

Generalised folds for nested datatypes

Richard Bird¹ and Ross Paterson²

Abstract. Nested datatypes generalise regular datatypes in much the same way that context-free languages generalise regular ones. Although the categorical semantics of nested types turns out to be similar to the regular case, the fold functions are more limited because they can only describe natural transformations. Practical considerations therefore dictate the introduction of a generalised fold function in which this limitation can be overcome. In the paper we show how to construct generalised folds systematically for each nested datatype, and show that they possess a uniqueness property analogous to that of ordinary folds. As a consequence, generalised folds satisfy fusion properties similar to those developed for regular datatypes. Such properties form the core of an effective calculational theory of inductive datatypes.

Keywords: Functional programming, program construction.

1. Introduction

A *nested* datatype is a parametrised inductive datatype whose declaration involves a change to the accompanying type parameter(s). In other words, the recursion is “nested” within a change of parameter. A formal definition is given in Section 3. Such types have also been called *non-regular* or *non-uniform*.

A simple example is provided by the type *Pow* of *power trees*, declared in Haskell by:

```
data Pow a  = Zero a | Succ (Pow (Pair a))  
type Pair a = a × a
```

(We have varied the syntax slightly, writing $a \times a$ where Haskell would use (a, a) .)
Elements of *Pow a* consist of pairs of pairs of pairs ... of values of type

¹ Programming Research Group, Oxford University, Wolfson Building, Parks Road, Oxford, OX1 3QD, UK.

² Department of Computer Science, City University, Northampton Square, London EC1V 0HB, UK.

a , where the depth of nesting is given by reading the associated constructor expression as a Peano numeral. For example,

$$\text{Succ}(\text{Succ}(\text{Zero}((1, 2), (3, 4))))$$

is a power tree of height 2. Each element therefore represents a perfectly balanced binary tree with labels in the leaves.

As this example suggests, nested datatypes can capture certain structural invariants in a way that regular datatypes cannot. Connelly and Lockwood Morris [CL95] use a nested type for modelling a generalisation of the trie data structure, and Okasaki [Oka98] puts nested datatypes to a variety of uses in the design of purely functional data structures. Bird and Paterson [BP99] use a nested datatype to describe de Bruijn notation for lambda expressions. Other examples of datatypes that have been shown very recently to be expressible by nested declarations include AVL trees, 2-3 trees, and square matrices.

The aim of this paper is to lay the groundwork for a useful calculational theory of nested datatypes; the work can be regarded as a successor to [BM98]. The categorical semantics is straightforward, but there is a complication. For a regular datatype, the fold function induced by initiality is a useful tool in the construction of practical programs. In nested datatypes, it isn't. The basic reason is that the argument of the fold is required to be a natural transformation, and the result of the fold is also a natural transformation. Consequently, functions that inspect the elements of a nested datatype cannot be defined as folds. This phenomenon motivates the introduction of a generalised fold operator in which these limitations can be overcome. The first contribution of the paper is to show how to define a generalised fold operator for each nested type.

In the calculational theory of regular inductive datatypes, everything hinges on the uniqueness property of the fold function: the fold is the unique function of its type satisfying a certain functional equation. From this flow the fusion laws of fold; such laws are of fundamental importance in reasoning about programs and improving their efficiency (see, e.g. [Bir98, Mal90b]). The second contribution of the paper is to show that generalised folds also possess a uniqueness property.

Having established the uniqueness property, we can set up appropriate fusion laws for reasoning about nested datatypes. The third contribution in this paper is to present such laws, and to discuss some of their consequences.

The rest of the paper is organised as follows. In Sections 2 and 3 we briefly review the standard theory of regular datatypes, and show how the theory extends to nested datatypes. These sections also contain the formal definitions of regular and nested datatypes. In Section 4 we introduce the generalised folds, and in Section 5 we present the fusion laws. In Section 6 we will justify our proofs of the fusion laws by showing that generalised folds are uniquely determined by their defining equations. Familiarity with basic category theory is assumed.

2. Semantics of non-parametrised datatypes

In a categorical setting, non-parametrised inductive datatypes are introduced by taking categorical least fixed points of appropriate functors (see e.g. [Lam70, BdM97, Hag87, Mal90a, MA86]). Given a suitable category \mathbf{C} , the least fixed point of a functor $F :: \mathbf{C} \rightarrow \mathbf{C}$, if it exists, is an object T of \mathbf{C} for which the arrow $\alpha :: F T \rightarrow T$ is an initial algebra in the category of F -algebras. The objects of this category are arrows $F A \rightarrow A$ of \mathbf{C} , and the arrows are homomorphisms

between algebras. Thus, for objects $f :: F A \rightarrow A$ and $g :: F B \rightarrow B$, an arrow $h :: f \rightarrow g$ is an arrow $h :: A \rightarrow B$ of \mathbf{C} such that $h \cdot f = g \cdot F h$.

The initiality of α means that to each algebra $f :: F A \rightarrow A$ there is associated a unique homomorphism $h :: \alpha \rightarrow f$, that is a function $h :: T \rightarrow A$ satisfying the equation

$$h \cdot \alpha = f \cdot F h.$$

The unique homomorphism h is denoted by *fold* f . It is a consequence of initiality that α is an isomorphism, whence $T \cong F T$; this explains the terminology ‘fixed point’.

One class of appropriate functors are the *polynomial* functors. In a category \mathbf{C} with products and coproducts, the polynomial functors $F :: \mathbf{C}^n \rightarrow \mathbf{C}$ are defined by the following grammar:

$$F ::= \underline{A}^n \mid + \mid \times \mid \Pi_i^n \mid F \cdot \langle F_1, \dots, F_n \rangle$$

The functor \underline{A}^n is an n -ary constant functor, delivering the type A for all arguments. The functor Π_i^n is an n -ary projection, selecting its i th argument. A special case is $Id = \Pi_1^1$. The superscripts n will be omitted in the rest of the paper. The functor $F \cdot \langle F_1, \dots, F_n \rangle$ denotes the composition of an n -ary functor F with n functors F_i , all of the same arity, so that

$$(F \cdot \langle F_1, \dots, F_n \rangle) A_1 \dots A_m = F (F_1 A_1 \dots A_m) \dots (F_n A_1 \dots A_m)$$

We omit the brackets when $n = 1$. We also write $F + G$ for $+ \cdot \langle F, G \rangle$ and similarly $F \times G$.

Provided the category \mathbf{C} has colimits of all countable chains (the category **Set** of sets and total functions is an example), every unary polynomial functor possesses a least fixed point.

As examples, we can introduce the datatype *Nat* of natural numbers as the least fixed point of $F = \underline{1} + Id$, lists of natural numbers as the least fixed point of $F = \underline{1} + \underline{Nat} \times Id$, and binary trees, with natural number labels, as the least fixed point of $F = \underline{1} + \underline{Nat} \times Id \times Id$.

Such declarations can be translated into Haskell fairly directly. For example, the datatype of lists of natural numbers can be declared by

```
data NatListF a = Nil | Cons (Nat × a)
newtype NatList = In (NatListF NatList)
```

The **data** declaration introduces the polynomial functor *NatListF*, and the **newtype** declaration introduces *NatList* as isomorphic to *NatListF NatList*. In Haskell, both type variables and type constructor variables are denoted using lower-case letters.

In general, Haskell programs may fail to terminate, so these datatypes conceptually include an additional ‘value’ representing non-termination, and the language defines how non-termination propagates through an expression. However, we shall be using a subset of Haskell in which all programs terminate, so that this is not an issue, and our programs may be modelled in our category \mathbf{C} . This means that we cannot use unrestricted recursion; we must show that any recursive definitions we use have unique solutions. An additional restriction in comparison with Haskell is that recursively defined types in this paper do not involve function types, so we are concerned solely with inductive types.

The constructor *In* corresponds to the initial algebra α . The converse function α^{-1} is given by the function *out*, defined in Haskell by

$$\begin{aligned} out & \quad \quad \quad :: NatList \rightarrow NatListF NatList \\ out (In\ x) & = x \end{aligned}$$

The type of lists of natural numbers can also be declared by

$$\mathbf{data}\ NatList = Nil \mid Cons (Nat \times NatList)$$

In this version, which is closer to the usual style of datatype declaration in functional programming, the polynomial functor $NatListF$ does not appear explicitly.

