Manufacturing Datatypes April 1999

RALF HINZE

Institut für Informatik III, Universität Bonn Römerstraße 164, 53117 Bonn, Germany (e-mail: ralf@informatik.uni-bonn.de)

Abstract

This paper describes a general framework for designing purely functional datatypes that automatically satisfy given size or structural constraints. Using the framework we develop implementations of different matrix types (eg square matrices) and implementations of several tree types (eg Braun trees, 2-3 trees). Consider, for instance, representing square $n \times n$ matrices. The usual representation using lists of lists fails to meet the structural constraints: there is no way to ensure that the outer list and the inner lists have the same length. The main idea of our approach is to solve in a first step a related, but simpler problem, namely to generate the multiset of all square numbers. In order to describe this multiset we employ recursion equations involving finite multisets, multiset union, addition and multiplication lifted to multisets. In a second step we mechanically derive datatype definitions from these recursion equations which enforce the 'squareness' constraint. The transformation makes essential use of polymorphic types.

Die ganze Zahl schuf der liebe Gott, alles Übrige ist Menschenwerk. — Leopold Kronecker

1 Introduction

Many information structures are defined by certain size or structural constraints. Take, for instance, the class of perfect leaf trees (Hinze, 1999a): a perfect leaf tree of height 0 is a leaf and a perfect leaf tree of height h + 1 is a node with two children, each of which is of height h. How can we represent perfect leaf trees of arbitrary height such that the structural constraints are enforced? The usual recursive representation of leaf trees is apparently not very helpful since there is no way to ensure that the children of a node have the same height. As another example, consider square $n \times n$ matrices (Okasaki, 1999). How do we represent square matrices such that the matrices are actually square? Again, the standard representation using lists of lists fails to meet the constraints: the outer list and the inner lists have not necessarily the same length. In this paper, we present a framework that allows to design representations of perfect leaf trees, square matrices, and many other information structures that automatically satisfy the given size or structural constraints.

Let us illustrate the main ideas by means of example. As a first example, we will devise a representation of Toeplitz matrices (Cormen et al., 1991) where a Toeplitz matrix is an $n \times n$ matrix (a_{ij}) such that $a_{ij} = a_{i-1,j-1}$ for $1 < i, j \leq n$. Clearly, to represent a Toeplitz matrix of size n + 1 it suffices to store 2 * n + 1 elements. Now, instead of designing a representation from the scratch we first solve a related, but apparently simpler problem, namely, to generate the set of all odd numbers. Actually, we will work with multisets instead of sets for reasons to be explained later. In order to describe multisets of natural numbers we employ systems of recursion equations. The following system, for instance, specifies the multiset of all odd numbers, ie the multiset which contains one occurrence of each odd number.

$$odd = \langle 1 \rangle \uplus \langle 2 \rangle + odd$$

Here, $\langle n \rangle$ denotes the singleton multiset which contains n exactly once, (H) denotes multiset union and (+) is addition lifted to multisets: $A + B = \langle a + b \mid a \leftarrow A; b \leftarrow B \rangle$. We agree upon that (+) binds more tightly than (\oiint) . Now, how can we turn the equation into a sensible datatype definition for Toeplitz matrices? The first thing to note is that we are actually looking for a datatype which is parameterized by the type of matrix elements. Such a type is also known as a *type constructor* or as a *functor*¹. An element of a parameterized type is called a *container*. The equation above has the following counterpart in the world of functors.

$$Odd = Id | (Id \times Id) \times Odd$$

Here, Id is the identity functor given by $Id \ a = a$. Furthermore, (|) and (×) denote disjoint sums and products lifted to functors, ie $(F_1 | F_2) \ a = F_1 \ a | F_2 \ a$ and $(F_1 \times F_2) \ a = F_1 \ a \times F_2 \ a$. Comparing the two equations we see that $\langle 1 \rangle$ corresponds to Id, (\uplus) corresponds to (|), and (+) corresponds to (×). This immediately implies that $Id \times Id$ corresponds to $\langle 1 \rangle + \langle 1 \rangle = \langle 2 \rangle$. The relationship is very tight: the functor corresponding to a multiset M contains, for each member of M, a container of that size. For instance, $Id \times Id$ corresponds to $\langle 1 \rangle + \langle 1 \rangle = \langle 2 \rangle$ as it contains one container of size 2; $Id | Id \times Id$ corresponds to $\langle 1 \rangle \uplus \langle 1 \rangle + \langle 1 \rangle = \langle 1, 2 \rangle$ as it contains one container of size 1 and another one of size 2.

Functor equations are written in a compositional style. To derive a datatype declaration from a functor equation we simply rewrite it into an applicative form—additionally adding constructor names and possibly making cosmetic changes.²

data Toeplitz
$$a = Corner \ a \mid Extend \ a \ a \ (Toeplitz \ a)$$

The left upper corner of a Toeplitz matrix is represented by *Corner* a; *Extend* r c m extends the matrix m by an additional row and an additional column, both of which are represented by elements.

