct, we have:

$$\begin{aligned}
& \text{sum} \bullet (+) \\
= & \quad \{ \text{property of summation} \} \\
& (+) \bullet \text{sum} \times \text{sum} \\
= & \quad \{ \times \text{ is a binary functor} \} \\
& (+) \bullet \text{sum} \times \text{id} \bullet \text{id} \times \text{sum} .
\end{aligned}$$

And thus we have found that  $\beta = (+) \bullet \text{sum} \times \text{id}$  satisfies  $\text{sum} \bullet (+) = \beta \bullet \text{id} \times \text{sum}$ .

Combining everything, we have found that

$$\text{sum} \bullet \text{concat} = (\text{zero} \nabla ((+) \bullet \text{sum} \times \text{id}))$$

or, expressed in a more familiar style:

$$\begin{aligned}
\text{sum} \bullet \text{concat} &= \text{foldr } \odot \ 0 \ \mathbf{where} \\
& \quad xs \odot y = \text{sum } xs + y
\end{aligned}$$

This derivation was not generic but specific for lists of lists. Meertens [37] shows how to do this generically, and also that the generic solution is no more complicated to obtain than this specific one, whilst being much more general.

**Exercise 2.14** We calculated above that  $\text{out} = (F \text{ in})$  satisfies  $\text{in} \bullet \text{out} = \text{id}_{\mu F}$ . Verify that  $\text{out} \bullet \text{in} = \text{id}_{F \mu F}$ .

□

**Exercise 2.15** Suppose that  $(A, \varphi)$  is an initial  $F$ -algebra. Prove that  $(A, \varphi)$  is isomorphic to  $(\mu F, \text{in}_F)$ . *Hint.* Consider the unique homomorphism  $h :: (A, \varphi) \rightarrow (\mu F, \text{in}_F)$ .

**Exercise 2.16** Consider the datatype  $\text{Bin } a$  for some arbitrary type  $a$ . The pattern functor for this type is  $F$  where  $Ff = \text{id}_a + (f \times f)$ . Catamorphisms over this type take the form  $(f \nabla \odot)$  where  $f$  is a function and  $\odot$  is a binary operator.

Define a catamorphism that counts the number of tips in a *Bin*. Define, in addition, a catamorphism that counts the number of joins in a *Bin*. Use the fusion rule for catamorphisms to determine a relation between the number of tips and the number of joins in a *Bin*. That is, derive the definition of a function  $f$  such that

$$f \cdot \text{NoOfTips} = \text{NoOfJoins} .$$

□

*Banana split* In this subsection we demonstrate the beauty of generic programming. We solve the following problem. Given are two catamorphisms  $([f]) :: \mu F \rightarrow a$  and  $([g]) :: \mu F \rightarrow b$ , and we need to have a function that returns the combined result of both. One solution is the program  $([f]) \triangle ([g])$ , but this can be inefficient since, computationally, the source data value is traversed twice, once for each of the two catamorphisms. So the question we want to solve is: can we combine these two into a single catamorphism  $([\chi])$ ?

This generic problem is motivated by our knowledge of specific cases. Take, for example, the problem of finding both the sum and the product of a list of numbers. The sum can of course be expressed as a catamorphism—it is the catamorphism  $([0 \nabla \text{add}])$ , where `add` is ordinary addition of real numbers—. Similarly the product function is a catamorphism, namely  $([1 \nabla \text{mul}])$ , where `mul` is the ordinary multiplication of real numbers. Equally obvious is that it should be possible to combine the sum and product of a list of numbers into one catamorphism. After all, the function  $sp = \text{sum} \triangle \text{product}$  is straightforward to express as a fold in Haskell:

$$\begin{aligned} sp &= \text{foldr } \odot \ e \ \mathbf{where} \\ &\quad x \odot (u, v) = (x + u, x \times v) \\ &\quad e \quad \quad = (0, 1) \end{aligned}$$

We can try to derive this special case in our calculus but more effective is to derive the solution to the generic problem. The benefit is not only that we then have a very general result that can be instantiated in lots of ways (one of which is the `sum` $\triangle$ `product` problem), but also that the derivation is much simpler because it omits irrelevant detail.

We begin the calculation of  $\chi$  as follows:

$$\begin{aligned} ([f]) \triangle ([g]) &= ([\chi]) \\ \equiv \quad \{ &\quad \text{There is a choice here. We can either use the} \\ &\quad \text{characterisation of } ([\chi]) \text{ or the characterisation of } f \triangle g. \\ &\quad \text{For no good reason, we choose the latter. } \} \\ ([f]) &= \text{exl} \cdot ([\chi]) \quad \wedge \quad ([g]) = \text{exr} \cdot ([\chi]) . \end{aligned}$$

This first step involves a difficult choice. At this point in time there is no reason why the use of one characterisation is preferable to the other (since both are equivalences). In fact, choosing to use the characterisation of  $([\chi])$  first does lead to a successful calculation of  $\chi$  of a similar length. We leave it as an exercise.

We now have to satisfy two conjuncts. Since the two conjuncts are symmetrical we proceed with just the first.

$$\begin{aligned}
& \llbracket f \rrbracket = \text{exl} \cdot \llbracket \chi \rrbracket \\
\Leftarrow & \quad \{ \text{Fusion} \} \\
& f \cdot F\text{exl} = \text{exl} \cdot \chi \\
\equiv & \quad \{ \bullet \quad \chi := \alpha \triangle \beta . \} \\
& f \cdot F\text{exl} = \text{exl} \cdot \alpha \triangle \beta \\
\equiv & \quad \{ \triangle \text{ computation} \} \\
& f \cdot F\text{exl} = \alpha .
\end{aligned}$$

The crucial step here (indicated by the bullet) is where we postulate the form of the solution, the motivation being the step that immediately follows.

In summary we have calculated that

$$\llbracket f \rrbracket = \text{exl} \cdot \llbracket \chi \rrbracket \Leftarrow \chi = \alpha \triangle \beta \wedge \alpha = f \cdot F\text{exl} .$$

Similarly,

$$\llbracket g \rrbracket = \text{exr} \cdot \llbracket \chi \rrbracket \Leftarrow \chi = \alpha \triangle \beta \wedge \beta = g \cdot F\text{exr} .$$

Putting everything together, we conclude that

$$\llbracket f \rrbracket \triangle \llbracket g \rrbracket = \llbracket (f \cdot F\text{exl}) \triangle (g \cdot F\text{exr}) \rrbracket .$$

This is affectionately called the *banana-split* theorem (because the brackets denoting a catamorphism look like bananas, and the  $\triangle$  operator is pronounced “split”).

**Exercise 2.17** Calculate  $\chi$  but start by using the characterisation of  $\llbracket f \rrbracket$ . In other words, calculate  $\chi$  as a solution of the equation

$$\llbracket f \rrbracket \triangle \llbracket g \rrbracket \cdot \text{in} = \chi \cdot F(\llbracket f \rrbracket \triangle \llbracket g \rrbracket) .$$

(You may find that you get a solution that is equivalent to the one above but not syntactically identical.)

□

*Type functors* In general, a binary functor gives rise to a new functor by a combination of parameterisation and constructing an initial algebra. For example, the binary pattern functor  $L$  that maps  $x$  and  $y$  to  $1+(x \times y)$  gives rise to the functor *List*. Such functors are called *type functors*. Here we show how this is done.

For greater clarity we will use an infix notation for binary functors. Suppose that  $\circledast$  is a binary functor, which we write as an infix operator. That is, for types  $a$  and  $b$ ,  $a \circledast b$  is a type and, for functions  $f :: a \rightarrow b$  and  $g :: c \rightarrow d$ ,  $f \circledast g$  is a function of type  $a \circledast c \rightarrow b \circledast d$ . Suppose  $a$  is an arbitrary type. Then the pair of mappings  $b \mapsto a \circledast b$  and  $f \mapsto \text{id}_a \circledast f$

is a functor (the functor formed by specialising the first operand of  $\otimes$  to the type  $a$ ). We denote this functor by  $(a\otimes)$  and call it a *parameterised* functor.

Now, since  $(a\otimes)$  is a unary functor, we can consider an initial  $(a\otimes)$ -algebra with carrier  $\mu(a\otimes)$ . Abstracting from  $a$  we have constructed a mapping from types to types. Let us introduce a special notation for this mapping:

$$\boxed{\tau(\otimes) = (a \mapsto \mu(a\otimes))}$$

So  $List = \tau(L)$ , with  $L$  the binary functor defined above.

For  $\tau(\otimes)$  to be a functor, we need, in addition to the action on types, an action on functions, which has to satisfy, for a function  $f :: a \rightarrow b$ ,

$$\tau(\otimes) f :: \tau(\otimes) a \rightarrow \tau(\otimes) b .$$

We derive a candidate for  $\tau(\otimes) f$  from type considerations. In the calculation, catamorphisms are  $(a\otimes)$  catamorphisms and  $\text{in}_{b\otimes}$  is an initial  $(b\otimes)$ -algebra.

$$\begin{aligned} & \tau(\otimes) f :: \tau(\otimes) a \rightarrow \tau(\otimes) b \\ \equiv & \quad \{ \text{definition of } \tau(\otimes) \text{ on types} \} \\ & \tau(\otimes) f :: \mu(a\otimes) \rightarrow \mu(b\otimes) \\ \Leftarrow & \quad \{ \bullet \tau(\otimes) f := \llbracket \varphi \rrbracket, \text{ typing rule for } \llbracket - \rrbracket \} \\ & \varphi :: a \otimes \mu(b\otimes) \rightarrow \mu(b\otimes) \\ \Leftarrow & \quad \{ \bullet \varphi := \text{in}_{b\otimes} \cdot \psi, \text{ type of in} \} \\ & \psi :: a \otimes \mu(b\otimes) \rightarrow b \otimes \mu(b\otimes) \\ \Leftarrow & \quad \{ f :: a \rightarrow b, \quad \text{id}_{\mu(b\otimes)} :: \mu(b\otimes) \rightarrow \mu(b\otimes), \\ & \quad \otimes \text{ respects typing} \} \\ & \psi = f \otimes \text{id}_{\mu(b\otimes)} . \end{aligned}$$

Performing the collected substitutions gives us this candidate definition

$$\boxed{\tau(\otimes) f = \llbracket \text{in}_{b\otimes} \cdot (f \otimes \text{id}_{\mu(b\otimes)}) \rrbracket}$$

Exercise 2.20 is to show that  $\tau(\otimes)$  respects composition and identities. According to the notational convention introduced earlier the action of  $\tau(\otimes)$  on functions can also be written  $\text{map}_{\tau(\otimes)}$ .

A final comment: The parameter  $a$  in a parameterised functor may actually be an  $n$ -tuple if functor  $\otimes$  is  $(n + 1)$ -ary, and then  $\tau(\otimes)$  is an  $n$ -ary functor. However, we only consider unary type functors, derived with  $\tau(\otimes)$  from binary functors in these lectures.

**Exercise 2.18** Consider the datatype  $Bool = \mu((1+1)^k)$ . Define  $false = \text{in}_{Bool} \cdot \text{inl}$ ,  $true = \text{in}_{Bool} \cdot \text{inr}$ . Examine and explain the meaning of the catamorphism  $\llbracket u \nabla v \rrbracket$  for  $Bool$ .

**Exercise 2.19 (cata-map fusion)** Derive a fusion rule of the form

$$([f]) \bullet (\tau(\circlearrowleft) g) = ([h]) .$$

Hint: instantiate the fusion rule for catamorphisms with  $F := (b\circlearrowleft)$ . Note also that  $\tau(\circlearrowleft)g$  is a catamorphism.

□

**Exercise 2.20** Complete the verification of the fact that  $\tau(\circlearrowleft)$  is a functor by showing that  $\tau(\circlearrowleft) \text{id}_a = \text{id}_{\tau(\circlearrowleft)a}$  and  $\tau(\circlearrowleft) (f \bullet g) = (\tau(\circlearrowleft) f) \bullet (\tau(\circlearrowleft) g)$ . (Hint: make use of exercise 2.19.)

□

**Exercise 2.21** Specialise the definition of  $\tau(\circlearrowleft) f$  for  $\circlearrowleft = L$ , the bifunctor giving the type functor  $List = \tau(L)$ , using  $\text{in} = \text{nil} \nabla \text{cons}$ , and verify that this is the familiar `map` function for lists. Also, instantiate your solution to exercise 2.19 and use it to express the sum of the squares of a list of numbers as a catamorphism. (That is, express the sum of a list of numbers as a catamorphism, and the list of squares of a list on numbers as a map. Then fuse the two functions together.)

*Regular Functors and Datatypes* We are now in a position to complete our discussion of the datatypes introduced in section 2.2 by giving a complete analysis of the definition of the *Rose* datatype. As we saw in exercise 2.12, its pattern functor is  $a\circlearrowleft z = a \times (List\ z)$ , or, in terms of the extraction functors `Par` and `Rec`,  $(\circlearrowleft) = \text{Par} \times (List\ \text{Rec})$ , which is not a polynomial functor, because of the appearance of the type functor *List*. Yet  $\tau(\circlearrowleft)$  is well defined. Incorporating type functors into the ways of constructing functors extends the class of polynomial functors to the class of *regular* functors.

A functor built only from constants, extractions, sums, products, composition and  $\tau()$  is called a *regular* functor. All the datatypes we have seen, including *List* and *Rose* are regular functors, and their constructor functions (combined together using the  $\nabla$  combinator) are initial algebras with respect to the pattern functors of the datatype.

This concludes the theory development. We have shown precisely what it means to say that a datatype is both an algebra and a functor.

## 2.5 A Simple Polytypic Program

We began section 2.2 with four representative examples of datatypes: *List*, *Maybe*, *Bin* and *Rose*. For each of these datatypes we can define a summation function that sums all the values stored in an instance of the datatype — assuming the values are numbers. Here is how one would do that in a non-generic programming style.

$$\begin{aligned}
\text{sum}_{List} \text{ nil} &= 0 \\
\text{sum}_{List} (\text{cons } u \text{ } us) &= u + \text{sum}_{List} us \\
\\
\text{sum}_{Maybe} \text{ none} &= 0 \\
\text{sum}_{Maybe} (\text{one } u) &= u \\
\\
\text{sum}_{Bin} (\text{tip } u) &= u \\
\text{sum}_{Bin} (\text{join } x \text{ } y) &= \text{sum}_{Bin} x + \text{sum}_{Bin} y \\
\\
\text{sum}_{Rose} (\text{fork } u \text{ } rs) &= u + \text{sum}_{List} (\text{map}_{List} \text{sum}_{Rose} rs)
\end{aligned}$$

We now want to replace all these definitions by a single generic definition  $\text{sum}_F$  for arbitrary unary functor  $F$ , which can be specialised to any of the above datatype constructors and many more by taking  $F$  to be  $List$ ,  $Maybe$ ,  $Bin$ , and so on. We do this by induction on the structure of the regular functors. That is, we define summation for a constant functor, for the extraction functors, for the composition of functors, for disjoint sum and cartesian product, and finally for a type functor. Let us begin with the type functors since this is where we see how to formulate the induction hypothesis.

For the type functor  $\pi(\otimes)$ , the requirement is to construct a function  $\text{sum}_{\pi(\otimes)}$  of type  $\mu(\mathbb{N}\otimes) \rightarrow \mathbb{N}$ . The obvious thing to do here is to define  $\text{sum}$  as a catamorphism,  $([f])$  say. In that case, the type requirement on  $f$  is that  $f :: \mathbb{N}\otimes\mathbb{N} \rightarrow \mathbb{N}$ . Note that the two arguments to the binary functor  $\otimes$  are both  $\mathbb{N}$ . This suggests the inductive hypothesis that there is a  $\text{sum}$  function of type  $F \mathbb{N} \rightarrow \mathbb{N}$  for all *unary* regular functors  $F$  obtained from an arbitrary non-constant  $n$ -ary regular functor by copying the (single) argument  $n$  times. We also need to define  $\text{sum}$  for the constant functor  $1$ . With this preparation, we can begin the analysis.

For the constant functor  $1$ , we define

$$\text{sum}_1 = 0 \text{ .}$$

This is because the sum of zero numbers is zero.

For the extraction functors, it is clear that

$$\text{sum}_{\text{Ex}} = \text{id}_{\mathbb{N}}$$

since the sum of a single number is that number itself.

For disjoint sum and cartesian product, we have:

$$\text{sum}_{F+G} = \text{sum}_F \nabla \text{sum}_G \quad \text{and} \quad \text{sum}_{F \times G} (x, y) = \text{sum}_F x + \text{sum}_G y \text{ .}$$

In the case of disjoint sum, either the  $\text{sum}_F$  function has to be applied, or the  $\text{sum}_G$  function, depending on the type of the argument. In the case of cartesian product, an element of an  $F \times G$  structure is a pair consisting of an element of an  $F$  structure and an element of a  $G$  structure, and the two sums have to be added together.