The function $fold$ may be defined in Haskell as follows:

$$\begin{aligned} fold & \quad \quad \quad :: (NatListF\ a \rightarrow a) \rightarrow NatList \rightarrow a \\ fold\ f & = f \cdot natlistF\ (fold\ f) \cdot out \end{aligned}$$

In Haskell, a free type variable, like a , is implicitly universally quantified, and may be instantiated to any type.

The function $natlistF$ is the functorial action of $NatListF$:

$$\begin{aligned} natlistF & \quad \quad \quad :: (a \rightarrow b) \rightarrow NatListF\ a \rightarrow NatListF\ b \\ natlistF\ f\ Nil & = Nil \\ natlistF\ f\ (Cons\ (n, x)) & = Cons\ (n, f\ x) \end{aligned}$$

For example, $sumNL = fold\ sumF$ sums a list, where

$$\begin{aligned} sumF & \quad \quad \quad :: NatListF\ Nat \rightarrow Nat \\ sumF\ Nil & = 0 \\ sumF\ (Cons\ (m, n)) & = m + n \end{aligned}$$

3. Semantics of parametrised datatypes

Next, we consider the extension of this theory to parametrised datatypes, like the following Haskell definition for lists:

$$\mathbf{data}\ List\ a = Nil \mid Cons (a \times List\ a)$$

As before, this may be rewritten in the equivalent form

$$\begin{aligned} \mathbf{data}\ Base\ a\ b & = Nil \mid Cons (a \times b) \\ \mathbf{newtype}\ List\ a & = In (Base\ a (List\ a)) \end{aligned}$$

To incorporate parametrised datatypes into the theory, one can proceed in essentially two different ways. The first method is standard, but the second extends to nested constructors. For simplicity, we will consider only datatypes parametrised by a single type variable; the extension to multiple parameters is straightforward.

3.1. First method: families of fixed points

The first method takes the fixed point of a binary functor B with respect to its second parameter. We begin with the observation that the partial application of a binary functor B to a type A yields a unary functor $B\ A$, whose least fixed point $T\ A$ can be constructed in the usual way. Then for each A there is an associated initial algebra α_A of type

$$\alpha_A \quad :: \quad B\ A\ (T\ A) \rightarrow T\ A.$$

The associated fold function takes an algebra $f :: B A Y \rightarrow Y$ as argument, and is uniquely defined by the assertion that for any arrow $h :: T A \rightarrow Y$, we have $h = \text{fold } f$ if and only if

$$h \cdot \alpha_A = f \cdot B \text{ id } h.$$

The datatype constructor T can be made into a functor by defining $T f :: T X \rightarrow T Y$ for an arrow $f :: X \rightarrow Y$ by

$$T f = \text{fold } (\alpha_Y \cdot B f \text{ id}).$$

Restated, this definition takes the form

$$T f \cdot \alpha_X = \alpha_Y \cdot B f (T f).$$

Hence α is a natural transformation from $B \cdot \langle Id, T \rangle$ to T .

The class of *regular* functors is defined as the closure of the class of polynomial bifunctors under least fixed point operations. For example, *List* is the least fixed point of a bifunctor $B = \underline{1} + \Pi_1 \times \Pi_2$, and the general tree constructor *Tree* is the least fixed point of a bifunctor $B = \Pi_1 \times (List \cdot \Pi_2)$.

For the *List* datatype constructor, the fold function is implemented by

$$\begin{aligned} \text{fold} &:: (Base\ a\ b \rightarrow b) \rightarrow List\ a \rightarrow b \\ \text{fold } f &= f \cdot \text{base id } (\text{fold } f) \cdot \text{out} \end{aligned}$$

where the function *base* implements the functorial action of *Base*:

$$\begin{aligned} \text{base} &:: (a \rightarrow c) \rightarrow (b \rightarrow d) \rightarrow Base\ a\ b \rightarrow Base\ c\ d \\ \text{base } f\ g\ Nil &= Nil \\ \text{base } f\ g\ (Cons\ (x, y)) &= Cons\ (f\ x, g\ y) \end{aligned}$$

Though it is defined in a slightly different way than usual, the function *fold* is essentially the same as the standard function *foldr* in functional programming. The functorial action of *List* (corresponding to the standard function *map*) may be defined using *fold*:

$$\begin{aligned} \text{list} &:: (a \rightarrow b) \rightarrow List\ a \rightarrow List\ b \\ \text{list } f &= \text{fold } (In \cdot \text{base } f \text{ id}) \end{aligned}$$

3.2. Second method: higher-order fixed points

The second method yields a larger class of datatype constructors. Instead of constructing fixed points in the category \mathbf{C} , we will work in the functor categories $\mathbf{C}^n \rightarrow \mathbf{C}$. The objects of $\mathbf{C}^n \rightarrow \mathbf{C}$ are n -ary functors over \mathbf{C} , and the arrows are natural transformations between functors. We shall write $F \rightarrow G$ for the set of natural transformations from F to G .

In this category, limits and colimits may be constructed pointwise from limits and colimits in \mathbf{C} [Mac71, V.3]. For example, the shorthand notations $F \times G$ and $F + G$ do in fact denote products and coproducts in this category³.

Moreover, we can define a functor $T :: \mathbf{C}^n \rightarrow \mathbf{C}$ as a fixed point of a suitable higher-order functor $F :: (\mathbf{C}^n \rightarrow \mathbf{C}) \rightarrow (\mathbf{C}^n \rightarrow \mathbf{C})$, via the same colimit construction used in \mathbf{C} . We shall refer to such functors as *hofunctors*. As before,

³ Everything we do here generalises in a straightforward way to functor categories $\mathbf{D} \rightarrow \mathbf{C}$, for an arbitrary category \mathbf{D} .

this means that $\alpha :: F T \rightarrow T$ is an initial object in the category of F -algebras. Note that α is a natural transformation by construction. Hence, for $f :: A \rightarrow B$, we have

$$T f \cdot \alpha_A = \alpha_B \cdot F T f.$$

The fold function induced by the initiality of α associates with each natural transformation $f :: F N \rightarrow N$, for some given functor N , a natural transformation $hfold f :: T \rightarrow N$ defined as the unique function h satisfying

$$h \cdot \alpha = f \cdot F h.$$

This is exactly the same definition as for *fold* on unparametrised datatypes, but installed at a higher level. One major difference though is that the functoriality of T is given by construction, not as an instance of the fold.

We define the class of *polynomial hofunctors* of the form $F X = P$, where P is a polynomial functor expression that may also include the functor variable X . The grammar for such functor expressions is

$$P ::= \underline{A}^n \mid + \mid \times \mid \Pi_i^n \mid P \cdot \langle P_1, \dots, P_n \rangle \mid X$$

A fixed point of a polynomial hofunctor is called a *nested* functor⁴. We consider three examples.

Example 3.1 The functor *List* is given as the least fixed point of the polynomial hofunctor $ListF X = \underline{1} + Id \times X$. Thus *List* is a nested functor which is also a regular functor. In Haskell, the type constructor *List* may be introduced by the declarations

```
data Base a b    = Nil | Cons (a × b)
type ListF x a  = Base a (x a)
newtype List a = In (ListF List a)
```

The **type** declaration introduces *ListF x a* as a synonym for *Base a (x a)*. The type constructor *List* is therefore exactly the same as before.

However, the fold function is different. This function may be expressed in Haskell as follows:

```
hfold    :: (∀ a. ListF n a → n a) → List b → n b
hfold f  = f · listF (hfold f) · out
```

The function *listF* expresses part of the functorial action of *ListF*: for each natural transformation $f :: X \rightarrow Y$, we have $listF f :: ListF X \rightarrow ListF Y$. In Haskell, we have

```
listF    :: (∀ a. x a → y a) → ListF x b → ListF y b
listF f  = base id f
```

The first arguments of *hfold* and *listF* are polymorphic functions, described using a *rank 2* type signature [McC84]. Such signatures are allowed in recent extensions to Haskell [Jon98, PL97]. Polymorphic functions correspond to natural transformations in a categorical setting.

The definition of *hfold* is equivalent to the earlier definition of the regular *fold*, but has a different type: the new version takes natural transformations to

⁴ Though we deal only with single recursive definitions here, our results could be extended to systems of multiple simultaneous definitions.

natural transformations. As a result, the new fold is less useful because individual list elements can never be inspected.