Of course, this is not the only implementation conceivable. Alternatively, we can

¹ Categorically speaking, a functor must satisfy additional conditions, see (Bird & de Moor, 1997). All the type constructors listed in this paper are functors in the category-theoretical sense.

 $^{^{2}}$ Examples are given in the functional language Haskell (Peyton Jones & Hughes, 1998).

define odd in terms of the set of all even numbers.

$$odd = \langle 1 \rangle + even$$
$$even = \langle 0 \rangle \uplus \langle 2 \rangle + even$$

As innocent as this variation may look it has the advantage that the left upper corner can be accessed in constant time as opposed to linear time with the first representation.

Easier still, we may define *odd* in terms of the natural numbers using the fact that each odd number is of the form 1 + n * 2 for some n.

$$\begin{array}{rcl} odd & = & \langle 1 \rangle + nat * \langle 2 \rangle \\ nat & = & \langle 0 \rangle \uplus \langle 1 \rangle + nat \end{array}$$

The first equation makes use of the multiplication operation, which is defined analogous to (+). To which operation on functors does multiplication correspond? We will see that under certain conditions to be spelled out later (*) corresponds to the composition of functors (·) given by $(F_1 \cdot F_2) \ a = F_1 \ (F_2 \ a)$. The functor equations derived from *odd* and *nat* are

$$Odd = Id \times Nat \cdot (Id \times Id)$$

 $Nat = K Unit | Id \times Nat$.

Here, K t denotes the constant functor given by K t a = t and *Unit* is the unit type containing a single element. Unsurprisingly, *Nat* models the ubiquitous datatype of polymorphic lists.

data Toeplitz
$$a$$
 = Toeplitz a (List (a, a))
data List a = Nil | Cons a (List a)

Thus, to store an even number of elements we simply use a list of pairs. This representation has the advantage that the list type can be easily replaced by a more efficient sequence type.

Next, let us apply the technique to design a representation of perfect leaf trees. The related problem is simple: we have to generate the multiset of all powers of 2.

$$power = \langle 1 \rangle \uplus power * \langle 2 \rangle$$

The corresponding functor equation is

$$Power = Id \mid Power \cdot (Id \times Id)$$

from which we can easily derive the following datatype definition.

data Perfect
$$a = Zero \ a \mid Succ \ (Perfect \ (a, a))$$

Thus, a perfect leaf tree of height 0 is a leaf and a perfect leaf tree of height h + 1 is a perfect leaf tree of height h, whose leaves contain pairs of elements. Note that

this definition proceeds bottom-up, whereas the definition given in the beginning proceeds top-down. The type Perfect is an example for a so-called nested datatype (Bird & Meertens, 1998): the recursive call of Perfect on the right-hand side is not a copy of the declared type on the left-hand side, ie the type recursion is nested. It is revealing to have a closer look at the types. The table below illustrates the construction of an element of type Perfect Int (\$ always refers to the expression in the preceding row).

(((1,2),(3,4)),((5,6),(7,8)))	::	(((Int, Int), (Int, Int)), ((Int, Int), (Int, Int)))
$Zero\$::	$Perfect \ (((Int, Int), (Int, Int)), ((Int, Int), (Int, Int)))$
Succ \$::	$Perfect \ ((Int, Int), (Int, Int))$
Succ \$::	Perfect (Int, Int)
Succ \$::	Perfect Int

We start with a nested pair of integers. Note that the type expression has the same size as the value expression. Using the constructor Zero the nested pair is turned into a leaf. Now, each application of Succ halves the size of the type expression. In each case the typechecker ensures that the elements are pairs of the same type.

As the final example, let us tackle the problem of representing square matrices. We soon find that the related problem of generating the multiset of all square numbers is not quite as easy as before. One could be tempted to define square = nat * nat. However, this does not work since the resulting multiset contains products of arbitrary numbers. Incidentally, nat * nat is related to $List \cdot List$, the lists of lists implementation we already depraved. We must somehow arrange that (*) is only applied to singleton multisets. A trick to achieve this is to first rewrite the definition of nat into a tail-recursive form.

nat = nat' (0) $nat' n = n \uplus nat' ((1) + n)$

The definition of nat' closely resembles the function $from :: Int \rightarrow [Int]$ given by $from \ n = n : from \ (n + 1)$, which generates the infinite list of successive integers beginning with n. Now, to obtain square numbers we simple replace n by n * n in the second equation.

square = square' (0) $square' n = n * n \uplus square' ((1) + n)$

Using this trick we are, in fact, able to enumerate the codomain of an arbitrary polynomial. Even more interesting, this trick is applicable to other representations of sequences, as well. But, we are skipping ahead. For now, let us determine the datatypes corresponding to *square* and *square'*. From the functor equations

we can derive the following datatype declarations.