For the composition of two functors, we have:

$$\text{sum}_{FG} = \text{sum}_F \bullet F \text{sum}_G \text{ .}$$



Here the argument is that an  $FG$  structure is an  $F$  structure of  $G$  structures. The function  $F\text{sum}_G$  applies  $\text{sum}_G$  to all the individual  $G$  structures, and then  $\text{sum}_F$  adds their values.

The final case is a type functor, which we have already discussed.

$$\text{sum}_{\tau(\emptyset)} = (\text{sum}_{\emptyset}) .$$

We leave it to the reader to check that application of the above rules results in the particular instances of  $\text{sum}$  given above.

### 3 PolyP

The previous chapter introduces datatypes and functions on datatypes such as the catamorphism. The formal language used to introduce datatypes and functions is the language of category theory. The language of category theory is not a programming language, and although the accompanying text mentions programming, it is impossible to ‘run’ catamorphisms. This chapter introduces PolyP, a programming language with which generic functions such as the catamorphism can be implemented. The name of PolyP is derived from ‘polytypic programming’, an alternative name for generic programming.

PolyP is an extension of (a subset of) the functional programming language Haskell. The extension introduces a new kind of (top level) definition, the `polytypic` construct, which is used to define functions by induction over pattern functors, which describe the structure of (a subset of) regular datatypes. PolyP is based on the initial algebra approach to datatypes and work in the Squiggol community on datatypes. It is a tool that supports polytypic programming, and as such it has spurred the development of new polytypic programs.

In Haskell, datatypes are defined by means of the `data` construct, examples of which have been given in chapter 2. PolyP extracts the pattern functor from a datatype definition, and uses this structure information to instantiate generic programs on particular datatypes. We will use the name *polytypic function* for a generic program in PolyP.

PolyP has a number of limitations. The datatypes PolyP can handle are a subset of the datatypes induced by the regular functors defined in the previous chapter: PolyP’s pattern functors are binary and the type functors are unary which means that it can only handle datatypes with one type argument. Furthermore, datatypes cannot be mutually recursive.

Information about PolyP and polytypic programming in general can be found on

<http://www.cs.chalmers.se/~patrikj/poly/>

The names of pattern functors in PolyP differ slightly from the names in the previous chapter. Section 3.1 introduces PolyP’s functor names. Section 3.2 gives an implementation of the polytypic function `sum` from section 2.5. Section 3.3 defines most of the basic polytypic concepts. Type checking of polytypic functions is explained in section 3.4. Since we will use a number of polytypic functions in the rest of these notes, section 3.5 gives more examples of polytypic functions, and section 3.6 introduces PolyLib: a library of polytypic functions.

#### 3.1 Regular Functors in PolyP

The previous chapter explains how datatypes are defined by means of pattern functors. A pattern functor is a regular functor, i.e., a polynomial functor possibly extended with a type functor. PolyP defines polytypic functions by induction over the pattern functor of a datatype. The names for the pattern functors constructors used in PolyP differs slightly from the names in the previous chapter. This section defines the syntax for pattern functors used in PolyP.

PolyP’s functors are specified by the following context-free grammar:

$f, g ::= f + g \mid f * g \mid \text{Empty} \mid \text{Par} \mid \text{Rec} \mid d @ g \mid \text{Const } t$

The following table relates this syntax to the functors introduced in the previous chapter.

+	*	Empty	Par	Rec	$d @ g$	Const $t$
+	$\times$	$1^k$	exl	exr	$a \mapsto b \mapsto d(g \ a \ b)$	$t^k$

$+$  and  $*$  are the standard sum and product functors lifted to act on functors. `Empty` is the constant binary version of functor  $1^k$ . `Par` and `Rec` are mentioned in chapter 2, and are `exl` and `exr`, respectively. Composition of functors  $d$  and  $g$  is denoted by  $d @ g$  and is only defined for a unary functor  $d$  and a binary functor  $g$ . Finally, `Const  $t$`  is the binary variant of  $t^k$ . The  $t$  stands for a monotype such as `Bool`, `Char` or `(Int, [Float])`.

In PolyP, as in Haskell, type functors (recursive datatypes) are introduced by the `data` construct. Every Haskell datatype constructor  $d$  is equal to  $\pi(f)$  for some pattern functor  $f$ . In PolyP this  $f$  is denoted by `FunctorOf  $d$` . A datatype  $d \ a$  is regular (satisfies `Regular  $d$` ) if it contains no function spaces, and if the argument of the type constructor  $d$  is the same on the left- and right-hand side of its definition. For each one parameter regular datatype  $d \ a$ , PolyP automatically generates `FunctorOf  $d$`  using roughly the same steps as those used manually in section 2.4. For example, for

```
data Error a = Error String | Ok a
data List a  = Nil | Cons a (List a)
data Bin a   = Tip a | Join (Bin a) (Bin a)
data Rose a  = Fork a (List (Rose a))
```

PolyP generates the following functors:

```
FunctorOf Error = Const String + Par
FunctorOf List  = Empty + Par * Rec
FunctorOf Bin   = Par + Rec * Rec
FunctorOf Rose  = Par * (List @ Rec)
```

Pattern functors are *only* constructed for datatypes defined by means of the `data` construct. If somewhere in a program a polytypic function is applied to a value of type `Error (List a)`, PolyP will generate an instance of the polytypic function on the datatype `Error b`, not on the type `(Error @ List) a`. This also implies that the functor  $d$  in the functor composition  $d @ g$  is always a type functor.

### 3.2 An Example: psum

PolyP introduces a new construct `polytypic` for defining polytypic functions by induction on the structure of a binary pattern functor:

$$\text{polytypic } p :: t = \text{case } f \text{ of } \{f_i \rightarrow e_i\}$$

where  $p$  is the name of the value being defined,  $t$  is its type,  $f$  is a functor variable,  $f_i$  are functor patterns and  $e_i$  are PolyP expressions. The explicit type in the `polytypic` construct is needed since we cannot in general infer the type from the cases.

The informal meaning is that we define a function that takes (a representation of) a pattern functor as its first argument. This function selects the expression in the first branch of the case matching the functor, and the expression may in turn use the polytypic function (on subfunctors). Thus the polytypic construct is a (recursive) template for constructing instances of polytypic functions given the pattern functor of a datatype. The functor argument of the polytypic function need not (and cannot) be supplied explicitly but is inserted by the compiler during type inference.

```

psumd :: Regular d => d Int -> Int
psumd = catad fsumFunctorOf d

polytypic fsumf :: f Int Int -> Int
= case f of
  g + h -> fsumg 'either' fsumh
  g * h -> \ (x,y) -> fsumg x + fsumh y
  Empty -> \x -> 0
  Par -> id
  Rec -> id
  d @ g -> psumd . (pmapd fsumg)
  Const t -> \x -> 0

```

Fig. 1. The definition of psum

As an example we take the function `psum` defined in figure 1. (The subscripts indicating the type are included for readability and are not part of the definition.) Function `psum` sums the integers in a structure with integers. It is the PolyP implementation of the function `sum` defined in section 2.5. The function `either :: (a -> c) -> (b -> c) -> Either a b -> c` (corresponding to  $\vee$ ) and datatype `Either a b` (corresponding to  $a+b$ ) are defined in Haskell's prelude. The definition of functions `cata` and `pmap` (the implementations in PolyP of the catamorphism and the map, see chapter 2) will be given later. When `psum` is used on an element of type `Bin Int`, the compiler performs roughly the following rewrite steps to construct the actual instance of `psum` for `Bin`:

$$\text{psum}_{\text{Bin}} \rightarrow \text{cata}_{\text{Bin}} \text{fsum}_{\text{FunctorOf Bin}}$$

It follows that we need an instance of `cata` for the type functor `Bin`, and an instance of function `fsum` for the pattern functor `FunctorOf Bin = Par + Rec * Rec`. For the latter instance, we use the definition of `fsum` to transform as follows:

$$\text{fsum}_{\text{FunctorOf Bin}} \rightarrow \text{fsum}_{\text{Par+Rec*Rec}} \rightarrow \text{fsum}_{\text{Par}} \text{'either'} \text{fsum}_{\text{Rec*Rec}}$$

We transform the functions `fsumPar` and `fsumRec*Rec` separately. For `fsumPar` we have

$$\text{fsum}_{\text{Par}} \rightarrow \text{id}$$

and for `fsumRec*Rec` we have

```
fsumRec*Rec
→ \ (x,y) -> fsumRec x + fsumRec y
→ \ (x,y) -> id x + id y
```

The last function can be rewritten into `uncurry (+)`, and thus we obtain the following function for summing a tree:

```
cataBin (id 'either' (uncurry (+)))
```

By expanding `cataBin` in a similar way we obtain a Haskell function for the instance of `psum` on `Bin`. The function we obtain is the same as the function `sumBin` defined in section 2.5.

### 3.3 Basic Polytypic Functions

In the definition of function `psum` we used functions like `cata` and `pmap`. This subsection defines these and other basic polytypic functions.

Since polytypic functions cannot refer to constructor names of specific datatypes, we introduce the predefined functions `out` and `inn`. Function `out` is used in polytypic functions instead of pattern matching on the constructors of a datatype. For example `out` on `Bin` is defined as follows:

```
outBin :: Bin a -> Either a (Bin a, Bin a)
outBin (Tip x) = Left x
outBin (Join l r) = Right (l,r)
```

Function `inn` is the inverse of function `out`. It collects the constructors of a datatype into a single constructor function.

```
out :: Regular d => FunctorOf d a (d a) <- d a
inn :: Regular d => FunctorOf d a (d a) -> d a
```

Function `inn` is an implementation of `in` from chapter 2. The following calculation shows that the type of `inn` really corresponds to the type of `in`:

```
FunctorOf d a (d a) -> d a
= {      d = τ(f) for some regular functor f. }
  FunctorOf τ(f) a (τ(f) a) -> τ(f) a
= {      Definition of FunctorOf }
  f a (τ(f) a) -> τ(f) a
= {      Definition of τ() }
  f a μ(f a) -> μ(f a)
```

PolyP generates definitions of `inn` and `out` for all datatypes.

As explained in chapter 2, a functor is a mapping between categories that preserves the algebraic structure of the category. Since a category consists of objects (types) and arrows

(functions), a functor consists of two parts: a definition on types, and a definition on functions. A pattern functor  $f$  in PolyP is a function that take two types and return a type. The part of the functor that takes two functions and returns a function is called  $fmap_f$ , see figure 2.

```

polytypic fmap_f :: (a -> c) -> (b -> d) -> f a b -> f c d
= \p r -> case f of
    g + h -> (fmap_g p r) +- (fmap_h p r)
    g * h -> (fmap_g p r) *- (fmap_h p r)
    Empty -> id
    Par -> p
    Rec -> r
    d @ g -> pmap_d (fmap_g p r)
    Const t -> id

(-*-) :: (a -> c) -> (b -> d) -> (a,b) -> (c,d)
(f *- g) (x,y) = (f x , g y)

(+--) :: (a -> c) -> (b -> d) -> Either a b -> Either c d
(f +- g) = either (Left . f) (Right . g)

```

Fig. 2. The definition of  $fmap$ .

Using  $fmap$  we can define the polytypic version of function map,  $pmap$ , as follows:

```

pmap    :: Regular d => (a -> b) -> d a -> d b
pmap f = inn . fmap f (pmap f) . out

```

where  $out$  takes the argument apart,  $fmap$  applies  $f$  to parameters and  $(pmap f)$  recursively to substructures and  $inn$  puts the parts back together again. Function  $pmap_d$  is the function action of the type functor  $d$ .

Function  $cata$  is also defined in terms of function  $fmap$ :

```

cata :: Regular d => (FunctorOf d a b -> b) -> (d a -> b)
cata f = f . fmap id (cata f) . out

```

Note that this definition is a copy of the computation rule for the catamorphism in section 2.4, with  $in$  on the left-hand side replaced by  $out$  on the right-hand side.

### 3.4 Type Checking Polytypic Functions

We want to be sure that functions generated by polytypic functions are type correct, so that no run-time type errors occur. For that purpose PolyP type checks definitions of polytypic functions. This subsection briefly discusses how to type check polytypic functions, the details of the type checking algorithm can be found in [25].

Functor expressions contain  $+$ ,  $*$ , etc., and such expressions have to be translated to real types. For this translation we interpret functor constructors as type synonyms:

```

type (f + g) a b = Either (f a b) (g a b)
type (f * g) a b = (f a b , g a b)
type Empty a b   = ()
type Par a b     = a
type Rec a b     = b
type (d @ g) a b = d (g a b)
type Const t a b = t

```

So, for example, interpreting the functors in the pattern functor for `List` as type synonyms, we have:

```

FunctorOf List a b
= {      FunctorOf List = Empty + Par * Rec }
  (Empty + Par * Rec) a b
= {      Type synonym for + }
  Either (Empty a b) ((Par * Rec) a b)
= {      Type synonyms for Empty and * }
  Either () (Par a b, Rec a b)
= {      Type synonyms for Par and Rec }
  Either () (a,b)

```

To infer the type of a polytypic definition from the types of the expressions in the `case` branches, higher-order unification would be needed. As general higher-order unification is undecidable we require inductive definitions of polytypic functions to be explicitly typed, and we only check that this type is valid. Given an inductive definition of a polytypic function

```

polytypic foo :: ... f ...
= case f of
  g + h -> bar
  ...

```

where `f` is a functor variable, the rule for type checking these definitions checks among other things that the declared type of function `foo`, with `g + h` substituted for `f`, is an instance of the type of expression `bar`. For all of the expressions in the branches of the `case` it is required that the declared type is an instance of the type of the expression in the branch with the left-hand side of the branch substituted for `f` in the declared type. The expression `g + h` is an abstraction of a type, so by substituting `g + h` (or any of the other abstract type expressions) for `f` in the type of `foo` we mean the following: substitute `g + h` for `f`, and rewrite the expression obtained thus by interpreting the functor constructors as type synonyms. As an example we take the case `g * h` in the definition of `fsum`:

```

polytypic fsum :: f Int Int -> Int
= case f of
  ...
  g * h -> \(x,y) -> fsum x + fsum y
  ...

```

The type of the expression  $\backslash(x,y) \rightarrow \text{fsum } x + \text{fsum } y$  is  $(r \text{ Int Int}, s \text{ Int Int}) \rightarrow \text{Int}$ . Substituting the functor to the left of the arrow in the case branch,  $g * h$ , for  $f$  in the declared type  $f \text{ Int Int} \rightarrow \text{Int}$  gives  $(g * h) \text{ Int Int} \rightarrow \text{Int}$ , and rewriting this type using the type rewrite rules, gives  $(g \text{ Int Int}, h \text{ Int Int}) \rightarrow \text{Int}$ . This type is  $\alpha$ -convertible to (and hence certainly an instance of) the type of the expression to the right of the arrow in the case branch, so this part of the polytypic function definition is type correct.

### 3.5 More Examples of Polytypic Functions

This section describes some polytypic functions that will be used in the sequel. These functions can be found in PolyLib, the library of PolyP. The next section gives an overview of PolyLib.

Function `flatten` takes a value of type `d a` and flattens it into a list of values of type `[a]`. It is defined using function `fflatten :: f a [a] -> [a]`, which takes a value `v` of type `f a [a]`, and returns the concatenation of all the values (of type `a`) and lists (of type `[a]`) occurring at the top level in `v`. The definition of `flatten` and `fflatten` is given in figure 3. As an example, we unfold the definition of `fflatten` when used on the type `List`

```

flattend :: Regular d => d a -> [a]
flattend = catad fflattenFunctorOf d

polytypic fflattenf :: f a [a] -> [a]
= case f of
  g + h -> either fflatteng fflattenh
  g * h -> \ (x,y) -> fflatteng x ++ fflattenh y
  Empty -> nil
  Par -> singleton
  Rec -> id
  d @ g -> concat . flattend . pmapd fflatteng
  Const t -> nil

nil x = []
singleton x = [x]

```

Fig. 3. The definition of `flatten` and `fflatten`.

a (remember that `FunctorOf List = Empty+Par*Rec`):

```

fflattenEmpty+Par*Rec
→ either fflattenEmpty fflattenPar*Rec
→ either nil (\ (x,y) -> fflattenPar x ++ fflattenRec y)
→ either nil (\ (x,y) -> id x ++ id y)
→ either nil (uncurry (++))

```

The expression `pequal eq x y` checks whether or not the values `x` and `y` are equivalent using the equivalence operator `eq` to compare the elements pairwise. It is defined in terms



of function `fequal eq (pequal eq)`, where the first argument, `eq`, compares parameters for equality and the second argument, `(pequal eq)`, compares the subterms recursively. The third and fourth arguments are the two (unfolded) terms to be compared. These functions are defined in figure 4.

```

polytypic fequalf :: (a -> b -> Bool) -> (c -> d -> Bool) ->
                    f a c -> f b d -> Bool
= \p r -> case f of
    g + h -> sumequal (fequalg p r) (fequalh p r)
    g * h -> prodequal (fequalg p r) (fequalh p r)
    Empty -> \_ _ -> True
    Par -> p
    Rec -> r
    d @ g -> pequald (fequalg p r)
    Const t -> (==)

pequal :: (a -> b -> Bool) -> d a -> d b -> Bool
pequal eq x y = fequal eq (pequal eq) (out x) (out y)

sumequal :: (a -> b -> Bool) -> (c -> d -> Bool) ->
           Either a c -> Either b d -> Bool
sumequal f g (Left x) (Left v) = f x v
sumequal f g (Right y) (Right w) = g y w
sumequal f g _ _ = False

prodequal :: (a -> b -> Bool) -> (c -> d -> Bool) ->
           (a,c) -> (b,d) -> Bool
prodequal f g (x,y) (v,w) = f x v && g y w

```

Fig. 4. The definition of `pequal` and `fequal`.