For example, if a fold is to produce a natural number, we would have to take $N = \underline{Nat}$, the constant functor that returns Nat on every type, so that $hfold f :: List \rightarrow \underline{Nat}$. This function will operate on lists of any type, and must satisfy the naturality property

$$hfold f = hfold f \cdot list k$$

for any function k , and thus cannot depend on elements of the list. We can define the length function, but little more. In the following section we will see how to overcome this limitation.

We have noted above that a functor $List$ is introduced automatically as part of the categorical semantics, but we have not yet defined its functorial action $list$ in Haskell. An inductive argument establishes that $list$ is uniquely determined by the naturality of $In :: ListF List \rightarrow List$. To express this naturality condition, we could extend the definition of $listF$ given above to the full functorial action of $ListF$, or simply expand $ListF$ in the definition of $List$, yielding the definition

$$\begin{aligned} list &:: (a \rightarrow b) \rightarrow List a \rightarrow List b \\ list f &= In \cdot base f (list f) \cdot out \end{aligned}$$

This is equivalent to the definition using $fold$ at the end of the previous section. Note, however, that $list$ cannot be defined using $hfold$.

Example 3.2 The type constructor $Nest$ is defined as the least fixed point of

$$NestF X = \underline{1} + Id \times (X \cdot Pair).$$

In Haskell the corresponding declarations are:

```
data Base a b = Nil | Cons (a × b)
type NestF x a = Base a (x (Pair a))
type Pair a = a × a
newtype Nest a = In (NestF Nest a)
```

An equivalent datatype is introduced by the ‘flattened’ declaration

$$\mathbf{data} \text{ Nest } a = Nil | Cons (a \times Nest (a \times a))$$

Unlike $List$, the datatype constructor $Nest$ is not a regular functor: occurrences of $Nest$ in its defining expression are “nested” within a change of type parameter from a to $Pair a$. In effect, values of type $Nest a$ correspond to lists in which the first element has type a , the second element has type $(a \times a)$, the third has type $((a \times a) \times (a \times a))$, and so on. A variant of $Nest$ was used in [Oka98] as a basis for an implementation of lists with an efficient indexing operation.

The fold function for $Nest$ may be implemented in Haskell by expanding $NestF$ in terms of the bifunctor $Base$:

$$\begin{aligned} hfold &:: (\forall a. Base a (n (Pair a)) \rightarrow n a) \rightarrow Nest b \rightarrow n b \\ hfold f &= f \cdot base id (hfold f) \cdot out \end{aligned}$$

Similarly, the mapping function $nest$ may be defined by

$$\begin{aligned}
\text{nest} & \quad \text{:: } (a \rightarrow b) \rightarrow \text{Nest } a \rightarrow \text{Nest } b \\
\text{nest } f & \quad = \text{In} \cdot \text{base } f (\text{nest } (\text{pair } f)) \cdot \text{out} \\
\text{pair} & \quad \text{:: } (a \rightarrow b) \rightarrow \text{Pair } a \rightarrow \text{Pair } b \\
\text{pair } f (x, y) & \quad = (f x, f y)
\end{aligned}$$

Unlike the case of *List*, the naturality of the argument f in $\text{hfold } f$ is crucial to the definition. The recursive occurrence of $\text{hfold } f$ in its defining expression is applied to elements of type $\text{Nest } (\text{Pair } a)$; consequently, f must have type

$$f \quad \text{:: } \text{Base } (\text{Pair}^k b) (n (\text{Pair}^{k+1} b)) \rightarrow n (\text{Pair}^k b)$$

for all $k \geq 0$. The type signature of f in the definition of hfold generalises this requirement by making f fully polymorphic.

Note that in the definition of *Nest*, the parameters of the recursive use do not themselves contain a use of *Nest*. We call such nested datatypes *linear*; almost all of the known practical examples of nested datatypes belong to this class. For an exception, see Section 5 of [BP99]. Our final example is a non-linear variant of *Nest*.

Example 3.3 The type functor *Host* may be defined as the least fixed point of

$$\text{HostF } X \quad = \quad \mathbf{1} + \text{Id} \times (X \cdot (\text{Id} \times X)),$$

The Haskell declaration is:

```

data Base a b    = Nil | Cons (a × b)
type HostF x a  = Base a (x (a × x a))
newtype Host a = In (HostF Host a)

```

This datatype is equivalent to the flattened version

```

data Host a    = Nil | Cons (a × Host (a × Host a))

```

In effect, values of type $\text{Host } a$ correspond to lists in which the first element is of type a , the second element is of type $a \times \text{Host } a$, the third element is of type

$$(a \times \text{Host } a) \times \text{Host } (a \times \text{Host } a),$$

and so on. We know of no practical use for this datatype.

The fold and map functions may be implemented by

$$\begin{aligned}
\text{hfold} & \quad \text{:: } (\forall a. \text{Base } a (n (a \times n a)) \rightarrow n a) \rightarrow \text{Host } b \rightarrow n b \\
\text{hfold } f & \quad = f \cdot \text{base } \text{id} (\text{hfold } f \cdot \text{host } (\text{id} \times \text{hfold } f)) \cdot \text{out} \\
\text{host} & \quad \text{:: } (a \rightarrow b) \rightarrow \text{Host } a \rightarrow \text{Host } b \\
\text{host } f & \quad = \text{In} \cdot \text{base } f (\text{host } (f \times \text{host } f)) \cdot \text{out} \\
\times & \quad \text{:: } (a \rightarrow c) \rightarrow (b \rightarrow d) \rightarrow a \times b \rightarrow c \times d \\
(f \times g) (x, y) & \quad = (f x, g y)
\end{aligned}$$

The functions *host* and \times implement the functorial actions of *Host* and \times , respectively.

4. Generalised folds

A simple fold takes as argument a natural transformation $f :: F N \rightarrow N$ and returns a natural transformation of type $T \rightarrow N$. Its definition has the form

$$hfold f \cdot \alpha = \Psi (hfold f)$$

where the form $\Psi :: (T \rightarrow N) \rightarrow F T \rightarrow N$ is defined by $\Psi h = f \cdot F h$.

Greater flexibility is achieved if we generalise the fold to return a natural transformation of type $T \cdot M \rightarrow N$. For example, both M and N could be taken to be the constant functor that returns the type of integers. Then we would have a way of reducing a T -structure of integers to a single integer. The generalised fold will take as arguments a series of natural transformations f, g_1, \dots, g_m , to be introduced below. A generalised fold

$$gfold f g_1 \cdots g_m :: T \cdot M \rightarrow N$$

will be defined by an equation of the form

$$gfold f g_1 \cdots g_m \cdot \alpha = \Psi (gfold f g_1 \cdots g_m)$$

for some form $\Psi :: (T \cdot M \rightarrow N) \rightarrow F T \cdot M \rightarrow N$. To develop this operator Ψ , observe that every polynomial hofunctor F can be expressed in the form

$$F X = B \cdot \langle Id, X \cdot F_1 X, \dots, X \cdot F_n X \rangle$$

for some $n \geq 0$ (if $n = 0$, then $F X = B \cdot Id = B$). The subsidiary hofunctors F_i may be decomposed similarly⁵.

For example, in the case of *List*, *Nest*, or *Host*, the functor B is given by

$$B X Y = 1 + X \times Y.$$

Now for $h :: T \cdot M \rightarrow N$, we require (expanding $F T \cdot M$):

$$\Psi h :: B \cdot \langle M, T \cdot F_1 T \cdot M, \dots, T \cdot F_n T \cdot M \rangle \rightarrow N.$$

The component types $T \cdot F_i T \cdot M$ are not suitable for arguments to h , so we introduce auxiliary forms $\Phi[F_i] h :: F_i T \cdot M \rightarrow M \cdot F_i N$, to be defined below. Then we have

$$h \cdot T (\Phi[F_i] h) :: T \cdot F_i T \cdot M \rightarrow N \cdot F_i N$$

and we can define

$$\Psi h = f \cdot B id (h \cdot T (\Phi[F_1] h)) \cdots (h \cdot T (\Phi[F_n] h))$$

for some natural transformation

$$f :: B \cdot \langle M, N \cdot F_1 N, \dots, N \cdot F_n N \rangle \rightarrow N$$

which must be supplied as a parameter to *gfold*.