type Matrix a	=	Matrix' Nil a
data Matrix' t a	=	Zero $(t (t a)) \mid Succ (Matrix' (Cons t) a)$
data Nil a	=	Nil
data Cons t a	=	$Cons \ a \ (t \ a)$

The type constructors Nil and Cons t correspond to K Unit and $Id \times f$. As an aside, note that Nil and Cons are obtained by decomposing the List datatype into a base and into a recursive case. Furthermore, note that Square' is not a functor but a higher-order functor as it takes functor to functors. Accordingly, Matrix' is a type constructor of kind $(* \to *) \to (* \to *)$. Recall that the kind system of Haskell specifies the 'type' of a type constructor (Jones, 1995). The '*' kind represents nullary constructors like Bool or Int. The kind $\kappa_1 \to \kappa_2$ represents type constructors that map type constructors of kind κ_1 to those of kind κ_2 . Though the type of square matrices looks daunting, it is comparatively easy to construct elements of that type. Here is a square matrix of size 3.

 $\begin{array}{c} Succ \; (Succ \; (Succ \; (Zero \; (Cons \; (Cons \; a_{11} \; (Cons \; a_{12} \; (Cons \; a_{13} \; Nil)))) \\ \; & (Cons \; (Cons \; a_{21} \; (Cons \; a_{22} \; (Cons \; a_{23} \; Nil))) \\ \; & (Cons \; (Cons \; a_{31} \; (Cons \; a_{32} \; (Cons \; a_{33} \; Nil))) \\ \; & (Nil)))))) \end{array}$

Perhaps surprisingly, we have essentially a list of lists! The only difference to the standard representation is that the size of the matrix is additionally encoded into a prefix of *Zero* and *Succ* constructors. It is this prefix that takes care of the size constraints. The following table shows the construction of $Succ^3$ (*Zero* m) in more detail (f^n a means f applied n times to a).

Roughly speaking, the outer applications of the value constructor *Cons* make sure that the inner lists have the same length and *Zero* checks that the inner lists have the same length as the outer list.

This completes the overview. The rest of the paper is organized as follows. Section 2 introduces multisets and operations on multisets. Furthermore, we show how to transform equations into a tail-recursive form. Section 3 explains functors and makes the relationship between multisets and functors precise. A multitude of examples is presented in Section 4: among other things we study random access lists, Braun trees, 2-3 trees, and square matrices. Finally, Section 5 reviews related work and points out directions for future work.

2 Multisets

A multiset of type (a) is a collection of elements of type a that takes account of their number but not of their order. In this paper, we will only consider multisets formed according to the following grammar.

 $M ::= \varnothing | \langle 0 \rangle | \langle 1 \rangle | (M \uplus M) | (M + M) | (M * M)$

Here, \varnothing denotes the empty multiset, (n) denotes the singleton multiset which contains n exactly once, (\boxminus) denotes multiset union, (+) and (*) are addition and multiplication lifted to multisets, ie $A \otimes B = (a \otimes b \mid a \leftarrow A; b \leftarrow B)$ for $\otimes \in \{+, *\}$. If the meaning can be resolved from the context, we abbreviate (n) by n. Furthermore, we agree upon that multiplication takes precedence over addition, which in turn takes precedence over multiset union.

Multisets are defined by *higher-order recursion equations*. Higher-order means that the equations may not only involve multisets, but also functions over multisets, function over functions over multisets etc. In this paper, we will, however, restrict ourselves to first-order equations. The exploration of higher-order kinds is the topic of future research. The meaning of higher-order recursion equations is given by the usual least fixpoints semantics.

A multiset is called *simple* iff it is either the empty multiset or a multiset containing a single element arbitrarily often. Simple multiset are denoted by lower case letters. A product A * B is called simple iff B is simple. For instance, nat * 2 is simple while nat * nat is not. We will see in Section 3 that only simple products correspond to compositions of functors. That is, nat * 2 corresponds to $Nat \cdot (Id \times Id)$ but nat * nat does not correspond to $Nat \cdot Nat$. For that reason, we confine ourselves to simple products when defining multisets.

A multiset is called *unique* iff each element occurs at most once. For instance, *pos* given by $pos = 1 \uplus 1 + pos$ is unique whereas $pos = 1 \uplus pos + pos$ denotes a non-unique multiset. Note that the first definition corresponds to non-empty lists and the second to leaf trees. The ability to distinguish between unique and non-unique representations is the main reason for using multisets instead of sets.

The multiset operations satisfy a variety of laws listed in Figure 1. The laws have been chosen so that they hold both for multisets and for the corresponding operations on functors. This explains why, for instance, a * b = b * a is restricted to simple sets: the corresponding property on functors, $F \cdot G = G \cdot F$, does not hold in general. Of course, for functors the equations state isomorphisms rather than equalities.