### 3.6 PolyLib: a Library of Polytypic Functions

Using different versions of PolyP (and its predecessors) we have implemented a number of polytypic programs. For example, we have implemented a polytypic equality function, a polytypic show function, and a polytypic parser. Furthermore, we have implemented some more involved polytypic programs for pattern matching, unification and rewriting. These polytypic programs use several basic polytypic functions, such as the relatively well-known `cata` and `pmap`, but also less well-known functions such as `propagate` and `thread`. We have collected these basic polytypic functions in the library of PolyP: PolyLib [27, app. B]. This paper describes the polytypic functions in PolyLib, motivates their presence in the library, and gives a rationale for their design. This section first introduces the format used for describing polytypic library functions, then it gives an overview of the contents of the library, followed by a description of each of the submodules in the library.

**Describing Polytypic Functions** The description of a polytypic function consists of (some of) the following components: its name and type; an (in)formal description of the function; other names the function is known by; known uses of the function; and its background and relationship to other polytypic functions. For example:

---

```
pmap :: (a -> b) -> d a -> d b
```

Function `pmap` takes a function `f` and a value `x` of datatype `d a`, and applies `f ...`

**Also known as:** `map` [31], `mapn` [29]. **Known uses:** Everywhere! **Background:** This was one of the first ...

---

A problem with describing a library of polytypic functions is that it is not completely clear how to *specify* polytypic functions. The most basic combinators have immediate category theoretic interpretations that can be used as a specification, but for more complicated combinators the matter is not all that obvious. Thus, we will normally not provide formal specifications of the library functions, though we try to give references to more in-depth treatments.

The polytypic functions in the library are only defined for regular datatypes `d a`. In the type this is indicated by adding a context `Regular d => ...`, but we will omit this for brevity.

**Library Overview** We have divided the library into six parts, see figure 5. The first part of the library contains powerful recursion combinators such as `map`, `cata` and `ana`. This part is the core of the library in the sense that it is used in the definitions of all the functions in the other parts. The second part deals with zips and some derivatives, such as the equality function. The third part consists of functions that manipulate monads (see section 4.1). The fourth and fifth parts consist of simpler (but still very useful) functions, like flattening and summing. The sixth part consists of functions that manipulate constructors and constructor names. The following sections describe each of these parts in more detail.

pmap, fmap, cata ana, hylo, para crush, fcrush (a) Recursion op's	pzip, fzip punzip, funzip pzipWith, pzipWith' pequal, fequal (b) Zips etc.	pmapM, fmapM, cataM anaM, hyloM, paraM propagate, cross thread, fthread (c) Monad op's	flatten, fflatten fl_par, fl_rec, conc (d) Flatten functions
	psum, size, prod pand, pall por, pany, pelem (e) Miscellaneous	constructorName, fconstructorName constructors, fconstructors constructor2Int, fconstructor2Int int2constructor, int2fconstructor (f) Constructor functions	

**Fig. 5.** Overview of PolyLib

## Recursion Operators

---

```
pmap :: (a -> b) -> d a -> d b
fmap :: (a -> c) -> (b -> d) -> f a b -> f c d
```

Function `pmap` takes a function `f` and a value `x` of datatype `d a`, and applies `f` recursively to all occurrences of elements of type `a` in `x`. With `d` as a functor acting on types, `pmapd` is the corresponding functor action on functions. Function `fmapf` is the corresponding functor action for a pattern functor `f`. **Also known as:** `map` [31], `mapn` [29]. In `charity` [13] `mapd f x` is written `d{f}(x)`. **Known uses:** Everywhere! Function `fmap` is used in the definition of `pmap`, `cata`, `ana`, `hylo`, `para` and in many other PolyLib functions. **Background:** The `map` function was one of the first combinators distinguished in the work of Bird and Meertens, [12, 35]. The traditional `map` in functional languages maps a function over a list of elements. The current Haskell version of `map` is overloaded:

```
map :: Functor f => (a->b) -> f a -> f b
```

and can be used as the polytypic `pmap` if instance declarations for all regular type constructors are given. Function `pmap` can be used to give default instances for the Haskell `map`.

---

---

```
cata :: (FunctorOf d a b -> b) -> (d a -> b)
ana  :: (FunctorOf d a b <- b) -> (d a <- b)
hylo :: (f a b -> b) -> (c -> f a c) -> (c -> b)
para :: (d a -> FunctorOf d a b -> b) -> (d a -> b)
```

Four powerful recursion operators on the type `d a`: The catamorphism, `cata`, “evaluates” a data structure by recursively replacing the constructors with functions. The typing of `cata` may seem unfamiliar but with the explanation of `FunctorOf` above it can be seen as equivalent to:

```
cata :: (f a b -> b) -> (τ(f) a -> b)
```

The anamorphism, `ana`, works in the opposite direction and builds a data structure. The hylomorphism, `hylo`, is the generalisation of these two functions that simultaneously builds and evaluates a structure. Finally, the paramorphism, `para`, is a generalised form of `cata` that gives its parameter function access not only to the results of evaluating the substructures, but also the structure itself. **Also known as:**

PolyLib	Functorial ML [9]	Squiggol	<code>charity</code> [13]
<code>cata i</code>	<code>fold<sub>1</sub> i</code>	$\langle i \rangle$	$\{   i   \}$
<code>ana o</code>	-	$\langle o \rangle$	$(   o   )$

Functions `cata` and `para` are instances of the Visitor pattern in [21]. **Known uses:** Very many polytypic functions are defined using `cata`: `pmap`, `crush`, `thread`, `flatten`,

propagate, and all our applications use it. Function `para` is used in `rewrite`. **Background:** The catamorphism, `cata`, is the generalisation of the Haskell function `foldr` and the anamorphism, `ana`, is the (category theoretic) dual. Catamorphisms were introduced by Malcolm [33,34]. A hylomorphism is the fused composition of a catamorphism and an anamorphism specified by: `hylo i o = cata i . ana o`. The paramorphism [36], `para`, is the elimination construct for the type `d a` from Martin-Löf type theory. It captures the recursion pattern of primitive recursive functions on the datatype `d a`.

---

---

```
crush  :: (a->a->a) -> a -> d a -> a
fcrush :: (a->a->a) -> a -> f a a -> a
```

The function `crush op e` takes a structure `x` and inserts the operator `op` from left to right between every pair of values of type `a` at every level in `x`. (The value `e` is used in empty leaves.) **Known uses:** within the library see section 3.6. Many of the functions in that section are then used in the different applications. **Background:** The definition of `crush` is found in [37]. For an associative operator `op` with unit `e`, `crush op e` can be defined as `foldr op e . flatten`. As `crush` has the same arguments as `fold` on lists it can be seen as an alternative to `cata` as the generalisation of `fold` to regular datatypes.

---

## Zips

---

```
pzip  :: (d a,d b) -> Maybe ( d (a,b) )
punzip :: d (a,b) -> (d a,d b)
fzip  :: (f a b,f c d) -> Maybe ( f (a,c) (b,d) )
funzip :: f (a,c) (b,d) -> (f a b,f c d)
```

Function `punzip` takes a structure containing pairs and splits it up into a pair of structures containing the first and the second components respectively. Function `pzip` is a partial inverse of `punzip`: it takes a pair of structures and zips them together to `Just` a structure of pairs if the two structures have the same shape, and to `Nothing` otherwise. **Also known as:** `zipm` [29], `zip.×.d` [23], **Known uses:** Function `fzip` is used in the definition of `pzipWith`. **Background:** The traditional function `zip`

```
zip  :: [a] -> [b] -> [(a,b)]
```

combines two lists and does not need the `Maybe` type in the result as the longer list can always be truncated. (In general such truncation is possible for all types that have a nullary constructor, but not for all regular types.) A more general (“doubly polytypic”) variant of `pzip`: `transpose` (called `zip.d.e` in [23])

```
transpose :: d (e a) -> e (d a)
```

was first described by Fritz Ruehr [43]. For a formal and relational definition, see Hoogendijk & Backhouse [23].

---

---

```
pzipWith :: ((a,b) -> Maybe c) -> (d a,d b) -> Maybe (d c)
pzipWith' :: (FunctorOf d c e -> e) -> ((d a,d b) -> e) ->
           ((a,b) -> c) -> (d a,d b) -> e
```

Function `pzipWith op` works like `pzip` but uses the operator `op` to combine the values from the two structures instead of just pairing them. As the zip might fail, we also give the operator a chance to signal failure by giving it a `Maybe`-type as a result.<sup>6</sup>

Function `pzipWith'` is a generalisation of `pzipWith` that can handle two structures of different shape. In the call `pzipWith' ins fail op`, `op` is used as long as the structures have the same shape, `fail` is used to handle the case when the two structures mismatch, and `ins` combines the results from the substructures. (The type of `ins` is the same as the type of the first argument to `cata`.) **Also known as:** `zipopm` [29]. **Known uses:** Function `pzipWith'` is used in the definition of equality, matching and even unification. **Background:** Function `pzipWith` is the polytypic variant of the Haskell function `zipWith`

```
zipWith :: (a->b->c) -> [a] -> [b] -> [(a,b)]
```

but `pzipWith'` is new. Function `pzip` is just `pzipWith Just`.

---

---

```
pequal :: (a->b->Bool) -> d a -> d b -> Bool
fequal :: (a->b->Bool) -> (c->d->Bool) -> f a c -> f b d -> Bool
```

The expression `pequal eq x y` checks whether or not the structures `x` and `y` are equivalent using the equivalence operator `eq` to compare the elements pairwise. **Known uses:** `fequal` is used in the unification algorithm to determine when two terms are top level equal. **Background:** An early version of a polytypic equality function appeared in [44]. Function `pequal` can be instantiated to give a default for the Haskell `Eq`-class for regular datatypes:

```
(==) :: Eq a => d a -> d a -> Bool
(==) = pequal (==)
```

In Haskell the equality function can be automatically derived by the compiler, and our polytypic equality is an attempt at moving that derivation out of the compiler into the prelude.

---

<sup>6</sup> The type constructor `Maybe` can be replaced by any monad with a zero, but we didn't want to clutter up the already complicated type with contexts.

## Monad Operations

---

```
pmapM  :: Monad m => (a -> m b) -> d a -> m (d b)
pmapMr :: Monad m => (a -> m b) -> d a -> m (d b)
fmapM  :: Monad m => (a->m c) -> (b->m d) -> f a b -> m (f c d)
cataM  :: Monad m => (FunctorOf d a b->m b) -> (d a -> m b)
anaM   :: Monad m => (b->m (FunctorOf d a b)) -> (b -> m (d a))
hyloM  :: Monad m => (f a b->m b) -> (c->m (f a c)) -> c -> m b
paraM  :: Monad m => (d a->FunctorOf d a b->m b) -> d a -> m b
```

Function `pmapM` is a variant of `pmap` that threads a monad `m` from left to right through a structure after applying its function argument to all elements in the structure. Function `pmapMr` is the same but for threading a monad `m` from right to left through a structure. For symmetry's sake, the library also contains a function `pmapM1`, which is equal to `pmapM`. Furthermore, the library also contains the left and right variants of functions like `cataM` etc. A monadic map can, for example, use a state monad to record information about the elements in the structure during the traversal. The other recursion operators are generalised in the same way to form even more general combinators. **Also known as:** traversals [29]. **Known uses:** in `unify` and in the parser. **Background:** Monadic maps and catormorphisms are described in [20]. Monadic anamorphisms and hylomorphisms are defined in [39]. The monadic map (also called active traversal) is closely related to `thread` (also called passive traversal):

```
pmapM f = thread . pmap f
thread = pmapM id
```

---

---

```
propagate :: d (Maybe a) -> Maybe (d a)
cross     :: d [a] -> [d a]
```

Function `propagate` propagates `Nothing` to the top level. Function `cross` is the cross (or tensor) product that given a structure `x` containing lists, generates a list of structures of the same shape. This list has one element for every combination of values drawn from the lists in `x`. These two functions can be generalised to `thread` any monad through a value. **Known uses:** `propagate` is used in the definition of `pzip`. **Background:** Function `propagate` is an instance of `transpose` [43], and both `propagate` and `cross` are instances of `thread` below.

---

---

```
thread :: Monad m => d (m a) -> m (d a)
fthread :: Monad m => f (m a) (m b) -> m (f a b)
```

Function `thread` is used to tie together the monad computations in the elements from left to right. **Also known as:** `distd` [20]. **Known uses:** Function `thread` can be used to define the monadic map: `pmapM f = thread . pmap f`. Function `fthread` is also used in the parser to thread the parsing monad through different structures. Function `thread` can be instantiated (with `d = []`) to the Haskell prelude function

```
accumulate :: Monad m => [m a] -> m [a]
```

but also orthogonally (with `m = Maybe`) to `propagate` and (with `m = []`) to `cross`.

---

## Flatten Functions

---

```
flatten :: d a -> [a]
fflatten :: f a [a] -> [a]
fl_par   :: f a b -> [a]
fl_rec   :: f a b -> [b]
```

Function `flatten x` traverses the structure `x` and collects all elements from left to right in a list. The other three function are variants of this for a pattern functor `f`. **Also known as:** `extractm,i` [29], `listify` [23]. **Known uses:** `fl_rec` is used in the unification algorithm to find the list of immediate subterms of a term. Function `fflatten` is used to define `flatten`

```
flatten = cata fflatten
```

**Background:** In the relational theory of polytypism [23] there is a membership relation `mem.d` for every relator (type constructor) `d`. Function `flatten` can be seen as a functional implementation of this relation:

```
a mem.d x ≡ a 'elem' (flattend x)
```

---

**Miscellaneous** A number of simple polytypic functions can be defined in terms of `crush` and `pmap`. For brevity we present this part of PolyLib below by providing only the name, the type and the definition of each function.

---

```

psum :: d Int -> Int           psum = crush (+) 0
prod  :: d Int -> Int           prod  = crush (*) 1
conc  :: d [a] -> [a]          conc  = crush (++) []
pand  :: d Bool -> Bool        pand  = crush (&&) True
por   :: d Bool -> Bool        por   = crush (||) False

size   :: d a -> Int           size   = psum . pmap (\_->1)
flatten :: d a -> [a]         flatten = conc . pmap (:[])
pall   :: (a->Bool) -> d a -> Bool pall p = pand . pmap p
pany   :: (a->Bool) -> d a -> Bool pany p = por . pmap p
pelem  :: Eq a => a -> d a -> Bool pelem x = pany (\y->x==y)

```

---

## Constructors

---

```

constructorName :: d a -> String
fconstructorName :: f a b -> String
constructors     :: [d a]
fconstructors    :: [f a b]
constructor2Int  :: d a -> Int
fconstructor2Int :: f a b -> Int
int2constructor  :: Int -> d a
int2fconstructor :: Int -> f a b

```

Function `constructorName` takes a value of type `d a` and returns its outermost constructor name. Function `constructors` returns a list with all the constructors of a datatype `d a`. For example, for the datatype `Bin` it returns `[Tip undefined, Join undefined undefined]`. The functions `constructor2Int` and `int2constructor` take constructors to integers and vice versa. **Known uses:** `constructorName` is used in `pshow`, the polytypic version of the derived `show` function in Haskell, `constructors` is used in showing, parsing and compressing values, and both `int2constructor` and `constructor2Int` in compressing values.

---



## 4 Generic Unification

This chapter presents a substantial application of the techniques that have been developed thus far. The topic is a generic unification algorithm.

Briefly, unification is the process of making two terms (such as arithmetic expressions or type expressions) equal by suitably instantiating the variables in the terms. It is very widely used in, for example, pattern matching, type checking and theorem proving. For those who haven't already encountered it, let us first give an informal explanation before giving a summary of the development of the generic algorithm.