It remains to define $\Phi[F]h$, which we do inductively on the form of F . For $F X = B \cdot \langle Id, X \cdot F_1 X, \dots, X \cdot F_n X \rangle$, we define

$$\Phi[F]h = g \cdot B id (h \cdot T (\Phi[F_1]h)) \cdots (h \cdot T (\Phi[F_n]h))$$

for some natural transformation

$$g :: B \cdot \langle M, N \cdot F_1 N, \dots, N \cdot F_n N \rangle \rightarrow M \cdot F N.$$

⁵ There may be several such decompositions, each giving rise to a different definition of *gfold*.

Each natural transformation g arising in this way is provided as an extra parameter to $gfold$.

Note that in the special case $M = Id$, we can set each $g = id$, obtaining the simple fold as a special case:

$$hfold f = gfold f id \cdot \dots \cdot id.$$

Another special case arises when both M and N are constant functors; the associated instance of $gfold$ is called a *reduction*.

Let us now consider the detailed instantiations for each of our three example datatypes.

Example 4.1 For *List* we have $List = F List$, where

$$\begin{aligned} F X &= Base \cdot \langle Id, X \cdot F_1 X \rangle \\ F_1 X &= Id. \end{aligned}$$

The generalised fold therefore takes the following form in Haskell:

$$\begin{aligned} gfold &:: (\forall a. Base (m a) (n a) \rightarrow n a) \rightarrow \\ &(\forall a. m a \rightarrow m a) \rightarrow \\ &List (m b) \rightarrow n b \\ gfold f g &= f \cdot base id (gfold f g \cdot list g) \cdot out \end{aligned}$$

Of the two additional ingredients, M and g , it is the presence of the functor M that is the crucial one; the second argument g is useful but its effect can be achieved by other means. The point about M is that it provides greater freedom of manoeuvre.

As noted above, we can take both M and N to be constant functors, delivering, say, the types x and y respectively. Then the associated instance of $gfold$ is a reduction of type

$$reduce :: (Base x y \rightarrow y) \rightarrow (x \rightarrow x) \rightarrow List x \rightarrow y.$$

Viewing $f :: Base x y \rightarrow y$ as a pair $e :: y$ and $\oplus :: x \times y \rightarrow y$, the effect of $reduce f g$ on a list $[a_0, \dots, a_n]$ is to produce the value

$$a_0 \oplus (g a_1 \oplus (g^2 a_2 \oplus \dots (g^n a_n \oplus e) \dots)).$$

The regular *fold* is the special case of a list reduction in which $g = id$. List reductions are an important component of the circuit design language Ruby [JS93]; see also [BdM97].

This analysis of reductions suggests that a generalised fold over lists can be decomposed into a simple fold after a *triangle*. The triangle with respect to a function g takes $[a_0, a_1, \dots, a_n]$ to $[a_0, g a_1, \dots, g^n a_n]$. This function can be expressed as a fold:

$$\begin{aligned} triangle &:: (a \rightarrow a) \rightarrow List a \rightarrow List a \\ triangle g &= fold (In \cdot base id (list g)). \end{aligned}$$

The claim is that

$$gfold f g = fold f \cdot triangle g.$$

This result will be proved in the following section. In fact, for regular datatypes, we can always decompose a generalised fold into the composition of two regular folds in this way.

Example 4.2 For $Nest$ we have $Nest = F\ Nest$, where

$$\begin{aligned} F\ X &= Base \cdot \langle Id, X \cdot F_1\ X \rangle \\ F_1\ X &= Pair. \end{aligned}$$

The generalised fold therefore has the following form in Haskell:

$$\begin{aligned} gfold &:: (\forall a. Base\ (m\ a)\ (n\ (Pair\ a)) \rightarrow n\ a) \rightarrow \\ &(\forall a. Pair\ (m\ a) \rightarrow m\ (Pair\ a)) \rightarrow \\ &Nest\ (m\ b) \rightarrow n\ b \\ gfold\ f\ g &= f \cdot base\ id\ (gfold\ f\ g \cdot nest\ g) \cdot out \end{aligned}$$

If we take M and N to be constant functors, delivering types x and y respectively, then we obtain an instance of $gfold$, again called a reduction, with type:

$$reduce :: (Base\ x\ y \rightarrow y) \rightarrow (Pair\ x \rightarrow x) \rightarrow Nest\ x \rightarrow y$$

In particular, we can sum a nest of integers by $sumN = reduce\ sumB\ plus$, where $plus\ (m, n) = m + n$ and

$$\begin{aligned} sumB &:: Base\ Int\ Int \rightarrow Int \\ sumB\ Nil &= 0 \\ sumB\ (Cons\ (m, n)) &= m + n \end{aligned}$$

A nest of integers cannot be summed using a fold over nests (at least, not simply) because $sumB$ is not a natural transformation of the right type. The summation function can be defined using a generalised fold because every function between two given types lifts to a natural transformation between constant functors.

Example 4.3 For $Host$ we have $Host = F\ Host$, where

$$\begin{aligned} F\ X &= Base \cdot \langle Id, X \cdot F_1\ X \rangle \\ F_1\ X &= \times \cdot \langle Id, X \cdot F_2\ X \rangle \\ F_2\ X &= Id. \end{aligned}$$

The generalised fold therefore has the following form in Haskell:

$$\begin{aligned} gfold &:: (\forall a. Base\ (m\ a)\ (n\ (a \times n\ a)) \rightarrow n\ a) \rightarrow \\ &(\forall a. m\ a \times n\ a \rightarrow m\ (a \times n\ a)) \rightarrow \\ &(\forall a. m\ a \rightarrow m\ a) \rightarrow \\ &Host\ (m\ b) \rightarrow n\ b \\ gfold\ f\ g_1\ g_2 &= f \cdot base\ id\ (gfold\ f\ g_1\ g_2 \cdot \\ &host\ (g_1 \cdot id \times (gfold\ f\ g_1\ g_2 \cdot host\ g_2))) \cdot out \end{aligned}$$

Taking M and N to be constant functors, delivering types x and y respectively, we can define a reduction over $Host$ as an instance of $gfold$ with type

$$reduce :: (Base\ x\ y \rightarrow y) \rightarrow (x \times y \rightarrow x) \rightarrow (x \rightarrow x) \rightarrow Host\ x \rightarrow y$$

In particular, we can sum a host of integers by $sumH = reduce\ sumB\ plus\ id$, using the functions $sumB$ and $plus$ defined above.

4.1. Remarks

As we have defined them, generalised folds can be criticised from two opposing points of view. One can argue that our folds are *too* general in that arguments

of type $M \rightarrow M$ could be dropped, since it is not clear what they contribute. On the other hand, one can also argue that our folds are too *specific*. For example, it is quite possible to assign the following type signature to the above definition of *gfold* for *Host*:

$$\begin{aligned} \mathit{gfold} \quad &:: (\forall a. \mathit{Base} (m a) (n (p a)) \rightarrow n a) \rightarrow \\ &(\forall a. m a \times n (q a) \rightarrow m (p a)) \rightarrow \\ &(\forall a. m a \rightarrow m (q a)) \rightarrow \\ &\mathit{Host} (m b) \rightarrow n b \end{aligned}$$

The idea is that $\mathit{gfold} f g h :: \mathit{Host} (m b) \rightarrow n b$ is a natural transformation, and different instances may be instantiated at different types, reflected by the type constructors q and p . The previous type of *gfold* is then the instance in which $q a = a$ and $p a = a \times n a$.

Even more generally, we can recall that the generalised fold is defined as the unique solution of an equation

$$x \cdot \alpha = \Psi x.$$

We can use any form $\Psi :: (T \cdot M \rightarrow N) \rightarrow (F T \cdot M \rightarrow N)$ for which this equation has a unique solution. As we shall see in Section 6, it suffices that Ψ is defined in such a way that any functor may be used in place of T . In particular, Ψ may not use α or α^{-1} to build or dismantle values of type T . We write

$$\Psi \quad :: \quad \forall X. (X \cdot M \rightarrow N) \rightarrow (F X \cdot M \rightarrow N).$$

Thus Ψ is a natural transformation between two contravariant functors.

As always, the more general forms have more degrees of freedom, but consequently offer less guidance to the program designer. We have found the version described in this section, with associated fusion laws presented below, to be a reasonable compromise in practice [BP99]. However, more experience is needed to determine which of the formulations turn out to be the most useful in practical program construction.