In the introduction we have transformed the recursive definition of the multiset of all natural numbers into a tail-recursive form. In the rest of this section we will study this transformation in more detail. A function $h :: \langle a \rangle \to \langle a \rangle$ on multisets is said to be a homomorphism iff $h \otimes = \emptyset$ and $h (A \uplus B) = h A \uplus h B$. For instance, h N = A + N * b is a homomorphism while g N = N + N is not. Let h_1, \ldots, h_n be homomorphisms, let A be a multiset, and let X be given by

$$X = A \uplus h_1 X \uplus \cdots \uplus h_n X .$$

Manufacturing Datatypes

$\begin{array}{l} \langle m \rangle + \langle n \rangle \\ \langle m \rangle * \langle n \rangle \end{array} =$	$= \begin{array}{l} \lambda m + n \\ m + n \\ m + n \end{array}$	$\begin{array}{rcl} A \uplus (B \uplus C) &=& (A \uplus B) \uplus C \\ A \uplus B &=& B \uplus A \end{array}$
$\begin{array}{rcl} A + (B + C) & = \\ A + B & = \\ 0 + A & = \end{array}$	(A+B) + C $B+A$ A	
$\begin{array}{rcl} A*(B*C) & = \\ a*b & = \\ 1*A & = \\ A*1 & = \end{array}$	(A * B) * C b * a A A	$\begin{array}{rcl} (A \uplus B) + C &=& A + C \uplus B + C \\ (A \uplus B) \ast C &=& A \ast C \uplus B \ast C \\ (A + B) \ast c &=& A \ast c + B \ast c \\ 0 \ast A &=& 0 \end{array}$

A, B, C are multisets a, b, c are simple multisets m, n are natural numbers

Fig. 1. Laws of the operations.

The definition of X is not tail-recursive as the recursive occurrences of X are nested inside function calls. Note that *nat* is an instance of this scheme with $A = \langle 0 \rangle$, n = 1, and $h_1 N = \langle 1 \rangle + N$. Now, the *tail-recursive variant* of X is f A with fgiven by

$$f N = N \uplus f (h_1 N) \uplus \cdots \uplus f (h_n N) .$$

The definition of f is called *tail-recursive* for obvious reasons. Note that nat' (0) is the tail-recursive variant of nat. The correctness of the transformation is implied by the following theorem.

Theorem 1 Let X :: (a), A :: (a), and $f :: (a) \to (a)$ be given as above, then X = f A.

3 Functors

In close analogy to multiset expressions we define the syntax of *functor* expressions by the following grammar.

$$F ::= K Void \mid K Unit \mid Id \mid (F \mid F) \mid (F \times F) \mid (F \cdot F)$$

Here, K t denotes the constant functor given by K t a = t, Void is the empty type, and Unit is the unit type containing a single element. By Id we denote the identity functor given by Id a = a; $F_1 \cdot F_2$ denotes functor composition given by $(F_1 \cdot F_2)$ $a = F_1$ $(F_2 \ a)$. Disjoint sums and products are defined pointwise: $(F_1 | F_2)$ $a = F_1 \ a | F_2 \ a$ and $(F_1 \times F_2)$ $a = F_1 \ a \times F_2 \ a$.

All these constructs can be easily defined in Haskell. First of all, we require the following type definitions.

type Unit	=	()
data Either $a_1 a_2$	=	Left $a_1 \mid Right \mid a_2$
$\mathbf{data}\;(a_1,a_2)$	=	(a_1,a_2)

The predefined types *Either* a_1 a_2 and (a_1, a_2) implement disjoint sums and prod-

ucts. The operations on functors are then defined by

Using these type constructors it is straightforward to translate a functor equation into a Haskell datatype definition. For reasons of readability, we will often define special instances of the general schemes writing Nil instead of K Unit or Cons t instead of Prod Id t.

The translation of multisets into functors is given by the following table.

m_1	m_2	Ø	205	215	$m_1 \uplus m_2$	$m_1 + m_2$	$m_1 * m_2$
f_1	f_2	K Void	K Unit	Id	$f_1 \mid f_2$	$f_1 \times f_2$	$f_1 \cdot f_2$

We say that F corresponds to M if F is obtained from M using this translation. In the rest of this section we will briefly sketch the correctness of the translation. Informally, the functor corresponding to a multiset M contains, for each member of M, a container of that size. This statement can be made precise using the framework of polytypic programming (Hinze, 1999b). Briefly, a polytypic function is one which is defined by induction on the structure of functor expressions. A simple example for a polytypic function is $sum\langle f \rangle :: f \mathbb{N} \to \mathbb{N}$ which sums a structure of natural numbers. To make the relationship between multisets and functors precise we furthermore require the function $fan\langle f \rangle :: a \to \langle f a \rangle$ which generates the multiset of all structures of type f a from a given seed of type a.