We explain the process in terms of a specific case before considering the generic version. Consider the datatype definition

```
data Expr = var V
          | number Nat
          | plus Expr Expr
          | times Expr Expr
```

This can be read as the datatype of abstract syntax trees for a context-free grammar

$$E ::= V \mid N \mid (E + E) \mid (E * E)$$

for *terms* like “ $((1+x)*3)$ ” when  $V$  produces variables and  $N$  produces numbers.

Another view is that a *term* of the datatype *Expr* is a tree with the constructors **var**, **number**, **plus** and **times** at the nodes, and numbers and variables at the leaves. In this view, the constructors are *uninterpreted*, which means that trees corresponding to equal but non-identical arithmetic expressions are considered different. For example, the trees corresponding to  $((1+x)*3)$  and  $(3+(x*3))$  are different. It is this view of terms as tree structures that is used in unification. Nevertheless, for ease of writing we shall use the concrete syntax of arithmetic expressions to write terms.

Now consider two terms, say  $((1+x)*3)$  and  $((y+z)*3)$ . “Unifying” these terms means substituting terms for the variables  $x$ ,  $y$  and  $z$  so that the terms become identical. One possibility, in this case, is to substitute  $z$  for  $x$  and  $1$  for  $y$ . After this substitution both terms become equal to  $((1+z)*3)$ . There are many other possibilities. For example, we could substitute  $1$  for all of  $x$ ,  $y$  and  $z$ , thus unifying the two terms in the term  $((1+1)*3)$ . This latter substitution is however less general than the former. Unification involves seeking a “most general” unifier for two given terms. Of course, some pairs of terms are not unifiable: a trivial example is the pair of terms  $0$  and  $1$ . These are not unifiable because they contain no variables. The pair of terms  $x$  and  $(1+x)$  is also not unifiable, but for a different reason: namely, the first term will always have fewer constructors than the second whatever substitution we make for  $x$ .

We have described unification for arithmetic expressions but unification is also used for other term algebras. A major application is in polymorphic type inference, as in most modern functional languages. In this application it is type expressions that are unified.

Suppose that a program contains the function application  $\mathbf{f} \ \mathbf{x}$ , and at that stage the term representing the type inferred for  $\mathbf{f}$  is  $(p \rightarrow q)$ , and for  $\mathbf{x}$  it is  $r$ . Then first  $p$  and  $r$  are unified. If that fails, there is a type error. Otherwise, let  $(p' \rightarrow q')$  be the result of applying the most general unifier to  $(p \rightarrow q)$ . That is the new type inferred for  $\mathbf{f}$ , while we get  $p'$  for  $\mathbf{x}$ , and  $q'$  for the application  $\mathbf{f} \ \mathbf{x}$ .

In a *generic* unification algorithm we make the term structure a parameter of the algorithm. So, one instance of the algorithm unifies arithmetic expressions, another type expressions. In order to formalise this we use  $F$  to denote a functor (the pattern functor of the constant terms we want to unify) and show how to extend  $F$  to a functor  $F^*$  such that  $F^*V$ , for type  $V$  of variables, is the set of all terms. We also define substitution of variables, and most general unifiers.

The functor  $F^*$  is (the functor part of) a *monad*. In the last ten years, monads have been recognised to be an important concept in many applications of functional programming. We therefore begin in section 4.1 by introducing the concept at first without reference to unification. There is much that can be said about monads but our discussion is brief and restricted to just what we need to present the unification algorithm. The monad  $F^*$  defined by an arbitrary functor  $F$  is then discussed along with the definition of a substitution.

The discussion of the unification algorithm proper begins in section 4.2. Here the discussion is also brief since we assume that the non-generic algorithm is known from the literature. In order to compare the calculational method of proof with traditional proofs, chapter 5 presents a generic proof of one aspect of the algorithm's correctness, namely that a non-trivial expression is not unifiable with any variable that occurs properly in it.

## 4.1 Monads and Terms

*Monads and Kleisli composition* A monad is a concept introduced in category theory that has since proved its worth in functional programming. It is a general concept which we introduce via a particular instance, the *Maybe* monad.

Suppose we have two functions

$$\begin{aligned} f &:: a \rightarrow \text{Maybe } b \\ g &:: b \rightarrow \text{Maybe } c \end{aligned}$$

Think of these total functions as modelling partial functions:  $f$  computes a  $b$ -value from an  $a$ -value, or fails, and likewise,  $g$  computes a  $c$ -value from a  $b$ -value, or fails. Can we combine these functions into a single function

$$g \diamond f :: a \rightarrow \text{Maybe } c$$

that combines the computations of  $f$  and  $g$  when both succeed, and fails when either of them fails? The types don't fit for normal composition, but here is how to do it:

$$\begin{aligned} (g \diamond f) \ x &= h (f \ x) \ \mathbf{where} \\ h \ \mathbf{none} &= \mathbf{none} \\ h \ (\mathbf{one} \ y) &= g \ y \end{aligned}$$

This form of composition is called *Kleisli composition*. Kleisli composition shares some pleasant properties with normal composition. First, the operation is associative:

$$f \diamond (g \diamond h) = (f \diamond g) \diamond h$$

for  $f$ ,  $g$  and  $h$  such that the expressions involved are well-typed. We may therefore drop the parentheses in chains of Kleisli compositions and write  $f \diamond g \diamond h$ . Moreover,  $\diamond$  has neutral element **one**, which we call the *Kleisli identity*:

$$\mathbf{one} \diamond f = f = f \diamond \mathbf{one} \quad .$$

Kleisli composition gives a convenient way to fit functions together that would not fit together with normal composition. Kleisli composition is not just possible for *Maybe*, but for many other functors as well. A functor with a Kleisli composition and Kleisli identity—that satisfy a number of laws to be discussed shortly—is called a *monad*. A trivial example is the functor **Id**: take normal function composition as its Kleisli composition. A less trivial example is the functor *Set*. For this functor, Kleisli composition takes the form

$$(f \diamond g)x = \{z \mid \exists(y:: y \in gx \wedge z \in fy)\} \quad .$$

Its Kleisli identity is the singleton former  $\{-\}$ . We shall encounter more monads later.

Formally, the triple  $(M, \diamond, \eta)$  is a monad, where  $M$  is a functor,  $\diamond$  and  $\eta$  are its Kleisli composition and Kleisli identity, if the following properties hold. First,  $\diamond$  is a function of polymorphic type

$$(b \rightarrow Mc) \times (a \rightarrow Mb) \rightarrow (a \rightarrow Mc)$$

and  $\eta$  is a function of polymorphic type

$$a \rightarrow Ma \quad .$$

Second,  $\diamond$  is associative with  $\eta$  as neutral element. Finally, the following rules are satisfied:

$$\begin{aligned} Mf \bullet (g \diamond h) &= (Mf \bullet g) \diamond h \\ (f \diamond g) \bullet h &= f \diamond (g \bullet h) \\ (f \bullet g) \diamond h &= f \diamond (Mg \bullet h) \end{aligned}$$

In fact, these equalities are automatically satisfied in all the monads that we consider here. They are consequences of the so-called *free theorem* for  $\diamond$ . Their validity depends on a property called (*polymorphic*) *parametricity* that is satisfied by Haskell restricted to total functions which we discuss in section 5.2.

**Exercise 4.1** Let  $(M, \diamond, \eta)$  be a monad. Express  $Mf$  in terms of Kleisli composition and identity. Define

$$\mathbf{mul} = \mathbf{id} \diamond \mathbf{id} :: MMa \rightarrow Ma$$

(The function **mul** is called the *multiplier* of the monad.) What is the function **mul** for the case  $M = \mathit{Set}$ ?

Prove that  $f \diamond g = \text{mul} \bullet Mf \bullet g$ . Also prove the following three equalities:

$$\text{mul} \bullet \text{mul} = \text{mul} \bullet M\text{mul}$$

$$\text{mul} \bullet \eta = \text{id} = \text{mul} \bullet M\eta .$$

□

*Terms with variables* Recall the datatype *Expr* introduced at the beginning of this section. We can regard it as a datatype for terms involving numbers, addition and multiplication to which has been added an extra alternative for variables.

Let  $F$  be the pattern functor corresponding to the definition of *Expr* without variables. Then  $\text{Expr} = \mu G$ , where  $G a = V + Fa$ . This can be done generically. Consider, for unary functor  $F$ , the unary functor  $V^{\kappa+F}$ . This, we recall, is defined by

$$(V^{\kappa+F})a = V + Fa$$

where  $a$  ranges over types, and

$$(V^{\kappa+F})f = \text{id}_V + Ff$$

where  $f$  ranges over functions. For fixed  $F$ , the mapping  $a \mapsto \mu(a^{\kappa+F})$  is a functor, namely the type functor  $\tau(\odot)$  of the bifunctor  $a \odot b = a + Fb$ . Denote this type functor by  $F^*$  (so  $F^*V = \mu(V^{\kappa+F})$ )<sup>7</sup>. Its action on functions is as follows. For  $f :: a \rightarrow b$ :

$$F^*f = ([a^{\kappa+F}; \text{in}_{b^{\kappa+F}} \bullet f + \text{id}])$$

Note that we have specified the pattern functor “ $a^{\kappa+F}$ ” inside the catamorphism brackets here since there is a possibility of confusion between different algebras. Note also that

$$(V^{\kappa+F})F^*V = V + FF^*V$$

so that

$$\text{in}_{V^{\kappa+F}} :: V + FF^*V \rightarrow F^*V .$$

Given a datatype  $\mu F$ , we can then extend it with variables by switching to  $F^*V$ . We define two embeddings by:

$$\begin{aligned} \text{embl}_V &:: V \rightarrow F^*V & \text{embr}_V &:: FF^*V \rightarrow F^*V \\ \text{embl}_V &= \text{in}_{V^{\kappa+F}} \bullet \text{inl} & \text{embr}_V &= \text{in}_{V^{\kappa+F}} \bullet \text{inr} \end{aligned}$$

The functor  $F^*$  forms the *substitution monad*  $(F^*, \diamond, \eta)$  with, for functions  $f :: a \rightarrow F^*b$  and  $g :: b \rightarrow F^*c$ ,

$$g \diamond f = ([b^{\kappa+F}; g \nabla \text{embr}_c]) \bullet f$$

<sup>7</sup> The star notation is used here to suggest a link with the Kleene star, denoting iteration in a regular algebra.  $F^*$  can be seen as iterating functor  $F$  an arbitrary number of times. More significantly, the notation highlights a formal link between monads and closure operators. See, for example, [3] for more details.

$$\eta = \text{embl} \ .$$

Note that the catamorphism in the definition of  $g \circ f$  has pattern functor  $b^k + F$ , as indicated by the parameter before the semicolon. We omit explicit mention of this information later, but it is vital to the correct use of the computation and other laws. In addition we omit type information on the initial algebra, although again it is vitally important.

**Exercise 4.2** Take  $F$  to be  $a \times$  for some type  $a$ . What is the type  $(a \times)^* 1$ ? What is the multiplier, what is Kleisli identity and what is Kleisli composition? (Hint: use exercise 4.1 for the last part of this exercise.)

□

**Exercise 4.3** Consider the case  $F = (1+)$ . Show that  $(1+)^* V \cong \mathbb{N} \times (V+1)$ . Specifically, construct an initial algebra.

$$\text{in} \quad :: \quad V + (1 + (\mathbb{N} \times (V+1))) \rightarrow \mathbb{N} \times (V+1)$$

and express catamorphisms on elements of type  $\mathbb{N} \times (V+1)$  in terms of catamorphisms on  $\mathbb{N}$ .

□

**Exercise 4.4** Verify that Kleisli composition as defined above is indeed associative and that  $\text{embl}$  is its neutral element.

□

*Assignments and Substitutions* An *assignment* is a mapping of variables to terms, for example  $\{x := (y+x), y := 0\}$ . An assignment can be *performed* on a term. This means a *simultaneous and systematic* replacement of the variables in the term by the terms to which they are mapped. For example, performed on the term  $(x+y)$  our example assignment gives  $((y+x)+0)$ . We model assignments as functions with the typing  $V \rightarrow F^*V$ . Because we want functions to be total, this means we also have to define the assignment for all variables in  $V$ . If  $V = \{x, y, z\}$ , we can make the above assignment total by writing it as  $\{x := (y+x), y := 0, z := z\}$ . Note that to the left of “:=” in an assignment we have an element of  $V$ , and to the right an element of  $F^*V$ . So to be precise, if assignment  $f$  has “ $z := z$ ”, this means that  $f z = \eta z$ . In particular, the (empty) *identity assignment* is  $\eta$ .

Given an assignment  $f :: V \rightarrow F^*V$ , we want to define the *substitution*  $\text{subst } f$  as a function performing  $f$  on a term. The result is again a term. The term consisting of the single variable  $x$  is  $\eta x$ . Applying  $\text{subst } f$  to it, the result should be  $f x$ . So

$$\begin{aligned} & (\text{subst } f) \bullet \eta \\ = & \quad \left\{ \begin{array}{l} \text{desired result} \\ f \end{array} \right\} \end{aligned}$$

$$\begin{aligned}
&= \{ \text{Kleisli identity} \} \\
&\quad f \diamond \eta \\
&= \{ \text{monad equality} \} \\
&\quad (f \diamond \text{id}) \bullet \eta .
\end{aligned}$$

Since  $\text{subst } f$  is clearly a catamorphism that distributes through constructors — for example,  $(\text{subst } f) (x+y) = ((\text{subst } f x) + (\text{subst } f y))$  — it is fully determined by its action on variables. We have found:

$$\begin{aligned}
\text{subst} &:: (V \rightarrow F^*V) \rightarrow (F^*V \rightarrow F^*V) \\
\text{subst } f &= f \diamond \text{id}
\end{aligned}$$

Two substitutions can always be merged into a single one:

$$\begin{aligned}
&(\text{subst } f) \bullet (\text{subst } g) \\
&= \{ \text{definition of subst} \} \\
&\quad (f \diamond \text{id}) \bullet (g \diamond \text{id}) \\
&= \{ \text{monad equalities} \} \\
&\quad f \diamond (\text{id} \bullet (g \diamond \text{id})) \\
&= \{ \text{id is identity of } \bullet \} \\
&\quad f \diamond (g \diamond \text{id}) \\
&= \{ \diamond \text{ is associative} \} \\
&\quad (f \diamond g) \diamond \text{id} \\
&= \{ \text{definition of subst} \} \\
&\quad \text{subst } (f \diamond g) .
\end{aligned}$$

## 4.2 Generic Unification

*Unifiers* Two terms  $x$  and  $y$  containing variables can be *unified* if there is some assignment  $f$  such that performing  $f$  on  $x$  gives the same result as performing  $f$  on  $y$ . For example, the two terms

$$(\text{u} + ((1 * \text{v}) * 2)) \quad \text{and} \quad ((\text{w} * \text{v}) + (\text{u} * 2))$$

can be unified by the assignment

$$\{\text{u} := (1 * (\text{z} + 3)), \text{v} := (\text{z} + 3), \text{w} := 1\}$$

into the unification

$$((1 * (\text{z} + 3)) + ((1 * (\text{z} + 3)) * 2))$$

Such a unifying assignment is called a *unifier* of the terms. Unifiers are not unique. Another unifier of the same two terms of the example is

$$\{\text{u} := (1 * \text{z}), \text{v} := \text{z}, \text{w} := 1\}$$

which results in the unification

$$((1*z)+((1*z)*2))$$

This last unification is more general. If  $f$  is a unifier, then, for any assignment  $h$ , the combined substitution  $h \diamond f$  is also a unifier, since

$$\begin{aligned} & h \diamond f \text{ is a unifier of } (x,y) \\ \equiv & \quad \{ \text{definition of unifier} \} \\ & \text{subst } (h \diamond f) x = \text{subst } (h \diamond f) y \\ \equiv & \quad \{ \text{combined substitutions} \} \\ & (\text{subst } h) (\text{subst } f x) = (\text{subst } h) (\text{subst } f y) \\ \Leftarrow & \quad \{ \text{cancel } (\text{subst } h) \} \\ & \text{subst } f x = \text{subst } f y \\ \equiv & \quad \{ \text{definition of unifier} \} \\ & f \text{ is a unifier of } (x,y) . \end{aligned}$$

In the example, the first, less general unifier, can be formed from the more general one by taking  $h = \{v := (z+3)\}$ . This notion of generality gives a pre-ordering on unifiers (and actually on all assignments): define

$$f \sqsubseteq g \equiv \exists(h \text{ such that } f = h \diamond g)$$

The relation  $\sqsubseteq$  is obviously transitive and reflexive, but in general not anti-symmetric. If two unifiers are equally general:  $f \sqsubseteq g \wedge g \sqsubseteq f$ , then  $f$  and  $g$  can be different. But they are to all intents and purposes equivalent: they differ at most in the choice of names for the variables in the result.