5. Fusion laws

The fusion laws for a nested datatype T come in two flavours: *fold-fusion* and *map-fusion*. The former provides conditions under which

$$k \cdot \mathit{gfold} f g_1 \cdots g_m = \mathit{gfold} f' g'_1 \cdots g'_m.$$

The map-fusion law is similar, providing conditions under which

$$\mathit{gfold} f g_1 \cdots g_m \cdot T k = \mathit{gfold} f' g'_1 \cdots g'_m.$$

In the standard theory of regular datatypes, map-fusion is a special case of fold-fusion because the functorial action of regular datatypes can be defined as a fold. This is not the case with the higher-order semantics of Section 3.

In deriving the fusion laws, we shall assume that $\mathit{gfold} f' g'_1 \cdots g'_m$ is the unique function h' satisfying its defining equation. This assumption will be justified in Section 6.

As in the previous section, we suppose $\alpha :: F T \rightarrow T$, where

$$F X = B \cdot \langle \mathit{Id}, X \cdot F_1 X, \dots, X \cdot F_n X \rangle.$$

To reduce clutter, we show the calculations only for the case $n = 1$; the generalisation is straightforward.

The map-fusion law is dealt with first.

5.1. Map-fusion laws

Abbreviating $gfold f g_1 \cdots g_m$ to h , and $gfold f' g'_1 \cdots g'_m$ to h' , our aim is to give conditions under which $h \cdot T k = h'$, where $k :: M' \rightarrow M$. As noted above, we assume that h' is the unique function of type $T \cdot M' \rightarrow N$ satisfying

$$h' \cdot \alpha = \Psi' h',$$

where Ψ' is obtained by replacing each f and g_i in Ψ with the corresponding f' or g'_i . Hence it is sufficient to give conditions under which

$$h \cdot T k \cdot \alpha = \Psi' (h \cdot T k).$$

The general scheme is given by the following calculation:

$$\begin{aligned} & h \cdot T k \cdot \alpha \\ = & \{\text{naturality of } \alpha\} \\ & h \cdot \alpha \cdot B k (T (F_1 T k)) \\ = & \{\text{definition of } h\} \\ & f \cdot B id (h \cdot T (\Phi[F_1]h)) \cdot B k (T (F_1 T k)) \\ = & \{\text{functoriality of } B\} \\ & f \cdot B k (h \cdot T (\Phi[F_1]h) \cdot T (F_1 T k)) \\ = & \{\text{functoriality of } T\} \\ & f \cdot B k (h \cdot T (\Phi[F_1]h \cdot F_1 T k)) \\ = & \{\text{assume } \Phi[F_1]h \cdot F_1 T k = k \cdot \Phi'[F_1](h \cdot T k)\} \\ & f \cdot B k (h \cdot T (k \cdot \Phi'[F_1](h \cdot T k))) \\ = & \{\text{functoriality of } T\} \\ & f \cdot B k (h \cdot T k \cdot T(\Phi'[F_1](h \cdot T k))) \\ = & \{\text{functoriality of } B\} \\ & f \cdot B k id \cdot B id (h \cdot T k \cdot T(\Phi'[F_1](h \cdot T k))) \\ = & \{\text{assume } f' = f \cdot B k id\} \\ & f' \cdot B id (h \cdot T k \cdot T(\Phi'[F_1](h \cdot T k))) \\ = & \{\text{definition}\} \\ & \Psi' (h \cdot T k). \end{aligned}$$

The first assumption in this calculation is the condition

$$\Phi[F_1]h \cdot F_1 T k = k \cdot \Phi'[F_1](h \cdot T k),$$

which can be simplified using the inductive definitions of Φ and Φ' with a similar calculation. Eventually we will arrive at subforms that are independent of h , having generated a series of side-conditions. We will spell out the details for our three example datatypes.

Example 5.1 For *List*, we have $\Phi[F_1]h = g :: M \rightarrow M$ and $F_1 \text{List} = \text{Id}$. Since $\Phi'[F_1](h \cdot \text{list } k) = g' :: M' \rightarrow M'$, we obtain

MAP-FUSION LAW FOR *List*.

$$\text{gfold } f \cdot g \cdot \text{list } k = \text{gfold } (f \cdot \text{base } k \text{ id}) \ g' \iff g \cdot k = k \cdot g'.$$

Taking the special case $g = \text{id}$ and $g' = \text{id}$, we have

$$\text{fold } f \cdot \text{list } k = \text{fold } (f \cdot \text{base } k \text{ id}).$$

Example 5.2 For *Nest*, we have $\Phi[F_1]h = g :: \text{Pair} \cdot M \rightarrow M \cdot \text{Pair}$ and $F_1 \text{Nest} = \text{Pair}$. Since $\Phi'[F_1](h \cdot \text{nest } k) = g' :: \text{Pair} \cdot M' \rightarrow M' \cdot \text{Pair}$, we obtain

MAP-FUSION LAW FOR *Nest*.

$$\text{gfold } f \cdot g \cdot \text{nest } k = \text{gfold } (f \cdot \text{base } k \text{ id}) \ g' \iff g \cdot \text{pair } k = k \cdot g'.$$

Example 5.3 For *Host*, we have

$$\begin{aligned} \Phi[F_1]h &= g_1 \cdot \text{id} \times (h \cdot \text{host } g_2) \\ \Phi'[F_1](h \cdot \text{host } k) &= g'_1 \cdot \text{id} \times (h \cdot \text{host } k \cdot \text{host } g'_2), \end{aligned}$$

where $g_1 :: M \times N \rightarrow M \cdot (\text{Id} \times N)$ and $g_2 :: M \rightarrow M$. The types of g'_1 and g'_2 are similar, except that M is replaced by M' . Since $F_1 \text{Host } k = k \times \text{host } k$, the condition for map-fusion takes the form

$$g_1 \cdot \text{id} \times (h \cdot \text{host } g_2) \cdot k \times \text{host } k = k \cdot g'_1 \cdot \text{id} \times (h \cdot \text{host } k \cdot \text{host } g'_2)$$

or equivalently

$$g_1 \cdot k \times \text{id} \cdot \text{id} \times (h \cdot \text{host } (g_2 \cdot k)) = k \cdot g'_1 \cdot \text{id} \times (h \cdot \text{host } (k \cdot g'_2)).$$

This equation follows from $g_1 \cdot k \times \text{id} = k \cdot g'_1$ and $g_2 \cdot k = k \cdot g'_2$. Hence

MAP-FUSION LAW FOR *Host*.

$$\begin{aligned} \text{gfold } f \ g_1 \ g_2 \cdot \text{host } k &= \text{gfold } (f \cdot \text{base } k \text{ id}) \ g'_1 \ g'_2 \\ &\iff g_1 \cdot k \times \text{id} = k \cdot g'_1 \wedge g_2 \cdot k = k \cdot g'_2. \end{aligned}$$

5.2. Fold-fusion laws

Again abbreviating $\text{gfold } f \ g_1 \cdots g_m$ to h and $\text{gfold } f' \ g'_1 \cdots g'_m$ to h' , our aim is to give conditions under which $k \cdot h_{M'} = h'$, where $k :: N \cdot M' \rightarrow N'$. As before, it is sufficient to give conditions under which

$$k \cdot h \cdot \alpha = \Psi'(k \cdot h).$$

The general scheme is given by the following calculation:

$$\begin{aligned} &k \cdot h \cdot \alpha \\ &= \quad \{\text{definition of } h\} \\ &\quad k \cdot f \cdot B \text{id } (h \cdot T(\Phi[F_1] h)) \\ &= \quad \{\text{assume } k \cdot f = f' \cdot B \text{id } (k \cdot N \ p_1) \text{ for some } p_1\} \end{aligned}$$

$$\begin{aligned}
& f' \cdot B \text{ id } (k \cdot N p_1) \cdot B \text{ id } (h \cdot T (\Phi[F_1] h)) \\
= & \quad \{\text{functoriality of } B\} \\
& f' \cdot B \text{ id } (k \cdot N p_1 \cdot h \cdot T (\Phi[F_1] h)) \\
= & \quad \{\text{naturality of } h :: T \cdot M \rightarrow N, \text{ and functoriality of } T\} \\
& f' \cdot B \text{ id } (k \cdot h \cdot T (M p_1 \cdot \Phi[F_1] h)) \\
= & \quad \{\text{assume } M p_1 \cdot \Phi[F_1] h = \Phi'[F_1](k \cdot h)\} \\
& f' \cdot B \text{ id } (k \cdot h \cdot T (\Phi'[F_1](k \cdot h))) \\
= & \quad \{\text{definition}\} \\
& \Psi' (k \cdot h).
\end{aligned}$$

The second assumption in this calculation, namely,

$$M p_1 \cdot \Phi[F_1] h = \Phi'[F_1](k \cdot h),$$

where $p_1 :: F_1 N \cdot M' \rightarrow M' \cdot F_1 N'$, is elaborated to simpler assumptions, using the inductive definitions of Φ and Φ' . We will spell out the details for our three examples.