Theorem 2

If the functor F corresponds to the multiset M and if M's definition only involves simple products, then $M = \langle sum \langle F \rangle \ a \mid a \leftarrow fan \langle F \rangle \ 1 \rangle$.

The following example shows that it is necessary to restrict products to simple products: if we compose the functors corresponding to (1, 2) and (1, 3) we obtain a functor which corresponds to (1, 2, 3, 4, 4, 6). In general, functor composition corresponds to the multiset operation (\circledast) given by

$$A \circledast B = \langle b_1 + \dots + b_a \mid a \leftarrow A; b_1 \leftarrow B; \dots; b_a \leftarrow B \rangle$$

We take a container of type A and fill each of its slots with a container of type B. Summing the sizes of the B containers yields the overall size. The operations (*) and (\circledast) coincide only for simple products, ie if the containers of type B all have equal size.

4 Examples

In this section we apply the framework to generate efficient implementations of vectors (aka lists or sequences or arrays) and matrices.

4.1 Lists

A vector or a sequence type contains containers of arbitrary size. The problem related to designing a sequence type is, of course, to generate the multiset of all natural numbers. Different ways to describe this set correspond to different implementations of vectors. Perhaps surprisingly, there is an abundance of ways to solve this problem. In the introduction we already encountered the most direct solution:

 $nat_0 = 0 \uplus 1 + nat_0$.

If we transform the corresponding functor equation

 $Nat_0 = K Unit | Id \times Nat_0$

into a Haskell datatype, we obtain the ubiquitous datatype of polymorphic lists.

data Vector $a = Nil \mid Cons \mid a \mid Vector \mid a$

As an example, the list representation of the vector (0, 1, 2, 3, 4, 5) is

 $Cons \ 0 \ (Cons \ 1 \ (Cons \ 2 \ (Cons \ 3 \ (Cons \ 4 \ (Cons \ 5 \ Nil)))))$.

The tail-recursive variant of nat_0 is given by

$$\begin{array}{rcl} nat_1 & = & nat_1' \ 0 \\ nat_1' \ n & = & n \uplus nat_1' \ (1+n) \end{array} .$$

From the functor equations

$$\begin{array}{rcl} Nat_1 & = & Nat_1' \ (K \ Unit) \\ Nat_1' \ f & = & t \mid Nat_1' \ (Id \ \times f) \end{array}$$

we can derive the following datatype definitions.

type Vector = Vector' Nil data Vector' t a = Zero (t a) | Succ (Vector' (Cons t) a)

Using this representation the vector (0, 1, 2, 3, 4, 5) is written somewhat lengthy as

Succ (Succ (Succ (Succ (Succ (Succ (Zero (Cons 0 (Cons 1 (Cons 2 (Cons 3 (Cons 4 (Cons 5 Nil)))))))))))).

Fortunately, we can simplify the definitions slightly. Recall that *Vector'* is a type constructor of kind $(* \rightarrow *) \rightarrow (* \rightarrow *)$. In this case the 'higher-orderness' is, however, not required. Noting that the first argument of *Vector'* is always applied to the second we can transform *Vector'* into a first-order functor of kind $* \rightarrow * \rightarrow *$.

```
type Vector = Vector' ()
data Vector' t a = Zero t \mid Succ (Vector' (a, t) a)
```

The two variants are related by $Vector'_{ho} t \ a = Vector'_{fo} (t \ a) \ a$ and $Vector'_{fo} t \ a = Vector'_{ho} (K \ t) \ a$. Note that the type Matrix' defined in the introduction is not amenable to this transformation since the first argument of Matrix' is used at different instances. Using the first-order definition (0, 1, 2, 3, 4, 5) is represented by

Succ (Succ (Succ (Succ (Succ (Succ (Succ (Zero (0, (1, (2, (3, (4, (5, ())))))))))))).

4.2 Random-access lists

The definition of nat_0 is based on the unary representation of the natural numbers: a natural number is either zero or the successor of a natural number. Of course, we can also base the definition on the binary number system: a natural number is either zero, even, or odd.

$$nat_2 = 0 \uplus nat_2 * 2 \uplus 1 + nat_2 * 2$$

Transforming the corresponding functor equation

1

$$Nat_2 = K Unit | Nat_2 \cdot (Id \times Id) | Id \times Nat_2 \cdot (Id \times Id)$$

into a Haskell datatype yields

data Vector a = Null | Zero (Vector (a, a)) | One a (Vector (a, a)).

Interestingly, this definition implements random-access lists (Okasaki, 1998), which support logarithmic access to individual vector elements. A random-access list is basically a sequence of perfect leaf trees of increasing height. The vector (0, 1, 2, 3, 4, 5), for instance, is represented by

Zero (One (0,1) (One ((2,3), (4,5)) Null)).