If two terms are unifiable at all, then among all unifiers there is a *most general unifier*. That term is commonly abbreviated to *mgu*. Clearly, any two mgu's are equivalent. In the example, the second unifier is an mgu.

*A generic shell for unification* We develop the unification algorithm in two stages. In this stage we give a generic “shell” in terms of type classes. In the second stage, we show how to make any regular functor into an instance of the classes involved.

Terms may have children, they may happen to be variables, and we should be able to see if superficially —at the top level of the term trees— the constructors are equal. As before, we assume a fixed type  $V$  for variables. Here are the corresponding class declarations:

```

class Children t where children      :: t → List t
                        mapChildren  :: (t → t) → (t → t)
class VarCheck t where varcheck     :: t → Maybe V
class TopEq t where topEq          :: t × t → Bool

class (Children t, VarCheck t, TopEq t) ⇒ Term t

```

We give a concrete instantiation as an example — illustrating some fine points at the same time. Let  $C$  be some type for representing constructors. Here is the datatype we will use to instantiate the classes:

```
data  $T = \text{Var } V \mid \text{Con } C (\text{List } T)$ 
```

First we make  $T$  into an instance of *Children*:

```
instance Children  $T$  where
  children (Var  $v$ )    = nil
  children (Con  $c$   $ts$ ) =  $ts$ 
  mapChildren  $f$  (Var  $v$ )    = Var  $v$ 
  mapChildren  $f$  (Con  $c$   $ts$ ) = Con  $c$  (List  $f$   $ts$ )
```

Note here that *mapChildren*  $f$  only maps function  $f$  over the *immediate* children of its argument. No recursion is involved.

Here is how  $T$  fits in the *VarCheck* class:

```
instance VarCheck  $T$  where
  varcheck (Var  $v$ )    = one  $v$ 
  varcheck (Con  $c$   $ts$ ) = none
```

For *TopEq* we assume that *eq* is an equality test on  $C$  and on  $V$ :

```
instance TopEq  $T$  where
  topEq ((Var  $v_0$ )    , (Var  $v_1$ )    ) = eq  $v_0$   $v_1$ 
  topEq ((Con  $c_0$   $ts_0$ ), (Con  $c_1$   $ts_1$ )) = (eq  $c_0$   $c_1$   $\wedge$  length  $ts_0$  = length  $ts_1$ )
  topEq ( -            , -            ) = false
```

Note that for this test the children of the terms are irrelevant. This is why we give it the name *topEq*.

Having made  $T$  an instance of the three superclasses of *Term*, we can now proudly announce:

```
instance Term  $T$ 
```

So much for this concrete instantiation. We continue with the generic problem. Here is a function to collect all subterms of a term in the *Term* class (or actually the *Children* class):

```
subTerms :: Children  $t \Rightarrow t \rightarrow \text{List } t$ 
subTerms  $x = \text{cons } x (\text{concat } (\text{List } \textit{subTerms} (\textit{children } x)))$ 
```

and here is a function that uses a list comprehension to collect all variables occurring in a term:

```
vars :: Term  $t \Rightarrow t \rightarrow \text{List } V$ 
vars  $x = [v \mid \text{one } v \leftarrow \text{List } \textit{varCheck} (\textit{subTerms } x)]$ 
```



Earlier we saw a treatment of assignments as functions. Here we introduce a class for assignments, so that it is also possible to make other concrete representations into instances. The parameter  $t$  stands for terms.

```
class Assig t where idAssig  :: V → t
                    modBind  :: V × t → ((V → t) → (V → t))
                    lookupIn :: (V → t) × V → Maybe t
```

The type  $F^*V$  can be made into a generic instance by:

```
instance Assig (F* V) where
  idAssig      = embl
  modBind (v, x) = (f ↦ (v' ↦ if eq v' v then x else f v'))
  lookupIn (f, v) = if eq (f v) (idAssig v) then none else one (f v)
```

in which we see both the Kleisli identity `embl` of the substitution monad, and one of the *Maybe* monad. The result `none` signifies that  $v$  is mapped to itself (embedded in the term world).

We have chosen a particular implementation for assignments: assignments are functions. If Haskell would allow multiple parameter type classes we could abstract from the particular implementation, and replace the occurrences of  $V \rightarrow t$  in the types of the functions of the class *Assig* by a type variable  $a$ . Thus we could obtain a more concrete instance of *Assig* by taking list of pairs  $(v, x)$ , with  $v$  a variable and  $x$  a term, instead of functions. Then *idAssig* is the empty list, *modBind* can simply cons the pair onto the list, and *lookupIn* looks for the first pair with the given variable and returns the corresponding term. If the given variable is not found, it fails. An efficient implementation of *Assig* would use balanced trees, or even better hash tables. With the class mechanism the implementation can be *encapsulated*, that is, hidden to the rest of the program, so that the program can first be developed and tested with a simple implementation. It can later be replaced by a more efficient sophisticated implementation without affecting the rest of the program. It should be clear that this is an important advantage.

*The unification algorithm proper* We give the algorithm — which is basically the algorithm found in the literature — without much explanation. As to notation, we use the monad  $(Maybe, \diamond, \eta)$ .

```
unify  :: (Term t, Assig t) ⇒ t × t → Maybe (V → t)
unify' :: (Term t, Assig t) ⇒ t × t → ((V → t) → Maybe (V → t))
```

The definition of *unify* is now simply to start up *unify'* with the empty assignment. The function *unify'* is defined as a higher order function, threading “assignment transformations” together with  $\diamond$ .

```

unify (x, y) = unify' (x, y) idAssig
unify' (x, y) = uni (varCheck x, varCheck y) where
  uni (none , none ) | topEq (x, y) = uniTerms (x, y)
                    | otherwise    = const none
  uni (one u , one v) | eq u v      =  $\eta$ 
  uni (one u , _   )              =  $u \mapsto y$ 
  uni (_     , one v)              =  $v \mapsto x$ 

  uniTerms (x, y) = threadList(List unify' (zip (children x) (children y)))

```

All the right-hand sides here are functions that return maybe an assignment, given an assignment. The function *threadList* is simply the list catamorphism with Kleisli composition:

```

threadList :: Monad m => List (a -> m a) -> (a -> m a)
threadList = foldr ( $\diamond$ )  $\eta$ 

```

The auxiliary operator ( $\mapsto$ ) should “*modBind*” its arguments into the unifier being collected, but there are two things to be taken care of. No binding may be introduced that would mean an infinite assignment. This is commonly called the *occurs check*. And if the variable is already bound to a term, that term must be unified with the new term, and the unifier obtained must be threaded into the assignment being collected.

```

( $\mapsto$ ) :: (Term t, Assig t) => V x t -> ((V -> t) -> Maybe (V -> t))
(v  $\mapsto$  x) s = if occursCheck (v, s, x)
              then none
              else case lookupIn (s, v) of
                none  -> ( $\eta \bullet \text{modBind}$  (v, x)) s
                one y -> (( $\eta \bullet \text{modBind}$  (v, x))  $\diamond$  unify' (x, y)) s

```

The following is a hack to implement the occurs check. This is basically a reachability problem in a graph — is there a cycle from *v* to itself?, or rather: are we about to create a cycle? We must take account both of the unifier collected already, and the new term. Because we know no cycles were created yet, the graph is more like a tree, so any search strategy terminates. The approach here is not optimally efficient, but in practice quite good with lazy evaluation (and horrible with eager evaluation). There exist linear-time solutions, but they require much more bookkeeping.

```

occursCheck :: (Term t, Assig t) => V x (V -> t) x t -> Bool
occursCheck (v, s, x) = v  $\in$  reachlist (vars x) where
  reachlist vs = vs ++ concat (List reachable vs)
  reachable v = reachlist (mayvars (lookupIn (s, v)))
  mayvars none = []
  mayvars (one y) = vars y

```

Here, *reachlist* collects the variables reachable from a *list* of variables, while *reachable* collects the variables reachable from a *single* variable.

*The generic Term instance* All we have to do now is make  $F^*V$  an instance of the *Term* class. That is surprisingly easy. For the *Children* class:

```
instance Children (F*V) where
  children      = ((nil •!) ∇ fl_rec) • out
  mapChildren f = in • (id_V + F f) • out
```

where `fl_rec` is defined in `PolyLib`, see Section 3.6. For the *VarCheck* class:

```
instance VarCheck (F*V) where
  varcheck = (one ∇ (none •!)) • out
```

For *TopEq* we use the fact that `fequal` tests on equality of functor structures. `fequal` is defined in `PolyLib`, see Section 3.6.

```
instance TopEq (F*V) where
  topEq (t, t') = fequal (==) (x ↦ y ↦ True) (out t) (out t')
```

For a complete implementation of the generic unification program, see [26].

## 5 From Functions to Relations

In the preceding chapter we have done what we ourselves have decried: we have presented an algorithm without even a verification of its correctness, let alone a construction of the algorithm from its specification. An excuse is that a full discussion of correctness would have distracted from the main goal of that chapter, which was to show how the generic form of the —known to be correct— algorithm is implemented. That is, however, only an excuse since, so far as we know, no proof of correctness of the generic algorithm has ever been constructed. In section 5.4 we remedy this lacuna partially by presenting one lemma in such a proof of correctness. To that end, however, we need to extend the programming calculus from total functions to relations.

### 5.1 Why Relations?

In a summer school on advanced *functional* programming, it may seem odd to want to introduce relations but there are several good reasons for making it an imperative. In the first place, specifications are typically relations, not total functions. The specification of the unification algorithm is a case in point since it embodies both nondeterminism and partiality. Nondeterminism is embodied in the requirement to compute a most general unifier, not *the* most general unifier. It would be infeasible to require the latter since, in general, there is no single most general unifier of two terms. Partiality is also present in the fact that a most general unifier may not exist. Partiality can be got around in the implementation by using the *Maybe* monad as we did here, but avoiding nondeterminism in the specification is undesirable.

A second reason for introducing relations is that termination arguments are typically based on well-founded relations. Our discussion of the correctness of the unification algorithm in section 5.4 is based on the construction of a well-founded relation, although in this case termination is not the issue at stake.

A third, compelling reason for introducing relations is that the “free theorem” for polymorphic functions alluded to above and discussed in detail below is based on *relations* on functions and necessitates an extension of the concept of functor to relations. Also, the most promising work we know of that aims to be precise about what is really meant by “generic” is that due to Hoogendijk [22] which is based on a relational semantics of higher-order polymorphism .

### 5.2 Parametric Polymorphism

Space does not allow us to consider the extension to relations in full depth and so we will have to make do with a brief account of the issues involved. For more detail see [11, 1]. We believe, nevertheless, that a discussion of generic programming would be incomplete without a summary of Reynolds’ [40] *abstraction theorem* which has since been popularised under the name “theorems for free” by Wadler [45]. (This summary is taken from [23] which may be consulted for additional references.)

Reynolds' goal in establishing the abstraction theorem was to give a precise meaning to the statement that a function is "parametrically polymorphic". Suppose we have a polymorphic function  $f$  of type  $T\alpha$  for all types  $\alpha$ . That is, for each type  $A$  there is an instance  $f_A$  of type  $TA$ . The action of  $T$  is extended—in a way to be made precise shortly—to binary relations, where if relation  $R$  has type  $A \sim B$ , relation  $TR$  has type  $TA \sim TB$ . Then *parametricity* of the polymorphism of  $f$  means that for any binary relation  $R$  of type  $A \sim B$  we have  $(f_A, f_B) \in TR$ . Reynolds' abstraction theorem is the theorem that any polymorphic function expressible in the language defined in his paper is parametric. Wadler called this a "theorem for free" because, as we show shortly, the parametricity of a polymorphic function predicts algebraic properties of that function just from knowing the type of the function! Another way of viewing the theorem is as a healthiness property of functions expressible in a programming language—a programming language that guarantees that all polymorphic functions are parametric is preferable to one that cannot do so.

In order to make the notion of parametricity completely precise, we have to be able to extend each type constructor  $T$  in our chosen programming language to a function  $R \mapsto TR$  from relations to relations. Reynolds did so for function spaces and product. For product he extended the (binary) type constructor  $\times$  to relations by defining  $R \times S$  for arbitrary relations  $R$  of type  $A \sim B$  and  $S$  of type  $C \sim D$  to be the relation of type  $A \times C \sim B \times D$  satisfying

$$((u, v), (x, y)) \in R \times S \equiv (u, x) \in R \wedge (v, y) \in S .$$

For function spaces, Reynolds extended the  $\rightarrow$  operator to relations as follows. For all relations  $R$  of type  $A \sim B$  and  $S$  of type  $C \sim D$  the relation  $R \rightarrow S$  is the relation of type  $(A \rightarrow C) \sim (B \rightarrow D)$  satisfying

$$(f, g) \in R \rightarrow S \equiv \forall(x, y:: (x, y) \in R \Rightarrow (fx, gy) \in S) .$$

Note that if we equate a function  $f$  of type  $A \rightarrow B$  with the relation  $f$  of type  $B \sim A$  satisfying

$$b = fa \equiv (b, a) \in f$$

then the definition of  $f \times g$ , for functions  $f$  and  $g$ , coincides with the definition of the cartesian product of  $f$  and  $g$  given in section 2.3. Thus, not only does Reynolds' definition extend the definition of product beyond types, it also extends the definition of the product functor. Note also that the relational composition  $f \bullet g$  of two functions is the same as their functional composition. That is,  $a = f(gc) \equiv (a, c) \in f \bullet g$ . So relational composition also extends functional composition. Note finally that  $h \rightarrow k$  is a *relation* even for *functions*  $h$  and  $k$ . It is the relation defined by

$$(f, g) \in h \rightarrow k \equiv \forall(x, y:: x = hy \Rightarrow fx = k(gy)) .$$

Simplified and expressed in point-free form this becomes:

$$(f, g) \in h \rightarrow k \equiv f \bullet h = k \bullet g .$$

Writing the relation  $h \rightarrow k$  as an infix operator makes the rule easy to remember:

$$f (h \rightarrow k) g \equiv f \bullet h = k \bullet g \text{ .}$$

An example of Reynolds' parametricity property is given by function application. The type of function application is  $(\alpha \rightarrow \beta) \times \alpha \rightarrow \beta$ . The type constructor  $T$  is thus the function mapping types  $A$  and  $B$  to  $(A \rightarrow B) \times A \rightarrow B$ . The extension of  $T$  to relations maps relations  $R$  and  $S$  to the relation  $(R \rightarrow S) \times R \rightarrow S$ . Now suppose  $@$  is any parametrically polymorphic function with the same type as function application. Then Reynolds' claim is that  $@$  satisfies

$$(@_{A,C} , @_{B,D}) \in (R \rightarrow S) \times R \rightarrow S$$

for all relations  $R$  and  $S$  of types  $A \sim B$  and  $C \sim D$ , respectively. Unfolding the definitions, this is the property that, for all functions  $f$  and  $g$ , and all  $c$  and  $d$ ,

$$\forall(x, y :: (x, y) \in R \Rightarrow (fx, gy) \in S) \wedge (c, d) \in R \Rightarrow (f@c , g@d) \in S \text{ .}$$

The fact that function application itself satisfies this property is in fact the basis of Reynolds' inductive proof of the abstraction theorem (for a particular language of typed lambda expressions). But the theorem is stronger because function application is *uniquely* defined by its parametricity property. To see this, instantiate  $R$  to the singleton set  $\{(c, c)\}$  and  $S$  to the singleton set  $\{(fc , fc)\}$ . Then, assuming  $@$  satisfies the parametricity property,  $(f@c , f@c) \in S$ . That is,  $f@c = fc$ . Similarly, the identity function is the unique function  $f$  satisfying the parametricity property  $(f_A , f_B) \in R \rightarrow R$  for all types  $A$  and  $B$  and all relations  $R$  of type  $A \sim B$ —the parametricity property corresponding to the polymorphic type,  $\alpha \rightarrow \alpha$  for all  $\alpha$ , of the identity function—, and the projection function  $\text{exl}$  is the unique function  $f$  satisfying the parametricity property  $(f_{A,B} , f_{C,D}) \in R \times S \rightarrow R$  for all types  $A, B, C$  and  $D$  and all relations  $R$  and  $S$  of types  $A \sim B$  and  $C \sim D$ , respectively—the parametricity property corresponding to the polymorphic type,  $\alpha \times \beta \rightarrow \alpha$  for all  $\alpha$  and  $\beta$ , of the  $\text{exl}$  function.

The import of all this is that certain functions can be *specified* by a parametricity property. That is, certain parametricity properties have unique solutions. Most parametricity properties do not have unique solutions however. For example, both the identity function on lists and the reverse function satisfy the parametricity property of function  $f$ , for all  $R :: A \sim B$ ,

$$(f_A , f_B) \in \text{List } R \rightarrow \text{List } R \text{ .}$$

Here  $\text{List } R$  is the relation holding between two lists whenever the lists have the same length and corresponding elements of the two lists are related by  $R$ .