Example 5.4 For *List* we have $\Phi[F_1] h = g :: M \rightarrow M$. Hence

FOLD-FUSION LAW FOR *List*. Given the typings $f :: \text{Base} \cdot \langle M, N \cdot \text{Pair} \rangle \rightarrow N$, $g :: M \rightarrow M$ and $p :: M' \rightarrow M'$, we have

$$k \cdot \text{gfold } f \ g = \text{gfold } f' \ (M p \cdot g) \iff k \cdot f = f' \cdot \text{base id } (k \cdot N p).$$

For example, consider the composition

$$\text{gfold } f \ \text{id} \cdot \text{gfold } \alpha \ g,$$

where $f :: \text{Base} \cdot \langle M, N \cdot \text{Pair} \rangle \rightarrow N$ and $g :: M \rightarrow M$. The function $\text{gfold } f \ \text{id} :: \text{List} \cdot M \rightarrow N$ is a fold, while the function $\text{gfold } \alpha \ g :: \text{List} \cdot M \rightarrow \text{List} \cdot M$ is called the *triangle* with respect to g . We have

$$\text{gfold } f \ \text{id} \cdot \alpha = f \cdot \text{base id } (\text{gfold } f \ \text{id} \cdot \text{list id}).$$

Hence an appeal to fold-fusion gives

$$\text{gfold } f \ \text{id} \cdot \text{gfold } \alpha \ g = \text{gfold } f \ g.$$

Thus every generalised fold on lists (and indeed any regular type) can be factored as a regular fold after a triangle. In fact, every triangle can be expressed as a regular fold:

$$\begin{aligned}
& \text{gfold } \alpha \ g \cdot \alpha \\
= & \quad \{\text{definition}\} \\
& \alpha \cdot \text{base id } (\text{gfold } \alpha \ g \cdot \text{list } g) \\
= & \quad \{\text{map-fusion}\} \\
& \alpha \cdot \text{base id } (\text{gfold } (\alpha \cdot \text{base } g \ \text{id}) \ g) \\
= & \quad \{\text{fold-fusion (backwards)}\} \\
& \alpha \cdot \text{base id } (g \cdot \text{gfold } \alpha \ g) \\
= & \quad \{\text{functoriality of } \text{base}\} \\
& \alpha \cdot \text{base id } g \cdot \text{base id } (\text{gfold } \alpha \ g).
\end{aligned}$$

Hence $gfold \alpha g = fold (\alpha \cdot base g id)$, and $gfold f g$ is a composition of regular folds. Note that we have not used any property of $base$ except functoriality, so this factorisation may be applied to the generalised fold on any regular type.

Example 5.5 For $Nest$, we have $\Phi[F_1]h = g :: Pair \cdot M \rightarrow M \cdot Pair$. Hence

FOLD-FUSION LAW FOR $Nest$. Given the typings

$$\begin{aligned} f &:: Base \cdot \langle M, N \cdot Pair \rangle \rightarrow N \\ g &:: Pair \cdot M \rightarrow M \cdot Pair \\ p &:: Pair \cdot M' \rightarrow M' \cdot Pair \end{aligned}$$

we have

$$k \cdot gfold f g = gfold f' (M p \cdot g) \Leftarrow k \cdot f = f' \cdot base id (k \cdot N p).$$

The fold-fusion law for $Nest$ takes the same form as for $List$, Though the types of g and p are different. Here is a simple example. The function $listify :: Nest \rightarrow List$ can be defined by

$$listify = gfold (\alpha \cdot base id (concat \cdot list duo)) id$$

where $\alpha :: Base \cdot \langle Id, List \rangle \rightarrow List$ is the initial list algebra, $duo :: Pair \rightarrow List$ converts a pair of values into a list of two values, and $concat :: List \cdot List \rightarrow List$ concatenates a list of lists into one list. Here $M = Id$ and $N = List$.

Furthermore, the function $sumL :: List \cdot \underline{Int} \rightarrow \underline{Int}$, which sums a list of integers, is given by

$$sumL = fold sumB,$$

where $fold$ is the fold function for lists.

Now consider the combination $sumL \cdot listify$. We can fuse the two functions provided

$$sumL \cdot \alpha \cdot base id (concat \cdot list duo) = f' \cdot base id (sumL \cdot list p),$$

for some $p :: Pair \cdot \underline{Int} \rightarrow \underline{Int}$ and $f' :: Base \cdot \langle \underline{Int}, \underline{Int} \rangle \rightarrow \underline{Int}$. For the proviso we can calculate:

$$\begin{aligned} &sumL \cdot \alpha \cdot base id (concat \cdot list duo) \\ = &\quad \{\text{definition of } sumL\} \\ &sumB \cdot base id sumL \cdot base id (concat \cdot list duo) \\ = &\quad \{\text{functoriality of } base\} \\ &sumB \cdot base id (sumL \cdot concat \cdot list duo) \\ = &\quad \{\text{property of } sumL \text{ and functoriality of } list\} \\ &sumB \cdot base id (sumL \cdot list (sumL \cdot duo)) \\ = &\quad \{\text{since } sumL \cdot duo = plus\} \\ &sumB \cdot base id (sumL \cdot list plus). \end{aligned}$$

Hence $sumL \cdot listify = gfold sumB plus$. The expression on the right is just the function $sumN$, given in the previous section, for summing a nest of integers directly.

Example 5.6 For $T = Host$, we have $\Phi[F_1]h = g_1 \cdot id \times (h \cdot T g_2)$, where $g_1 :: M \times N \rightarrow M \cdot (Id \times N)$ and $g_2 :: M \rightarrow M$. For fold-fusion we require $k \cdot f = f' \cdot base\ id\ (k \cdot N\ p)$, where $p :: M' \times (N \cdot M') \rightarrow M' \cdot (Id \times N')$ is required to satisfy the equation

$$M\ p \cdot \Phi[F_1]h = \Phi'[F_1](k \cdot h).$$

This condition expands to

$$M\ p \cdot g_1 \cdot id \times (h \cdot T\ g_2) = g'_1 \cdot id \times (k \cdot h \cdot T\ g'_2).$$

To eliminate the dependence on h , suppose

$$M\ p \cdot g_1 = g'_1 \cdot id \times (k \cdot N\ q)$$

for some $q :: M' \rightarrow M'$. Then we can argue:

$$\begin{aligned} & M\ p \cdot g_1 \cdot id \times (h \cdot T\ g_2) \\ = & \quad \{\text{assumption, and functoriality of } \times\} \\ & g'_1 \cdot id \times (k \cdot N\ q \cdot h \cdot T\ g_2) \\ = & \quad \{\text{naturality of } h :: T \cdot M \rightarrow N \text{ and functoriality of } T\} \\ & g'_1 \cdot id \times (k \cdot h \cdot T(M\ q \cdot g_2)) \end{aligned}$$

Hence the fold-fusion law can be put in the form

FOLD-FUSION LAW FOR *Host*. Given the typings

$$\begin{aligned} f & :: Base \cdot \langle M, N \cdot (Id \times N) \rangle \rightarrow N \\ g_1 & :: M \times N \rightarrow M \cdot (Id \times N) \\ g_2 & :: M \rightarrow M \\ q & :: M' \rightarrow M', \end{aligned}$$

we have

$$k \cdot gfold\ f\ g_1\ g_2 = gfold\ f'\ g'_1\ (M\ q \cdot g_2)$$

provided that, for some $p :: M' \times (N \cdot M') \rightarrow M' \cdot (Id \times N')$,

$$k \cdot f = f' \cdot base\ id\ (k \cdot N\ p) \quad \text{and} \quad M\ p \cdot g_1 = g'_1 \cdot id \times (k \cdot N\ q).$$