The sequence of Zero and One constructors encodes the size of the vector in binary representation (with the least significant bit first): we have $(011)_2 = 6$. The representation of a vector of size 11 is depicted in Figure 2(a). Note that the representation is not unique because of leading zeros: the empty sequence, for example, can be represented by Null, Zero Null, Zero (Zero Null) etc. There are at least two ways to repair this defect. The following definition ensures that the leading digit is always a one.

$$\begin{array}{rcl} nat_3 & = & 0 \uplus pos_3 \\ pos_3 & = & 1 \uplus pos_3 \ast 2 \uplus 1 + pos_3 \ast 2 \end{array}$$

More elegantly, one can define a zeroless representation (Okasaki, 1998) which employs the digits 1 and 2 instead of 0 and 1. We call this variant of the binary number system 1-2 system.

$$nat_4 = 0 \uplus 1 + nat_4 \ast 2 \uplus 2 + nat_4 \ast 2$$

This alternative has the further advantage that accessing the *i*-th element runs in $O(\log i)$ time (Okasaki, 1998).

4.3 Fork-node trees

Now, let us transform nat_3 into a tail-recursive form.

 $nat_5 = 0 \uplus pos'_5 1$ $pos'_5 n = n \uplus pos'_5 (n * 2) \uplus pos'_5 (1 + n * 2)$

Note that we may replace n * 2 by 2 * n = n + n if pos'_5 is called with a simple multiset as in $pos'_5 1$. The corresponding functor equations look puzzling.

In order to improve the readability of the derived datatypes let us define idioms for 2 * n = n + n and 1 + 2 * n = 1 + n + n.

data Fork
$$t a = Fork (t a) (t a)$$

data Node $t a = Node a (t a) (t a)$

These definitions assume that t is a simple functor. The alternative definitions **newtype** Fork t a = Fork (t (a, a)) and **data** Node t a = Node a (t (a, a)), which correspond to n * 2 and 1 + n * 2, work for arbitrary functors but are more awkward to use. Building upon Fork and Node the Haskell datatypes read

A vector of size n is represented by a complete binary tree of height $\lfloor \log_2 n \rfloor + 1$. A node in the *i*-th level of this tree is labelled with an element iff the *i*-th digit in the binary decomposition of n is one. The lowest level, which corresponds to a leading one, always contains elements. To the best of the author's knowledge this data structure, which we baptize fork-node trees for want of a better name, has not been described elsewhere. Our running example, the vector (0, 1, 2, 3, 4, 5), is represented by

NonEmpty (One (Zero (Base (Fork (Node 0 (Id 1) (Id 2)) (Node 3 (Id 5) (Id 5)))))).

Again, the size of the vector is encoded into the prefix of constructors: replacing *NonEmpty* and *One* by 1 and *Zero* by 0 yields the binary decomposition of the size with the most significant bit first. Figure 2(b) shows a sample vector of 11 elements. The vector elements are stored in left-to-right preorder: if the tree has a root, it contains the first element; the elements in the left tree precede the elements in the right tree. This layout is, however, by no means compelling. Alternatively, one could store the elements in level order. This choice facilitates the extension of a vector at the front but complicates accessing a vector element.

As always for vector types we can 'firstify' the type definitions.

The representation of (0, 1, 2, 3, 4, 5) now consists of nested pairs and triples.

```
NonEmpty (One (Zero (Base ((0,1,2), (3,4,5)))))
```

Finally, let us remark that the tail-recursive variant of nat_4 , which is based on the 1-2 system, yields a similar tree shape: a node on the *i*-th level contains *d* elements where *d* is the *i*-th digit in the 1-2 decomposition of the vector's size.

4.4 Rightist right-perfect trees

The definition of nat_2 is based on the fact that all natural numbers can be generated by shifting (n * 2) and setting the least significant bit (1 + n * 2). The following definition sets bits at arbitrary positions by repeatedly shifting a one.

$$\begin{array}{rcl} nat_{6} & = & nat_{6}' \ 1 \\ nat_{6}' \ p & = & 0 \uplus \ nat_{6}' \ (p*2) \uplus \ p + nat_{6}' \ (p*2) \end{array}$$

Of course, the two definitions are not unrelated, we have

$$nat_2 * p = nat'_6 p , \qquad (1)$$

ie $nat_6' p$ generates all multiples of p. In the *i*-th level of recursion the parameter of nat_6' equals $p * 2^i$ if the initial call was $nat_6' p$. Now, transforming the corresponding functor equations, which assume that f is simple,

$$\begin{array}{rcl} Nat_{6} & = & Nat_{6}' \ Id \\ Nat_{6}' \ f & = & f \mid Nat_{6}' \ (f \times f) \mid f \times Nat_{6}' \ (f \times f) \end{array}$$

into Haskell datatypes yields

This datatype implements higher-order random-access lists (Hinze, 1998). If we 'firstify' the type constructor Vector', we obtain the first-order variant as defined in Section 4.2. For a discussion of the tradeoffs we refer the interested reader to (Hinze, 1998). The vector (0, 1, 2, 3, 4, 5) is represented by

Zero (One (Fork (Id 0) (Id 1)) (One (Fork (Fork (Id 2) (Id 3)) (Fork (Id 4) (Id 6))) Null)).