*Free Theorem for Monads* Let us show the abstraction theorem at work on Kleisli composition. Kleisli composition is a polymorphic function of type

$$(b \rightarrow Mc) \times (a \rightarrow Mb) \rightarrow (a \rightarrow Mc)$$

for all types  $a$ ,  $b$  and  $c$ . If it is parametrically polymorphic then it satisfies the property that, for all relations  $R$ ,  $S$  and  $T$  and all functions  $f_0$ ,  $f_1$ ,  $g_0$  and  $g_1$ , if

$$((f_0, g_0), (f_1, g_1)) \in (S \rightarrow MT) \times (R \rightarrow MS)$$

then

$$(f_0 \diamond g_0, f_1 \diamond g_1) \in R \rightarrow MT .$$

This assumes that we have shown how to extend the functor  $M$  to relations. For our purposes here, we will only need to instantiate  $R$ ,  $S$  and  $T$  to functions, and it simplifies matters greatly if we use the point-free definition of  $h \rightarrow k$  given above. Specifically, we have, for all functions  $h$ ,  $k$  and  $l$ ,

$$\begin{aligned} & ((f_0, g_0), (f_1, g_1)) \in (k \rightarrow Ml) \times (h \rightarrow Mk) \\ \equiv & \quad \{ \text{definition of } \times \} \\ & f_0 (k \rightarrow Ml) f_1 \wedge g_0 (h \rightarrow Mk) g_1 \\ \equiv & \quad \{ \text{point-free definition of } \rightarrow \text{ for functions} \} \\ & f_0 \cdot k = Ml \cdot f_1 \wedge g_0 \cdot h = Mk \cdot g_1 . \end{aligned}$$

In this way, we obtain the property that for all functions  $f_0$ ,  $f_1$ ,  $g_0$ ,  $g_1$ ,  $h$ ,  $k$  and  $l$ , if

$$(1) \quad f_0 \cdot k = Ml \cdot f_1 \wedge g_0 \cdot h = Mk \cdot g_1$$

then

$$(2) \quad (f_0 \diamond g_0) \cdot h = Ml \cdot (f_1 \diamond g_1) .$$

With its seven free variables, this is quite a complicated property. More manageable properties can be obtained by instantiating the functions in such a way that the premise becomes true. An easy way to do this is to reduce the premise to statements of the form

$$f_i = \dots \wedge g_j = \dots ,$$

where  $i$  and  $j$  are either 0 or 1, by instantiating suitable combinations of  $h$ ,  $k$  and  $l$  to the identity function. For instance, by instantiating  $h$  and  $k$  to the identity function the premise (1) reduces to

$$f_0 = Ml \cdot f_1 \wedge g_0 = g_1 .$$

Substituting the right sides for  $f_0$  and  $g_0$  in the conclusion (2) together with the identity function for  $h$  and  $k$ , we thus obtain

$$(Ml \cdot f_1) \diamond g_1 = Ml \cdot (f_1 \diamond g_1) .$$

for all functions  $l$ ,  $f_1$  and  $g_1$ . This is the first of the “free theorems” for Kleisli composition listed in section 4.1.

**Exercise 5.1** Derive the other two “free theorems” stated in section 4.1 from the above parametricity property. Investigate other properties obtained by setting combinations of  $f_0$ ,  $f_1$ ,  $g_0$ ,  $g_1$  to the identity function.

□

**Exercise 5.2** Instantiating  $M$  to the identity functor we see that the free theorem for Kleisli composition predicts that any parametrically polymorphic function with the same type as (ordinary) function composition is associative. Can you show that function composition is uniquely defined by its parametricity property?

□

**Exercise 5.3** Derive the free theorem for catamorphisms from the polymorphic type of  $f \mapsto ([f])$ . Show that the fusion law is an instance of the free theorem.

□

### 5.3 Relators

As we have argued, an extension of the calculus of datatypes to relations is desirable from a practical viewpoint. In view of Reynolds' abstraction theorem, it is also highly desirable from a theoretical viewpoint, at least if one's goal is to develop generic programming. We have also shown how the product functor is extended to relations. In a relational theory of datatypes, all functors are extended to relations in such a way that when restricted to functions all their algebraic properties remain unchanged. Functors extended in this way are called *relators*.

The formal framework for this extension is known as an *allegory*. An allegory is a category with additional structure, the additional structure capturing the most essential characteristics of relations. The additional axioms are as follows. First of all, relations of the same type are ordered by the *partial order*  $\subseteq$  and composition is monotonic with respect to this order. That is,

$$S_1 \bullet T_1 \subseteq S_2 \bullet T_2 \quad \Leftarrow \quad S_1 \subseteq S_2 \quad \wedge \quad T_1 \subseteq T_2 \quad .$$

Secondly, for every pair of relations  $R, S :: A \sim B$ , their *intersection (meet)*  $R \cap S$  exists and is defined by the following universal property, for each  $X :: A \sim B$ ,

$$X \subseteq R \quad \wedge \quad X \subseteq S \quad \equiv \quad X \subseteq R \cap S \quad .$$

Finally, for each relation  $R :: A \sim B$  its *converse*  $R^\cup :: B \sim A$  exists. The converse operator satisfies the requirements that it is its own Galois adjoint, that is,

$$R^\cup \subseteq S \quad \equiv \quad R \subseteq S^\cup \quad ,$$

and is contravariant with respect to composition,

$$(R \bullet S)^\cup = S^\cup \bullet R^\cup \quad .$$

All three operators of an allegory are connected by the *modular law*, also known as Dedekind's law [41]:

$$R \bullet S \cap T \subseteq (R \cap T \bullet S^\cup) \bullet S \quad .$$



Now, a *relator* is a monotonic functor that commutes with converse. That is, the functor  $F$  is a relator iff,

$$(3) \quad FR \bullet FS = F(R \bullet S) \quad \text{for each } R :: A \sim B \text{ and } S :: B \sim C,$$

$$(4) \quad F\text{id}_A = \text{id}_{FA} \quad \text{for each } A,$$

$$(5) \quad FR \subseteq FS \iff R \subseteq S \quad \text{for each } R :: A \sim B \text{ and } S :: A \sim B,$$

$$(6) \quad (FR)^\cup = F(R^\cup) \quad \text{for each } R :: A \sim B.$$

*Relators extend functors* A design requirement which led to the above definition of a relator [4, 5] is that a relator should extend the notion of a functor but in such a way that it coincides with the latter notion when restricted to functions. Formally, relation  $R :: A \sim B$  is everywhere defined or *total* iff

$$\text{id}_B \subseteq R^\cup \bullet R,$$

and relation  $R$  is single-valued or *simple* iff

$$R \bullet R^\cup \subseteq \text{id}_A.$$

A *function* is a relation that is both total and simple. It is easy to verify that total and simple relations are closed under composition. Hence, functions are closed under composition too. In other words, the functions form a sub-category. For an allegory  $\mathcal{A}$ , we denote the sub-category of functions by  $Map(\mathcal{A})$ . Moreover, it is easily shown that our definition guarantees that relators preserve simplicity and totality, and thus functionality of relations.

Having made the shift from categories to allegories, the extension of the functional theory of datatypes in chapter 2 is surprisingly straightforward (which is another reason why not doing it is short-sighted). The extension of the disjoint sum functor to a disjoint sum relator can be done in such a way that all the properties of  $+$  and  $\nabla$  remain valid, as is the case for the extension of the theory of initial algebras, catamorphisms and type functors. For example, catamorphisms with relations as arguments are well-defined and satisfy the fusion property, the map-fusion property etc. There is, however, one catch — the process of dualising properties of disjoint sum to properties of cartesian product is not valid. Indeed, almost all of the properties of cartesian product that we presented are not valid, in the form presented here, when the variables range over arbitrary relations. (The banana split theorem is a notable exception.)

An example of what goes wrong is the fusion law. Consider  $\text{id} \triangle \text{id} \bullet R$  and  $R \triangle R$ , where  $R$  is a relation. If  $R$  is functional —that is, if for each  $y$  there is at most one  $x$  such that  $(x, y) \in R$  then these two are equal. This is an instance of the fusion law presented earlier. However, if  $R$  is not functional then they may not be equal. Take  $R$  to be, for example, the relation  $\{(0, 0), (1, 0)\}$  in which both 0 and 1 are related to 0. Then,

$$\text{id} \triangle \text{id} \bullet R = \{((0, 0), 0), ((1, 1), 0)\}$$

whereas

$$R \triangle R = \{((0,0), 0), ((1,1), 0), ((0,1), 0), ((1,0), 0)\} .$$

The relation  $\text{id} \triangle \text{id}$  is the *doubling* relation: it relates a pair of values to a single value whereby all the values are equal. Thus,  $\text{id} \triangle \text{id} \bullet R$  relates a pair of *equal* values to 0. On the other hand,  $R \triangle R$  relates a pair of values to a single value, whereby each component of the pair is related by  $R$  to the single value. The difference thus arises from the nondeterminism in  $R$ .

In conclusion, extending the functional theory of datatypes to relations is desirable but not without pitfalls. The pitfalls are confined, however, to the properties of cartesian product. We give no formal justification for this. The reader will just have to trust us that in the ensuing calculations, where one or more argument is a relation, that the algebraic properties that we exploit are indeed valid.

*Membership* We have argued that a datatype is not just a mapping from types to types but also a functor. We have now argued that a datatype is a relator. For the correctness of the generic unification algorithm we also need to know that a membership relation can be defined on a datatype.

The full theory of membership and its consequences has been developed by Hoogendijk and De Moor [24, 22]. Here we give only a very brief account.

Let  $F$  be a relator. A *membership relation* on  $F$  is a parametrically polymorphic relation  $\text{mem}$  of type  $a \sim Fa$  for all  $a$ . Parametricity means that for all relations  $R$ ,

$$\text{mem} \bullet FR \supseteq R \bullet \text{mem} .$$

In fact,  $\text{mem}$  is required to be the largest parametrically polymorphic relation of this type.

The existence of a membership relation captures the idea that a datatype is a structured repository of information. The relation  $\text{mem}_a$  holds between a value  $x$  of type  $a$  and an  $F$ -structure of  $a$ 's if  $x$  is stored somewhere in the  $F$ -structure. The parametericity property expresses the fact that determining membership is independent of the type  $a$ , and the fact that  $\text{mem}$  is the largest relation of its type expresses the idea that determining membership is independent of the position in the data structure at which a value is stored.

The parametericity property has the following consequence which we shall have occasion to use. For all (total) functions  $f$  of type  $a \rightarrow b$ ,

$$f \bullet \text{mem}_a = \text{mem}_b \bullet Ff .$$

## 5.4 Occurs-in

This section contains a proof of the generic statement that two expressions are not unifiable if one occurs in the other. We define a (generic) relation `occurs_properly_in` and we then show that `occurs_properly_in` is indeed a proper ordering on expressions (that is, if expression  $x$  `occurs_properly_in` expression  $y$  then  $x$  and  $y$  are different). We also show

that the `occurs_properly_in` relation is invariant under substitution. Thus, if expression  $x$  `occurs_properly_in` expression  $y$  no substitution can unify them. To show that `occurs_properly_in` is proper we define a (generic) function `size` of type  $F^*V \rightarrow \mathbb{N}$  and we show that `size` is preserved by the relation `occurs_properly_in`. The definition of `size` involves a restriction on the relator  $F$  which is used to guarantee correctness of the algorithm<sup>8</sup>.

**Definition 7** The relation `occurs_properly_in` of type  $F^*V \sim F^*V$  is defined by

$$\text{occurs\_properly\_in} = (\text{mem} \cdot \text{embr}_V^\cup)^+ .$$

(Recall that `mem` is the membership relation of  $F$  and that  $\text{embr}_V = \text{in}_{V^{\kappa+F}} \cdot \text{inr}$  where  $(F^*V, \text{in}_{V^{\kappa+F}})$  is an initial algebra.) Informally, the relation  $\text{embr}_V^\cup$  (which has type  $F^*V \sim F^*V$ ) destructs an element of  $F^*V$  into an  $F$  structure and then `mem` identifies the data stored in that  $F$  structure. Thus  $\text{mem} \cdot \text{embr}_V^\cup$  destructs an element of  $F^*V$  into a number of immediate subcomponents. Application of the transitive closure operation repeats this process thus breaking the structure down into all its subcomponents.

□

In our first lemma we show that the `occurs_properly_in` relation is closed under substitutions. That is, for all substitutions  $f$ ,

$$x \text{ occurs\_properly\_in } y \Rightarrow (fx) \text{ occurs\_properly\_in } (fy) .$$

The property is formulated without mention of the points  $x$  and  $y$  and proved using point-free relation algebra.

**Lemma 8** For all substitutions  $f$ ,

$$\text{occurs\_properly\_in} \subseteq f^\cup \cdot \text{occurs\_properly\_in} \cdot f .$$

**Proof** Suppose  $f$  is a substitution. That is,  $f = g \circ \text{id}$  for some  $g$ . Since `occurs_properly_in` is the transitive closure of the relation  $\text{mem} \cdot \text{embr}_V^\cup$  it suffices to establish two properties: first, that  $f^\cup \cdot \text{occurs\_properly\_in} \cdot f$  is transitive and, second,

$$\text{mem} \cdot \text{embr}_V^\cup \subseteq f^\cup \cdot \text{occurs\_properly\_in} \cdot f .$$

The first of these is true for all functions  $f$  (i.e. relations  $f$  such that  $f \cdot f^\cup \subseteq \text{id}$ ). (To be precise, if  $R$  is a transitive relation and  $f$  is a function then  $f^\cup \cdot R \cdot f$  is transitive.) We leave its simple proof to the reader. The second is proved as follows:

$$\begin{aligned} & f^\cup \cdot \text{occurs\_properly\_in} \cdot f \\ \supseteq & \quad \{ \quad R^+ \supseteq R \quad \} \\ & f^\cup \cdot \text{mem} \cdot \text{embr}_V^\cup \cdot f \\ \supseteq & \quad \{ \quad \text{embr}_V \text{ is a function, definition of } \text{embr}_V \quad \} \\ & f^\cup \cdot \text{mem} \cdot \text{embr}_V^\cup \cdot f \cdot \text{in} \cdot \text{inr} \cdot \text{embr}_V^\cup \end{aligned}$$

<sup>8</sup> A more general proof [7] using the generic theory of  $F$ -reductivity [15, 14, 16] avoids this assumption and, indeed, avoids the introduction of the size function altogether.

$$\begin{aligned}
&= \{ f = g \diamond \text{id} = ([g \nabla \text{embr}_V]), \text{computation} \} \\
&\quad f^\cup \cdot \text{mem} \cdot \text{embr}_{V^\cup} \cdot g \nabla \text{embr}_V \cdot \text{id} + Ff \cdot \text{inr} \cdot \text{embr}_{V^\cup} \\
&= \{ \text{computation} \} \\
&\quad f^\cup \cdot \text{mem} \cdot \text{embr}_{V^\cup} \cdot \text{embr}_V \cdot Ff \cdot \text{embr}_{V^\cup} \\
&= \{ \text{embr}_{V^\cup} \cdot \text{embr}_V = \text{id} \} \\
&\quad f^\cup \cdot \text{mem} \cdot Ff \cdot \text{embr}_{V^\cup} \\
&\supseteq \{ \text{parametricity of mem} \} \\
&\quad \text{mem} \cdot Ff^\cup \cdot Ff \cdot \text{embr}_{V^\cup} \\
&\supseteq \{ F \text{ is a relator and } f \text{ is a total function.} \\
&\quad \text{Thus, } Ff^\cup \cdot Ff \supseteq \text{id} \} \\
&\quad \text{mem} \cdot \text{embr}_{V^\cup} .
\end{aligned}$$

□

We now define a function  $\text{size}$  of type  $F^*V \rightarrow \mathbb{N}$  by

$$\text{size} = ([\text{zero} \nabla (\text{succ} \cdot \Sigma \text{mem})]) .$$

Here,  $\Sigma$  is the *summation quantifier*. That is, for an arbitrary relation  $R$  with target  $\mathbb{N}$ ,

$$(\Sigma R)x = \Sigma(m: m R x: m) .$$

The assumption in the definition of  $\text{size}$  is that  $F$  is finitely branching: that is, for each  $F$  structure  $x$ , the number of  $m$  such that  $m \text{ mem } x$  is finite.

Expressed in terms of points, the next lemma says that if a term  $x$  occurs properly in a term  $y$  then the size of  $x$  is strictly less than the size of  $y$ .