6. Uniqueness of generalised folds

Our aim in this section is to show that generalised folds are the unique solutions of their defining equations, which have the form

$$gfold\ f\ g_1 \cdots g_m \cdot \alpha_M = \Psi(gfold\ f\ g_1 \cdots g_m)$$

for some form $\Psi :: \forall X. (X \cdot M \rightarrow N) \rightarrow (F\ X \cdot M \rightarrow N)$. It turns out that it is easier to solve a more general problem. Suppose we are given functors $F :: \mathbf{C} \rightarrow \mathbf{C}$ and $L :: \mathbf{C} \rightarrow \mathbf{D}$ for categories \mathbf{C} and \mathbf{D} , for which F has least fixed point $\alpha :: F\ T \rightarrow T$, and a natural transformation $\Psi :: \forall A. (L\ A \rightarrow B) \rightarrow (L(F\ A) \rightarrow B)$. We seek conditions on F and L to ensure there is a unique $x :: L\ T \rightarrow B$ such that

$$x \cdot L\ \alpha = \Psi\ x. \tag{1}$$

Note that the naturality of Ψ implies that

$$\Psi(f \cdot Lx) = \Psi f \cdot L(Fx). \quad (2)$$

Some instances of this general scheme are:

- The *gfold* over nested types is the special case where F and L are higher order functors, with $L = (\cdot M)$. The action of this functor on natural transformations, like α , yields instances α_M .
- In the special case $L = Id$, equation (2) implies

$$\Psi x = \Psi id \cdot Fx$$

so that Ψ is determined by $\Psi id :: FB \rightarrow B$ and the generalised fold reduces to an ordinary fold.

- Further instances are defining equations of the form

$$x \cdot L\alpha = f \cdot Gx \cdot g$$

for functions $f :: GB \rightarrow B$ and $g :: L \cdot F \rightarrow G \cdot L$. For example, *zip*-like functions [FSZ94] are of this form, with $L = \times$ (a binary functor) and $g :: F_1 \times F_2 \rightarrow G \cdot \times$.

Furthermore, this uniqueness property leads to a general fusion law, stating that for any natural transformation $\Phi :: \forall A. (LA \rightarrow B) \rightarrow (L'A \rightarrow B)$,

$$\Phi(\text{gfold } \Psi) = \text{gfold } \Psi' \iff \Phi \cdot \Psi = \Psi' \cdot \Phi.$$

All the previous fusion laws are instances of this scheme.

We present two solutions to this problem. The first depends on the details of the colimit construction of the fixed point, but works for a larger class of functors L than the second, more abstract, approach. For example, the first approach works for $L = \times$, while the second does not.

6.1. First method: colimits

Theorem 1. Suppose the functors F and L preserve colimits of chains (all those we can define in Haskell do), and L preserves initiality. Then equation (1) has a unique solution.

Proof. Suppose the colimit used to construct the fixed point T of F consists of functions $e_n :: F^n 0 \rightarrow T$. By construction, the isomorphism $\alpha :: FT \cong T$ satisfies

$$\alpha \cdot F e_n = e_{n+1} \quad (3)$$

for each n . Since L preserves initiality, there is a unique arrow $z :: L0 \rightarrow B$. We shall show that

$$x \cdot L\alpha = \Psi x \iff \forall n. x \cdot L e_n = \Psi^n z.$$

Since the arrows $L e_n$ comprise a colimit, this establishes the existence and uniqueness of x .

First, we prove

$$\forall n. x \cdot L e_n = \Psi^n z \iff x \cdot L\alpha = \Psi x$$

by induction on n . The base case is immediate from the initiality of $L0$. The induction step is

$$\begin{aligned}
& x \cdot L e_{n+1} \\
= & \quad \{\text{equation (3)}\} \\
& x \cdot L(\alpha \cdot F e_n) \\
= & \quad \{\text{functor}\} \\
& x \cdot L\alpha \cdot L(F e_n) \\
= & \quad \{\text{hypothesis}\} \\
& \Psi x \cdot L(F e_n) \\
= & \quad \{\text{equation (2)}\} \\
& \Psi(x \cdot L e_n) \\
= & \quad \{\text{induction hypothesis}\} \\
& \Psi(\Psi^n z)
\end{aligned}$$

To establish the reverse implication, we reason

$$\begin{aligned}
& x \cdot L\alpha = \Psi x \\
\equiv & \quad \{\text{isomorphism}\} \\
& x = \Psi x \cdot L\alpha^{-1} \\
\Leftarrow & \quad \{\text{colimit}\} \\
& \forall n. \Psi x \cdot L\alpha^{-1} \cdot L e_n = \Psi^n z \\
\Leftarrow & \quad \{\text{by cases of } n \text{ (see below)}\} \\
& \forall n. x \cdot L e_n = \Psi^n z
\end{aligned}$$

The case $n = 0$ is immediate from the initiality of $L0$. The case $n = m + 1$ is

$$\begin{aligned}
& \Psi x \cdot L\alpha^{-1} \cdot L e_{m+1} \\
= & \quad \{\text{equation (3)}\} \\
& \Psi x \cdot L\alpha^{-1} \cdot L\alpha \cdot L(F e_m) \\
= & \quad \{\text{functor, isomorphism}\} \\
& \Psi x \cdot L(F e_m) \\
= & \quad \{\text{equation (2)}\} \\
& \Psi(x \cdot L e_m) \\
= & \quad \{\text{hypothesis}\} \\
& \Psi(\Psi^m z)
\end{aligned}$$

completing the proof. \square

6.2. Second method: adjoints

Another approach is to note that we seek a function of type $L T \rightarrow B$, whereas a fold can supply a function of type $T \rightarrow B'$. Recall that an adjunction between L and a functor R defines an isomorphism

$$\varphi \quad :: \quad \forall A, B. (L A \rightarrow B) \cong (A \rightarrow R B). \quad (4)$$

Setting $A = T$, this yields a correspondence between generalised folds and ordinary folds.

Theorem 2. If L has a right adjoint, then equation (1) has a unique solution.

Proof. Let φ denote the isomorphism defined by the adjunction, as above. We will be applying it to both sides of equation (1). We calculate

$$\begin{aligned}
& \varphi(\Psi x) \\
= & \quad \{\text{isomorphism}\} \\
& \varphi(\Psi(\varphi^{-1}(\varphi x))) \\
= & \quad \{\text{naturality of } \varphi^{-1}, \text{ introducing the counit } \varepsilon = \varphi^{-1} id\} \\
& \varphi(\Psi(\varepsilon \cdot L(\varphi x))) \\
= & \quad \{\text{naturality of } \Psi\} \\
& \varphi(\Psi \varepsilon \cdot L(F(\varphi x))) \\
= & \quad \{\text{naturality of } \varphi\} \\
& \varphi(\Psi \varepsilon) \cdot F(\varphi x)
\end{aligned}$$

Thus we have

$$\begin{aligned}
& x \cdot L\alpha = \Psi x \\
\equiv & \quad \{\text{isomorphism}\} \\
& \varphi(x \cdot L\alpha) = \varphi(\Psi x) \\
\equiv & \quad \{\text{naturality of } \varphi\} \\
& \varphi x \cdot \alpha = \varphi(\Psi x) \\
\equiv & \quad \{\text{above calculation}\} \\
& \varphi x \cdot \alpha = \varphi(\Psi \varepsilon) \cdot F(\varphi x) \\
\equiv & \quad \{\text{fold uniqueness}\} \\
& \varphi x = \text{fold}(\varphi(\Psi \varepsilon)) \\
\equiv & \quad \{\text{isomorphism}\} \\
& x = \varphi^{-1}(\text{fold}(\varphi(\Psi \varepsilon)))
\end{aligned}$$

□

Note that this version is strictly weaker than the lower-level approach using colimits: any functor with a right adjoint necessarily preserves all colimits, while \times is an example of a functor that preserves initiality and colimits of chains, but not coproducts.

To apply this theorem to generalised folds over nested datatypes, it remains to construct a right adjoint to $(\cdot M)$. This is a well-studied problem in category theory, where it is known as the right Kan extension [Mac71, X.3], and is known to exist if M is defined on a small complete category. It is also possible to eliminate the size condition by restricting the class of functors M , say to nested functors.