Interestingly, using a slight generalization of Theorem 1 we can transform nat'_6

into a tail-recursive form, as well.

The function nat'_7 is related to nat_2 by

$$n + nat_2 * p = nat'_7 n p . (2)$$

Assuming that p is simple we get the following functor equations

$$\begin{array}{rcl} Nat_7 & = & Nat'_7 \ (K \ Unit) \ Id \\ Nat'_7 \ f \ p & = & f \ | \ Nat'_7 \ f \ (p \times p) \ | \ Nat'_7 \ (f \times p) \ (p \times p) \end{array}$$

from which we can easily derive the datatype definitions below.

This datatype implements rightist right-perfect trees or RR-trees (Dielissen & Kaldewaij, 1995) where the offsprings of the nodes on the left spine form a sequence of perfect trees of decreasing height. Note that if we change *Prod* t p to *Prod* p t in the last line we obtain *leftist left-perfect trees*. Here is the vector (0, 1, 2, 3, 4, 5) written as an RR-tree.

Reading the constructors *Even* and Odd as digits (LSB first) gives the size of the vector. A sample vector of size 11 is shown in Figure 2(c). The 'firstification' of *Vector'* is left as an exercise to the reader.

4.5 Braun trees

Let us apply the framework to design a representation of *Braun trees* (Braun & Rem, 1983). Braun trees are node-oriented trees which are characterized by the following balance condition: for all subtrees, the size of the left subtree is either exactly the size of the right subtree, or one element larger. In other words, a Braun tree of size 2 * n + 1 has two children of size n and a Braun tree of size 2 * n + 2 has a left child of size n + 1 and a right child of size n. This motivates the following definition.

 The arguments of braun' are always two successive natural numbers. From the corresponding functor equations

we can derive the following datatype definitions.

$$\begin{array}{rcl} \textbf{data} \ Bin \ t_1 \ t_2 \ a & = & Bin \ (t_1 \ a) \ a \ (t_2 \ a) \\ \textbf{type} \ Braun & = & Braun' \ (K \ Unit) \ Id \\ \textbf{data} \ Braun' \ t \ t' \ a & = & Null \ (t \ a) \\ & & | & One \ (Braun' \ (Bin \ t \ t) \ (Bin \ t' \ t) \ a) \\ & & | & Two \ (Braun' \ (Bin \ t' \ t) \ (Bin \ t' \ t') \ a) \end{array}$$

Interestingly, Braun trees are based on the 1-2 number system (MSB first). The vector (0, 1, 2, 3, 4, 5), for instance, is represented as follows.

```
Two (Two (Null (Bin (Bin (Id 0) 1 (Id 2)) 3 (Bin (Id 4) 5 (K ())))))
```

Figure 2(d) displays the representation of a vector of 11 elements. R. Paterson has described a similar implementation (personal communication).

4.6 2-3 trees

Up to now we have mainly considered unique representations where the shape of a data structure is completely determined by the number of elements it contains. Interestingly, unique representations are not well-suited for implementing search trees: one can prove a lower bound of $\Omega(\sqrt{n})$ for insertion and deletion in this case (Snyder, 1977). For that reason, popular search tree schemes such as 2-3 trees (Aho *et al.*, 1983), red-black trees (Guibas & Sedgewick, 1978), or AVL-trees (Adel'son-Vel'skiĭ & Landis, 1962) are always based on non-unique representations. Let us consider how to implement, say, 2-3 trees. The other search tree schemes can be handled analogously. The definition of 2-3 trees is similar to that of perfect leaf trees: a 2-3 tree of height 0 is a leaf and a 2-3 tree of height h + 1 is a node with either two or three children, each of which is a 2-3 tree of height h. This similarity suggests to model 2-3 trees as follows.

$$tree23 = tree23' 0$$

$$tree23' N = N \uplus tree23' (N+1+N \uplus N+1+N+1+N)$$

Note that contrary to previous definitions the parameter of the auxiliary function does not range over simple sets. The corresponding functor equations



Fig. 2. Different representations of a vector with 11 elements.

give rise to the following datatype definitions.

type Tree23 a	=	Tree23' Nil a
data Tree23' t a	=	$Zero~(t~a) \mid Succ~(Tree 23'~(Node 23~t)~a)$
data Node23 t a	=	$Node2 \ (t \ a) \ a \ (t \ a) \ \ Node3 \ (t \ a) \ a \ (t \ a) \ a \ (t \ a)$

The vector (0, 1, 2, 3, 4, 5) has three different representations; one alternative is

Algorithms for insertion and deletion are described in (Hinze, 1998).