### Lemma 9

$$\text{occurs\_properly\_in} \subseteq \text{size}^\cup \cdot < \cdot \text{size} .$$

**Proof** Note that  $\text{occurs\_properly\_in}$  and  $<$  are both transitive relations. This suggests that we use the leapfrog rule:

$$a \cdot b^* \subseteq c^* \cdot a \Leftrightarrow a \cdot b \subseteq c \cdot a$$

which is easily shown to extend to transitive closure:

$$a \cdot b^+ \subseteq c^+ \cdot a \Leftrightarrow a \cdot b \subseteq c \cdot a .$$

We have:

$$\begin{aligned}
&\text{occurs\_properly\_in} \subseteq \text{size}^\cup \cdot < \cdot \text{size} \\
&\equiv \{ \text{size is a total function,} \\
&\quad \text{definition of occurs\_properly\_in} \} \\
&\text{size} \cdot (\text{mem} \cdot \text{embr}_{V^\cup})^+ \subseteq < \cdot \text{size} \\
&\Leftrightarrow \{ < \text{ is transitive. Thus, } < = <^+ .
\end{aligned}$$

$$\begin{aligned}
& \text{Leapfrog rule } \} \\
& \text{size} \bullet \text{mem} \bullet \text{embr}_V^U \subseteq < \bullet \text{size} \\
\equiv & \{ \text{embr}_V \text{ is a total function } \} \\
& \text{size} \bullet \text{mem} \subseteq < \bullet \text{succ} \bullet \text{embr}_V \\
\equiv & \{ \text{definition of size, embr}_V \text{ and computation } \} \\
& \text{size} \bullet \text{mem} \subseteq < \bullet \text{succ} \bullet \Sigma \text{mem} \bullet F \text{size} \\
\equiv & \{ < \bullet \text{succ} = \leq \} \\
& \text{size} \bullet \text{mem} \subseteq \leq \bullet \Sigma \text{mem} \bullet F \text{size} \\
\Leftarrow & \{ \text{property of natural numbers: for all } R, R \subseteq \leq \bullet \Sigma R \\
& \text{That is, } m R x \Rightarrow m \leq \Sigma(m: m R x: m). \} \\
& \text{size} \bullet \text{mem} \subseteq \text{mem} \bullet F \text{size} \\
\equiv & \{ \text{size is a total function,} \\
& \text{parametricity of mem for functions } \} \\
& \text{true} .
\end{aligned}$$

□

**Corollary 10** Suppose  $F$  is a finitely branching relator. Then

$$x \text{ occurs\_properly\_in } y \Rightarrow x \neq y .$$

**Proof** By the above lemma,

$$x \text{ occurs\_properly\_in } y \Rightarrow \text{size } x < \text{size } y .$$

Thus, since  $m < n \Rightarrow m \neq n$ ,

$$x \text{ occurs\_properly\_in } y \Rightarrow x \neq y .$$

□

**Corollary 11** If  $x \text{ occurs\_properly\_in } y$  then  $x$  and  $y$  are not unifiable.

**Proof** By lemma 8, if  $x \text{ occurs\_properly\_in } y$  then,  $fx \text{ occurs\_properly\_in } fy$ , for every substitution  $f$ . Thus, for every substitution  $f$ ,  $fx \neq fy$ .

□

**Exercise 5.4** Take  $F$  to be  $(1+)$ . What is  $\text{occurs\_properly\_in}$ ? Show that the relation is proper. (Note that the membership relation for  $(1+)$  is  $\text{inr}^U$ .)

Take  $F$  to be  $a \times$  for some fixed  $a$ . What is  $\text{occurs\_properly\_in}$ ?

□

## 6 Solutions to Exercises

**1.1** Take  $\otimes$  to be set intersection,  $\oplus$  to be set union,  $\mathbf{0}$  to be the empty set and  $\mathbf{1}$  to be the universe of all colours. The initial value of  $a[i, j]$  is the singleton set containing the edge colour as its element

□

**2.5**

$\text{map}_{Error} f (\text{error } s) = \text{error } s$

$\text{map}_{Error} f (\text{ok } x) = \text{ok } (fx)$

$\text{map}_{Drawing} f (\text{above } x y) = \text{above } (\text{map}_{Drawing} f x) (\text{map}_{Drawing} f y)$

$\text{map}_{Drawing} f (\text{beside } x y) = \text{beside } (\text{map}_{Drawing} f x) (\text{map}_{Drawing} f y)$

$\text{map}_{Drawing} f (\text{atom } x) = \text{atom } (fx)$

□

**2.6**

$$\begin{aligned}
 & (f \nabla g) \Delta (h \nabla k) = (f \Delta h) \nabla (g \Delta k) \\
 \equiv & \quad \{ \quad \Delta\text{-characterisation} \quad \} \\
 & f \nabla g = \text{exl} \cdot (f \Delta h) \nabla (g \Delta k) \quad \wedge \quad h \nabla k = \text{exr} \cdot (f \Delta h) \nabla (g \Delta k) \\
 \equiv & \quad \{ \quad \nabla\text{-fusion} \quad \} \\
 & f \nabla g = (\text{exl} \cdot (f \Delta h)) \nabla (\text{exl} \cdot (g \Delta k)) \\
 & \wedge h \nabla k = (\text{exr} \cdot (f \Delta h)) \nabla (\text{exr} \cdot (g \Delta k)) \\
 \equiv & \quad \{ \quad \text{injectivity of } \nabla \quad \} \\
 & f = \text{exl} \cdot (f \Delta h) \quad \wedge \quad g = \text{exl} \cdot (g \Delta k) \\
 & \wedge h = \text{exr} \cdot (f \Delta h) \quad \wedge \quad k = \text{exr} \cdot (g \Delta k) \\
 \equiv & \quad \{ \quad \Delta\text{-computation} \quad \} \\
 & \text{true} \quad .
 \end{aligned}$$

□

**2.7** The most obvious example is multiplication and division in ordinary arithmetic. (Indeed this is where the two-dimensional notation is commonly used.) Addition and subtraction also abide with each other.

Examples in the text are: disjoint sum and composition, and cartesian product and composition. (Indeed all binary functors abide with composition.)

The example used by Hoare was conditionals. The binary operator **if**  $p$ , where  $p$  is a proposition, (which has two statements as arguments) abides with **if**  $q$ , where  $q$  is also a proposition.

□

**2.11** First, the  $\nabla$ -+ fusion rule:

$$\begin{aligned}
 & f \nabla g \cdot h + k = (f \cdot h) \nabla (g \cdot k) \\
 \equiv & \quad \{ \quad \nabla \text{ characterisation} \quad \}
 \end{aligned}$$

$$\begin{aligned}
& f \nabla g \cdot h+k \cdot \text{inl} = f \cdot h \cdot \text{inl} \wedge f \nabla g \cdot h+k \cdot \text{inr} = g \cdot k \cdot \text{inr} \\
\equiv & \quad \{ \text{computation rules (applied four times)} \} \\
& \text{true} .
\end{aligned}$$

Second, the identity rule:

$$\begin{aligned}
& \text{id} + \text{id} \\
= & \quad \{ \text{definition of } + \} \\
& \text{inl} \nabla \text{inr} \\
= & \quad \{ \text{above} \} \\
& \text{id} .
\end{aligned}$$

□

**2.12** The pattern functor for **Bin** is  $\text{Exl} + (\text{Exr} \times \text{Exr})$  and for **Rose** is  $(\text{Exl} \times (\text{List } \text{Exr}))$ . That is, for **Bin** it is the binary functor mapping  $a$  and  $z$  to  $a + (z \times z)$ , which is polynomial, and for **Rose** it is the binary functor mapping  $a$  and  $z$  to  $a \times (\text{List } z)$ , which is not polynomial.

□

### 2.13

$$\begin{aligned}
& \text{even} \cdot \text{zero} \nabla \text{succ} = \text{true} \nabla \text{not} \cdot 1 + \text{even} \\
\equiv & \quad \{ \nabla \text{ fusion and } \nabla + \text{ fusion,} \\
& \quad \text{definition of functor } +1 \} \\
& (\text{even} \cdot \text{zero}) \nabla (\text{even} \cdot \text{succ}) = (\text{true} \cdot \text{id}_1) \nabla (\text{not} \cdot \text{even}) \\
\equiv & \quad \{ \text{true} \cdot \text{id}_1 = \text{true}, \nabla \text{ is injective} \} \\
& \text{even} \cdot \text{zero} = \text{true} \wedge \text{even} \cdot \text{succ} = \text{not} \cdot \text{even} \\
\equiv & \quad \{ \text{extensionality, identifying values zero and true} \\
& \quad \text{with functions zero and true with domain } 1 \} \\
& \text{even}(\text{zero}) = \text{true} \wedge \forall (n :: \text{even}(\text{succ } n) = \text{not}(\text{even } n)) .
\end{aligned}$$

□

### 2.14

$$\begin{aligned}
& \text{out} \cdot \text{in} \\
= & \quad \{ \text{definition of out} \} \\
& ([F \text{in}]) \cdot \text{in} \\
= & \quad \{ \text{computation rule} \} \\
& F \text{in} \cdot F([F \text{in}]) \\
= & \quad \{ F \text{ is a functor} \} \\
& F(\text{in} \cdot ([F \text{in}]))
\end{aligned}$$

$$\begin{aligned}
&= \{ \text{definition of out} \} \\
&\quad F(\text{in} \bullet \text{out}) \\
&= \{ \text{in} \bullet \text{out} = \text{id}_{\mu F} \} \\
&\quad F \text{id}_{\mu F} \\
&= \{ F \text{ is a functor} \} \\
&\quad \text{id}_{F \mu F} .
\end{aligned}$$

□

**2.16** We have

$$\text{NoOfTips} = ([1^{\kappa} \nabla \text{add0}])$$

where  $\text{add0}(m, n) = m+n$ , and

$$\text{NoOfJoins} = ([0^{\kappa} \nabla \text{add1}])$$

where  $\text{add1}(m, n) = m+n+1$ . Now,

$$\begin{aligned}
&f \bullet \text{NoOfTips} = \text{NoOfJoins} \\
\Leftarrow &\{ \text{definitions and fusion} \} \\
&f \bullet 1^{\kappa} \nabla \text{add0} = 0^{\kappa} \nabla \text{add1} \bullet \text{id} + (f \times f) \\
\equiv &\{ \text{fusion} \} \\
&(f \bullet 1^{\kappa}) \nabla (f \bullet \text{add0}) = 0^{\kappa} \nabla (\text{add1} \bullet f \times f) \\
\equiv &\{ \text{injectivity} \} \\
&f \bullet 1^{\kappa} = 0^{\kappa} \quad \wedge \quad f \bullet \text{add0} = \text{add1} \bullet f \times f \\
\equiv &\{ \text{pointwise definitions, for all } m \text{ and } n \} \\
&f1 = 0 \quad \wedge \quad f(m+n) = fm+1+fn \\
\Leftarrow &\{ \text{arithmetic, for all } m \} \\
&fm = m-1 .
\end{aligned}$$

We conclude that there is always one less join in a *Bin* than there are tips.

□

**2.17**

$$\begin{aligned}
&([f]) \triangle ([g]) \bullet \text{in} = \chi \bullet F([f] \triangle [g]) \\
\equiv &\{ \triangle \text{ fusion} \} \\
&([f] \bullet \text{in}) \triangle ([g] \bullet \text{in}) = \chi \bullet F([f] \triangle [g]) \\
\equiv &\{ \text{catamorphism computation} \} \\
&(f \bullet F([f])) \triangle (g \bullet F([g])) = \chi \bullet F([f] \triangle [g]) \\
\equiv &\{ \triangle \text{ characterisation} \} \\
&f \bullet F([f]) = \text{exl} \bullet \chi \bullet F([f] \triangle [g]) \\
\wedge &g \bullet F([g]) = \text{exr} \bullet \chi \bullet F([f] \triangle [g]) .
\end{aligned}$$



Once again, we continue with just one of the conjuncts, the other being solved by symmetry.

$$\begin{aligned}
& f \bullet F([f]) = \text{exl} \bullet \chi \bullet F([f] \triangle [g]) \\
\equiv & \quad \{ \text{postulate } \chi = \alpha \triangle \beta \} \\
& f \bullet F([f]) = \text{exl} \bullet \alpha \triangle \beta \bullet F([f] \triangle [g]) \\
\equiv & \quad \{ \triangle \text{ computation} \} \\
& f \bullet F([f]) = \alpha \bullet F([f] \triangle [g]) \\
\equiv & \quad \{ \text{postulate } \alpha = f \bullet \gamma \} \\
& f \bullet F([f]) = f \bullet \gamma \bullet F([f] \triangle [g]) \\
\Leftarrow & \quad \{ F \text{ respects composition, } \triangle \text{ computation} \} \\
& \gamma = F \text{exl} .
\end{aligned}$$

Combining the two postulates with the final statement, we get

$$([ \chi ]) = ([f] \triangle [g]) \Leftarrow \chi = (f \bullet F \text{exl}) \triangle (g \bullet F \text{exr}) .$$

□

**2.19** Substituting  $(a \circlearrowleft)$  for  $F$  in the catamorphism rule we get the rule:

$$h \bullet ([\varphi]) = ([\psi]) \Leftarrow h \bullet \varphi = \psi \bullet \text{id} \circlearrowleft h .$$

This is the fusion rule used below.

$$\begin{aligned}
& ([f]) \bullet (\tau(\circlearrowleft) g) = ([h]) \\
\equiv & \quad \{ \tau(\circlearrowleft) g = (\text{in} \bullet g \circlearrowleft \text{id}) \} \\
& ([f]) \bullet (\text{in} \bullet g \circlearrowleft \text{id}) = ([h]) \\
\Leftarrow & \quad \{ \text{fusion rule} \} \\
& ([f]) \bullet \text{in} \bullet g \circlearrowleft \text{id} = h \bullet \text{id} \circlearrowleft ([f]) \\
\equiv & \quad \{ \text{catamorphism computation} \} \\
& f \bullet \text{id} \circlearrowleft ([f]) \bullet g \circlearrowleft \text{id} = h \bullet \text{id} \circlearrowleft ([f]) \\
\equiv & \quad \{ \circlearrowleft \text{ is a binary functor. Thus,} \\
& \quad \text{id} \circlearrowleft ([f]) \bullet g \circlearrowleft \text{id} = g \circlearrowleft ([f]) = g \circlearrowleft \text{id} \bullet \text{id} \circlearrowleft ([f]) \} \\
& f \bullet g \circlearrowleft \text{id} \bullet \text{id} \circlearrowleft ([f]) = h \bullet \text{id} \circlearrowleft ([f]) \\
\Leftarrow & \quad \{ \text{cancellation} \} \\
& f \bullet g \circlearrowleft \text{id} = h .
\end{aligned}$$

We have thus established the rule:

$$([f]) \bullet (\tau(\circlearrowleft) g) = ([f \bullet g \circlearrowleft \text{id}]) .$$

□

2.20 First,

$$\begin{aligned}
& \tau(\otimes) \text{id}_a \\
= & \quad \{ \text{definition} \} \\
& (\text{in} \bullet \text{id} \otimes \text{id}) \\
= & \quad \{ \otimes \text{ respects identities,} \\
& \quad \text{identity is the unit of composition} \} \\
& (\text{in}) \\
= & \quad \{ \text{identity rule} \} \\
& \text{id}_{\tau(\otimes) a} .
\end{aligned}$$

Second,

$$\begin{aligned}
& \tau(\otimes) (f \bullet g) \\
= & \quad \{ \text{definition} \} \\
& (\text{in} \bullet (f \bullet g) \otimes \text{id}) \\
= & \quad \{ \text{id} = \text{id} \bullet \text{id}, \otimes \text{ respects composition} \} \\
& (\text{in} \bullet f \otimes \text{id} \bullet g \otimes \text{id}) \\
= & \quad \{ \text{exercise 2.19} \} \\
& (\text{in} \bullet f \otimes \text{id}) \bullet (\tau(\otimes) g) \\
= & \quad \{ \text{definition} \} \\
& (\tau(\otimes) f) \bullet (\tau(\otimes) g) .
\end{aligned}$$

□

4.1 To express  $Mf$  we use the last of the three monad equalities:

$$\begin{aligned}
& Mf \\
= & \quad \{ \text{identities} \} \\
& \eta \diamond (Mf \bullet \text{id}) \\
= & \quad \{ \text{monad equality} \} \\
& (\eta \bullet f) \diamond \text{id} .
\end{aligned}$$

Using  $\text{mul} = \text{id} \diamond \text{id}$ , we obtain that, for the functor  $Set$ ,

$$\text{mul } x = \{z \mid \exists(y:: z \in y \wedge y \in x)\} .$$

The equalities are proven as follows: First,

$$\begin{aligned}
& \text{mul} \bullet M\text{mul} \\
= & \quad \{ \text{mul} = \text{id} \diamond \text{id} \} \\
& (\text{id} \diamond \text{id}) \bullet M\text{mul} \\
= & \quad \{ \text{2nd monad equality, id is identity of composition} \}
\end{aligned}$$

$$\begin{aligned}
& \text{id} \diamond M \text{mul} \\
= & \quad \{ \text{id is identity of composition,} \\
& \quad \text{3rd monad equality, id is identity of composition} \} \\
& \text{mul} \diamond \text{id} \\
= & \quad \{ \text{mul} = \text{id} \diamond \text{id, Kleisli composition is associative,} \\
& \quad \text{mul} = \text{id} \diamond \text{id} \} \\
& \text{id} \diamond \text{mul} \\
= & \quad \{ \text{id is identity of composition, 2nd monad equality} \} \\
& (\text{id} \diamond \text{id}) \bullet \text{mul} \\
= & \quad \{ \text{mul} = \text{id} \diamond \text{id} \} \\
& \text{mul} \bullet \text{mul} .
\end{aligned}$$

Second,

$$\begin{aligned}
& \text{mul} \bullet \eta \\
= & \quad \{ \text{mul} = \text{id} \diamond \text{id, 2nd monad equality} \} \\
& \text{id} \diamond \eta \\
= & \quad \{ \eta \text{ is unit of Kleisli composition} \} \\
& \text{id} .
\end{aligned}$$

Third,

$$\begin{aligned}
& \text{mul} \bullet M\eta \\
= & \quad \{ \text{mul} = \text{id} \diamond \text{id, 2nd monad equality} \} \\
& \text{id} \diamond M\eta \\
= & \quad \{ \text{id is identity of composition, 3rd monad equality,} \\
& \quad \text{id is identity of composition} \} \\
& \eta \diamond \text{id} \\
= & \quad \{ \eta \text{ is unit of Kleisli composition} \} \\
& \text{id} .
\end{aligned}$$

□

**4.2**  $(a \times)^* 1$  is *List a*. The Kleisli identity is the function mapping  $x$  to  $[x]$ . The multiplier is the function `concat` that concatenates a list of lists to a list, preserving the order of the elements. The Kleisli composition  $g \diamond f$  first applies  $f$  to a value  $x$  of type  $a$ , which results in a list of  $b$ 's. Then  $g$  is mapped to all the elements of this list, and the resulting list of lists of  $c$ 's is flattened to a list of  $c$ 's.