We can define this adjunction in Haskell by introducing a *continuation* type $Cont\ M\ N$, and constructing an isomorphism pair between $T \cdot M \rightarrow N$ and $T \rightarrow Cont\ M\ N$. The functor $Cont$ is declared as

```
newtype Cont m n a = MkCont (\forall b. (a -> m b) -> n b)
```

Thus, *Cont* wraps a higher-order polymorphic function. For fixed M and N , the functorial action of *Cont* is defined by

$$\begin{aligned} \text{cont} & \quad :: (a \rightarrow b) \rightarrow \text{Cont } m \ n \ a \rightarrow \text{Cont } m \ n \ b \\ \text{cont } f \ (\text{MkCont } g) & = \text{MkCont } (\lambda k. g(k \cdot f)) \end{aligned}$$

One half of the isomorphism is the function

$$\begin{aligned} \text{toCont} & \quad :: \text{Functor } t \Rightarrow (\forall a. t \ (m \ a) \rightarrow n \ a) \rightarrow t \ b \rightarrow \text{Cont } m \ n \ b \\ \text{toCont } f \ x & = \text{MkCont } (\lambda k. f \ (\text{map}_t \ k \ x)) \end{aligned}$$

The function map_t is the functorial action of t . The subscript is added for clarity; it is not used in Haskell.

The other half of the isomorphism pair is

$$\begin{aligned} \text{fromCont} & \quad :: \text{Functor } t \Rightarrow (\forall a. t \ a \rightarrow \text{Cont } m \ n \ a) \rightarrow t \ (m \ b) \rightarrow n \ b \\ \text{fromCont } f & = \text{applyto } \text{id} \cdot f \end{aligned}$$

where the function $\text{applyto } k$ applies a continuation to k :

$$\begin{aligned} \text{applyto} & \quad :: (a \rightarrow m \ b) \rightarrow \text{Cont } m \ n \ a \rightarrow n \ b \\ \text{applyto } k \ f & = \text{unCont } f \ k \\ \text{unCont} & \quad :: \text{Cont } m \ n \ a \rightarrow (a \rightarrow m \ b) \rightarrow n \ b \\ \text{unCont } (\text{MkCont } x) & = x \end{aligned}$$

Then toCont and fromCont is an isomorphism pair, meaning that

$$\begin{aligned} \text{toCont} \cdot \text{fromCont} & = \text{id} \\ \text{fromCont} \cdot \text{toCont} & = \text{id}. \end{aligned}$$

The first identity is verified by the following calculation:

$$\begin{aligned} & \text{toCont } (\text{fromCont } f) \ x \\ = & \quad \{\text{definitions of } \text{toCont} \text{ and } \text{fromCont}\} \\ & \text{MkCont } (\lambda k. \text{applyto } \text{id} \ (f \ (\text{map}_t \ k \ x))) \\ = & \quad \{\text{definition of } \text{applyto}\} \\ & \text{MkCont } (\lambda k. \text{unCont } (f \ (\text{map}_t \ k \ x)) \ \text{id}) \\ = & \quad \{\text{naturality of } f\} \\ & \text{MkCont } (\lambda k. \text{unCont } (\text{cont } k \ (f \ x)) \ \text{id}) \\ = & \quad \{\text{definition of } \text{cont}\} \\ & \text{MkCont } (\lambda k. (\lambda k'. \text{unCont } (f \ x) \ (k' \cdot k)) \ \text{id}) \\ = & \quad \{\text{beta reduction; identity}\} \\ & \text{MkCont } (\lambda k. \text{unCont } (f \ x) \ k) \\ = & \quad \{\text{eta reduction; definition of } \text{unCont}\} \\ & f \ x. \end{aligned}$$

For the second identity, we calculate:

$$\begin{aligned} & \text{fromCont } (\text{toCont } f) \ x \\ = & \quad \{\text{definitions of } \text{fromCont}, \text{applyto}\} \\ & \text{unCont } (\text{toCont } f \ x) \ \text{id} \\ = & \quad \{\text{definitions of } \text{toCont} \text{ and } \text{unCont}\} \end{aligned}$$

$$\begin{aligned}
& (\lambda k. f (map_t k x)) id \\
= & \quad \{\text{beta reduction}\} \\
& f (map_t id x) \\
= & \quad \{\text{functoriality of } map_t\} \\
& f x.
\end{aligned}$$

7. Final remarks

The subject of nested datatypes and what they might provide for the practical programmer is still in its infancy. Most of the published work so far (in particular, [Oka98]) uses nested types only as a conceptual tool; the results are translated into programming terms by embedding nested types in regular ones in a systematic way. One reason that nested types have not been used directly is that, until recently, there has been no language support for defining functions over such types. Recent versions of Haskell permit rank-2 type signatures and, as we have seen, such a move is necessary to implement functions over nested types.

Defining functions over nested datatypes by explicit recursion is complicated, and even more error prone than using a similar style with regular types. The real purpose of concentrating on the fold function for a datatype is that it provides a structured approach to inductive functional programming. Such a concentration has added force when the datatype is nested.

Acknowledgements

Peter Freyd, Andy Pitts, James Worrell and Dominic Hughes directed us to the appropriate category theory for the adjunction in Section 6.2. Anonymous referees made several helpful suggestions.

References

- [BdM97] R. Bird and O. de Moor. *Algebra of Programming*. Prentice Hall, 1997.
- [Bir98] Richard Bird. *Introduction to Functional Programming using Haskell*. Prentice Hall, second edition, 1998.
- [BM98] Richard S. Bird and Lambert Meertens. Nested datatypes. In *Mathematics of Program Construction*, volume 1422 of *Lecture Notes in Computer Science*, pages 52–67. Springer, 1998.
- [BP99] Richard Bird and Ross Paterson. De Bruijn notation as a nested datatype. *Journal of Functional Programming*, 9(1):77–91, January 1999.
- [CL95] R.H. Connelly and F. Lockwood Morris. A generalisation of the trie data structure. *Mathematical Structures in Computer Science*, 5(3):381–418, 1995.
- [FSZ94] Leonidas Fegaras, Tim Sheard, and Tong Zhou. Improving programs which recurse over multiple inductive structures. In *ACM SIGPLAN Workshop on Partial Evaluation and Semantics-Based Program Manipulation*, June 1994.
- [Hag87] T. Hagino. *Category Theoretic Approach to Data Types*. PhD thesis, Laboratory for Foundations of Computer Science, University of Edinburgh, UK, 1987. Technical Report ECS-LFCS-87-38.
- [Hoo97] Paul Hoogendijk. *A Generic Theory of Data Types*. PhD thesis, Eindhoven University of Technology, 1997.
- [Jon98] Mark Jones. A technical summary of the new features in Hugs 1.3c, 1998. unpublished.

- [JS93] Geraint Jones and Mary Sheeran. Designing arithmetic circuits by refinement in Ruby. In R. Bird, C. Morgan, and J. Woodcock, editors, *Mathematics of Program Construction*, volume 669 of *Lecture Notes in Computer Science*, pages 208–232. Springer, 1993.
- [Lam70] J. Lambek. Subequalizers. *Canadian Mathematical Bulletin*, 13:337–349, 1970.
- [MA86] E.G. Manes and M.A. Arbib. *Algebraic Approaches to Program Semantics*. Texts and Monographs in Computing Science. Springer, 1986.
- [Mac71] Saunders Mac Lane. *Categories for the Working Mathematician*. Springer, New York, 1971.
- [Mal90a] Grant R. Malcolm. *Algebraic Data Types and Program Transformation*. PhD thesis, Rijksuniversiteit Groningen, 1990.
- [Mal90b] Grant R. Malcolm. Data structures and program transformation. *Science of Computer Programming*, 14(2-3):255–279, 1990.
- [McC84] N. J. McCracken. The typechecking of programs with implicit type structure. In *Semantics of Data Types*, volume 173 of *Lecture Notes in Computer Science*, pages 301–315. Springer, 1984.
- [Oka98] Chris Okasaki. *Purely Functional Data Structures*. Cambridge University Press, 1998.
- [PL97] S. Peyton Jones and J. Launchbury. Explicit quantification in Haskell, 1997. See: <http://www.dcs.gla.ac.uk/people/personal/simonpj/>.