□

**4.3** Since  $(1+)^* \emptyset = \mathbb{N}$  we obtain from the fusion theorem that  $(1+)^* V = \mathbb{N} \times (V+1)$ . Specifically,

$$\begin{aligned}
& \mathbb{N} \times (V+1) \text{ is an initial } X :: V+(1+X) \text{ algebra} \\
\Leftarrow & \quad \{ \text{fusion, } \mathbb{N} \text{ is an initial } 1+ \text{ algebra} \} \\
& \forall (X :: (1+X) \times (V+1) \cong V+(1+(X \times (V+1)))) \\
\Leftarrow & \quad \{ \text{rig} \} \\
& \text{true} .
\end{aligned}$$

The witness to the last step,  $\text{rig}$ , is the inverse of a natural isomorphism  $\text{rig}$  of type

$$Y+(1+(X \times (Y+1))) \rightarrow (1+X) \times (Y+1) .$$

It is easily constructed:

$$\text{rig} = ((\text{inl} \bullet!) \triangle \text{inl}) \nabla ((\text{inl} \triangle \text{inr}) \nabla (\text{inr} \times \text{id})) .$$

The initial algebra is  $\text{in}_{(1+)^*} = (\text{zero} \nabla \text{succ}) \times \text{id} \bullet \text{rig}$ .

□

#### 4.4

$$\begin{aligned}
& f \diamond (g \diamond h) = (f \diamond g) \diamond h \\
\equiv & \quad \{ \text{definition} \} \\
& ((f \nabla \text{embr}) \bullet (g \diamond h)) = (((f \diamond g) \nabla \text{embr}) \bullet h) \\
\Leftarrow & \quad \{ \text{definition of } f \diamond g, \text{ cancel } \bullet h \} \\
& ((f \nabla \text{embr}) \bullet (g \nabla \text{embr})) = (((f \diamond g) \nabla \text{embr}) \\
\Leftarrow & \quad \{ \text{fusion, definition of embr} \} \\
& ((f \nabla \text{embr}) \bullet g \nabla (\text{in} \bullet \text{inr})) = (f \diamond g) \nabla (\text{in} \bullet \text{inr}) \bullet \text{id} + F((f \nabla \text{embr})) \\
\equiv & \quad \{ \text{fusion properties of disjoint sum,} \\
& \quad \nabla \text{ is injective} \} \\
& ((f \nabla \text{embr}) \bullet g = f \diamond g \\
& \wedge ((f \nabla \text{embr}) \bullet \text{in} \bullet \text{inr} = \text{in} \bullet \text{inr} \bullet F((f \nabla \text{embr}))) \\
\equiv & \quad \{ \text{definition of } f \diamond g, \text{ computation laws} \} \\
& \text{true} .
\end{aligned}$$

The verification that  $\text{embr}$  is its neutral element is a straightforward use of the computation rules.

□

**5.1** Substituting the identity function for  $h$  and  $l$ , we get

$$f_0 \diamond (Mk \bullet g_1) = (f_0 \bullet k) \diamond g_1 .$$

Substituting the identity function for  $k$  and  $h$ , we get

$$(f_0 \diamond g_0) \bullet l = f_1 \diamond (g_0 \bullet l) .$$

□

**5.2** Suppose  $\circ$  is a function that has the same polymorphic type as function composition. Then, if it satisfies the parametricity property of composition, it is the case that, for all relations  $R, S$  and  $T$  and all functions  $f_0, f_1, g_0$  and  $g_1$ , if

$$(f_0, f_1) \in S \rightarrow R \quad \wedge \quad (g_0, g_1) \in T \rightarrow S$$

then

$$(f_0 \circ f_1, g_0 \circ g_1) \in T \rightarrow R .$$

Take  $R$  to be the singleton set  $\{(f(gc), f(gc))\}$ ,  $S$  to be the singleton set  $\{(gc, gc)\}$  and  $T$  to be the singleton set  $\{(c, c)\}$ , where  $f$  and  $g$  are two functions, and  $c$  is some value such that  $f(gc)$  is defined. Then  $(f, f) \in S \rightarrow R$  and  $(g, g) \in T \rightarrow S$ . So  $(f \circ g, f \circ g) \in T \rightarrow R$ . That is,  $(f \circ g)(c) = f(gc)$ . Thus, by extensionality,  $f \circ g = f \bullet g$ . The parametricity property does indeed uniquely characterise function composition!

□

**5.3** The type of an  $F$ -catamorphism is

$$(Fa \rightarrow a) \rightarrow (\mu F \rightarrow a) .$$

The free theorem is thus that, for all relations  $R$  and all functions  $f$  and  $g$ , if

$$(f, g) \in FR \rightarrow R$$

then

$$((f]), ([g])) \in \text{id}_{\mu F \rightarrow R} .$$

Taking  $R$  to be a function  $h$  and use the point-free definition of  $\rightarrow$ , this is the statement that

$$f \bullet Fh = h \bullet g \quad \Rightarrow \quad ([f]) = h \bullet ([g]) .$$

□

**5.4** Instantiating  $F$  to  $(1+)$  we get

$$\begin{aligned} & \text{occurs\_properly\_in}_{1+} \\ = & \quad \{ \text{definition} \} \\ & (\text{mem}_{1+} \bullet (\text{in}_{(1+)^*} \bullet \text{inl})^{\cup})^+ \\ = & \quad \{ \text{mem}_{1+} = \text{inr}^{\cup}, \quad \text{in}_{(1+)^*} = (\text{zero} \vee \text{succ}) \times \text{id} \bullet \text{rig} \} \\ & ((\text{zero} \vee \text{succ}) \times \text{id} \bullet \text{rig} \bullet \text{inl} \bullet \text{inr})^{\cup+} \\ = & \quad \{ \text{definition of rig, computation} \} \\ & (\text{succ} \times \text{id})^{\cup+} . \end{aligned}$$

A pair  $(m, x)$  “occurs properly in” a pair  $(n, y)$  if  $m < n$  and  $x = y$ . This particular instance of `occurs_properly_in` is thus proper in the sense that if  $u$  “occurs properly in”  $v$  then  $u$  and  $v$  are not equal.

$F^*1$  is `List a`, membership is the projection `exr` and `occurs_properly_in` is the relation “is a (proper) tail of”.

□

## References

1. C.J. Aarts, R.C. Backhouse, P. Hoogendijk, T.S. Voermans, and J. van der Woude. A relational theory of datatypes. Available via World-Wide Web at <http://www.win.tue.nl/cs/wp/papers>, September 1992.
2. Lennart Augustsson Cayenne, a language with dependent types. This volume, 1999.
3. R. C. Backhouse, M. Bijsterveld, R. van Geldrop, and J.C.S.P. van der Woude. Category theory as coherently constructive lattice theory. Department of Mathematics and Computing Science, Eindhoven University of Technology. Working document. Available via World-Wide Web at <http://www.win.tue.nl/cs/wp/papers>, Last revision: March 1997, 146 pages, 1995.
4. R.C. Backhouse, P. de Bruin, P. Hoogendijk, G. Malcolm, T.S. Voermans, and J. van der Woude. Polynomial relators. In M. Nivat, C.S. Rattray, T. Rus, and G. Scollo, editors, *Proceedings of the 2nd Conference on Algebraic Methodology and Software Technology, AMAST'91*, pages 303–326. Springer-Verlag, Workshops in Computing, 1992.
5. R.C. Backhouse, P. de Bruin, G. Malcolm, T.S. Voermans, and J. van der Woude. Relational catamorphisms. In Möller B., editor, *Proceedings of the IFIP TC2/WG2.1 Working Conference on Constructing Programs from Specifications*, pages 287–318. Elsevier Science Publishers B.V., 1991.
6. R.C. Backhouse and B.A. Carré. Regular algebra applied to path-finding problems. *Journal of the Institute of Mathematics and its Applications*, 15:161–186, 1975.
7. Roland Backhouse. Fixed point calculus applied to generic programming: Part 1. In Zoltan Esik, editor, *Proceedings, Workshop on Fixed Points in Computer Science*, August 1998.
8. Roland C. Backhouse, J.P.H.W. van den Eijnde, and A.J.M. van Gasteren. Calculating path algorithms. *Science of Computer Programming*, 22(1–2):3–19, 1994.
9. G. Bellè, C.B. Jay, and E. Moggi. Functorial ML. In *PLILP96*, volume 1140 of *LNCS*. Springer-Verlag, 1996.
10. Richard Bird, Oege de Moor, and Paul Hoogendijk. Generic functional programming with types and relations. *J. of Functional Programming*, 6(1):1–28, January 1996.
11. Richard S. Bird and Oege de Moor. *Algebra of Programming*. Prentice-Hall International, 1996.
12. R.S. Bird. An introduction to the theory of lists. In M. Broy, editor, *Logic of Programming and Calculi of Discrete Design*. Springer-Verlag, 1987. NATO ASI Series, vol. F36.
13. Robin Cockett and Tom Fukushima. About Charity. Yellow Series Report No. 92/480/18, Dep. of Computer Science, Univ. of Calgary, 1992.
14. H. Doornbos. *Reductivity arguments and program construction*. PhD thesis, Eindhoven University of Technology, Department of Mathematics and Computing Science, June 1996.
15. Henk Doornbos and Roland Backhouse. Induction and recursion on datatypes. In B. Möller, editor, *Mathematics of Program Construction, 3rd International Conference*, volume 947 of *LNCS*, pages 242–256. Springer-Verlag, July 1995.
16. Henk Doornbos and Roland Backhouse. Reductivity. *Science of Computer Programming*, 26(1–3):217–236, 1996.
17. R.W. Floyd. Algorithm 97. Shortest Path. *Comm. ACM*, 5(6):345, June 1962.
18. Maarten M. Fokkinga. *Law and Order in Algorithmics*. PhD thesis, Universiteit Twente, The Netherlands, 1992.
19. Maarten M. Fokkinga. Datatype laws without signatures. *Mathematical Structures in Computer Science*, 6:1–32, 1996.
20. M.M. Fokkinga. Monadic maps and folds for arbitrary datatypes. Memoranda Informatica 94-28, University of Twente, June 1994.
21. E. Gamma, R. Helm, R. Johnson, and J. Vlissides. *Design Patterns – Elements of Reusable Object-Oriented Software*. Addison-Wesley, 1995.
22. Paul Hoogendijk. *A Generic Theory of Datatypes*. PhD thesis, Department of Mathematics and Computing Science, Eindhoven University of Technology, 1997.
23. Paul Hoogendijk and Roland Backhouse. When do datatypes commute? In Eugenio Moggi and Giuseppe Rosolini, editors, *Category Theory and Computer Science, 7th International Conference*, volume 1290 of *LNCS*, pages 242–260. Springer-Verlag, September 1997.
24. Paul Hoogendijk and Oege de Moor. What is a datatype? Technical Report 96/16, Department of Mathematics and Computing Science, Eindhoven University of Technology, 1996. Submitted to *Science of Computer Programming*. Available via World-Wide Web at <http://www.win.tue.nl/cs/wp/papers>.
25. P. Jansson and J. Jeuring. PolyP - a polytypic programming language extension. In *POPL '97: The 24th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 470–482. ACM Press, 1997.

26. P. Jansson and J. Jeuring. Functional pearl: Polytypic unification. *Journal of Functional Programming*, 1998. In press.
27. Patrik Jansson. Functional polytypic programming — use and implementation. Technical report, Chalmers Univ. of Tech., Sweden, 1997. Lic. thesis. Available from <http://www.cs.chalmers.se/~patrikj/lic/>.
28. C.B. Jay. A semantics for shape. *Science of Computer Programming*, 25(251–283), 1995.
29. C.B. Jay, G. Bellè, and E. Moggi. Functorial ML. Extended version of [9] in press for *Journal of Functional Programming '98*, 1998.
30. C.B. Jay and J.R.B. Cockett. Shapely types and shape polymorphism. In D. Sannella, editor, *ESOP '94: 5th European Symposium on Programming*, pages 302–316. Springer Verlag Lecture Notes in Computer Science, April 1994.
31. J. Jeuring. Polytypic pattern matching. In *Conference Record of FPCA '95, SIGPLAN-SIGARCH-WG2.8 Conference on Functional Programming Languages and Computer Architecture*, pages 238–248, 1995.
32. J. Jeuring and P. Jansson. Polytypic programming. In J. Launchbury, E. Meijer, and T. Sheard, editors, *Proceedings of the Second International Summer School on Advanced Functional Programming Techniques*, pages 68–114. Springer-Verlag, 1996. LNCS 1129.
33. G. Malcolm. *Algebraic data types and program transformation*. PhD thesis, Groningen University, 1990.
34. G. Malcolm. Data structures and program transformation. *Science of Computer Programming*, 14(2–3):255–280, October 1990.
35. L. Meertens. Algorithmics – towards programming as a mathematical activity. In *Proceedings of the CWI Symposium on Mathematics and Computer Science*, pages 289–334. North-Holland, 1986.
36. L. Meertens. Paramorphisms. *Formal Aspects of Computing*, 4(5):413–424, 1992.
37. Lambert Meertens. Calculate polytypically! In Herbert Kuchen and S. Doaitse Swierstra, editors, *Proceedings of the Eighth International Symposium PLILP '96 Programming Languages: Implementations, Logics and Programs*, volume 1140 of *Lecture Notes in Computer Science*, pages 1–16. Springer Verlag, 1996.
38. Oege de Moor and Ganesh Sittampalam. Generic program transformation. This volume, 1999.
39. A. Pardo. Monadic corecursion —definition, fusion laws, and applications—. *Electronic Notes in Theoretical Computer Science*, 11, 1998.
40. J.C. Reynolds. Types, abstraction and parametric polymorphism. In R.E. Mason, editor, *IFIP '83*, pages 513–523. Elsevier Science Publishers, 1983.
41. J. Riguet. Relations binaires, fermetures, correspondances de Galois. *Bulletin de la Société Mathématique de France*, 76:114–155, 1948.
42. B. Roy. Transitivité et connexité. *C.R. Acad. Sci.*, 249:216, 1959.
43. Fritz Ruehr. *Analytical and Structural Polymorphism Expressed Using Patterns Over Types*. PhD thesis, University of Michigan, 1992.
44. Tim Sheard. Automatic generation and use of abstract structure operators. *ACM TOPLAS*, 13(4):531–557, 1991.
45. P. Wadler. Theorems for free! In *4'th Symposium on Functional Programming Languages and Computer Architecture, ACM, London*, September 1989.
46. S. Warshall. A theorem on boolean matrices. *J. ACM*, 9:11–12, 1962.