Fig. 3. The representation of a 6×6 matrix based on fork-node trees.

4.7 Matrices

Let us finally design representations of square matrices and rectangular matrices. In the introduction we have already discussed the central idea: we take a tail-recursive definition of the natural numbers (or of the positive numbers)

$$\begin{array}{rcl} X & = & f & a \\ f & n & = & n \uplus f & (h_1 & n) \uplus \cdots \uplus f & (h_n & n) \end{array}$$

and replace n by n * n in the second equation:

$$square = square' a$$

$$square' n = n * n \uplus square' (h_1 n) \uplus \dots \uplus square' (h_n n)$$

This transformation works provided a is a simple multiset and the h_i preserve simplicity. These conditions hold for all of the examples above with the notable exception of 2-3 trees. As a concrete example, here is an implementation of square matrices based on fork-node trees.

The representation of a 6×6 matrix is shown in Figure 3.

Rectangular matrices are equally easy to implement. In this case we replace n by

nat * n in the second equation:

$$rect = rect' a$$

 $rect' n = nat * n \uplus rect' (h_1 n) \uplus \cdots \uplus rect' (h_n n)$.

Alternatively, one may use the following scheme.

This representation requires more constructors than the first one $(n^2 + 1 \text{ instead of } n+1)$. On the positive side, it can be easily generalized to higher dimensions.

5 Related and future work

This work is inspired by a recent paper of C. Okasaki (Okasaki, 1999) who derives representations of square matrices from exponentiation algorithms. He shows, in particular, that the tail-recursive version of the fast exponentiation gives rise to an implementation based on rightist right-perfect trees. Interestingly, the simpler implementation based on fork-node trees is not mentioned. The reason is probably that fast exponentiation algorithms typically process the bits from least to most significant bit while fork-node trees and Braun trees are based on the reverse order. The relationship between number systems and data structures is explained at great length in (Okasaki, 1998). The development in Section 3 can be seen as putting this design principle on a formal basis.

Directions for future work suggest themselves. It remains to adapt the standard vector and matrix algorithms to the new representations. Some preparatory work has been done in this respect. In (1998) the author shows how to adapt search tree algorithms to nested representations of search trees using constructor classes. It is conceivable that this approach can be applied to matrix algorithms, as well. Furthermore, many functions like *map*, *listify*, *sum* etc can be generated automatically using the technique of polytypic programming (Hinze, 1999b). On the theoretical side, it would be interesting to investigate the expressiveness of the framework and of higher-order polymorphic types in general. Which class of multisets can be described using higher-order recursion equations? For instance, it appears to be impossible to specify the multisets of all prime numbers. Do higher-order kinds increase the expressiveness?

References

- Adel'son-Vel'skiĭ, G.M., & Landis, Y.M. (1962). An algorithm for the organization of information. *Doklady akademiia nauk SSSR*, **146**, 263-266. English translation in Soviet Math. Dokl. 3, pp. 1259-1263.
- Aho, Alfred V., Hopcroft, John E., & Ullman, Jeffrey D. (1983). Data structures and algorithms. Addison-Wesley Publishing Company.

- Bird, Richard, & de Moor, Oege. (1997). Algebra of programming. London: Prentice Hall Europe.
- Bird, Richard, & Meertens, Lambert. (1998). Nested datatypes. Pages 52-67 of: Jeuring, J. (ed), Fourth international conference on mathematics of program construction, MPC'98, Marstrand, Sweden. Lecture Notes in Computer Science, vol. 1422. Springer Verlag.
- Braun, W., & Rem, M. (1983). A logarithmic implementation of flexible arrays. Memorandum MR83/4, Eindhoven University of Technology.
- Cormen, Thomas H., Leiserson, Charles E., & Rivest, Ronald L. (1991). Introduction to algorithms. Cambridge, Massachusetts: The MIT Press.
- Dielissen, Victor J., & Kaldewaij, Anne. (1995). A simple, efficient, and flexible implementation of flexible arrays. Pages 232-241 of: Third international conference on mathematics of program construction (MPC'95). Lecture Notes in Computer Science, vol. 947. Springer Verlag.
- Guibas, Leo J., & Sedgewick, Robert. (1978). A diochromatic framework for balanced trees. Pages 8-21 of: Proceedings of the 19th annual symposium on foundations of computer science. IEEE Computer Society.
- Hinze, Ralf. 1998 (December). Numerical representations as higher-order nested datatypes. Tech. rept. IAI-TR-98-12. Institut für Informatik III, Universität Bonn.
- Hinze, Ralf. 1999a (March). Perfect trees and bit-reversal permutations. Tech. rept. IAI-TR-99-4. Institut für Informatik III, Universität Bonn.
- Hinze, Ralf. 1999b (March). Polytypic functions over nested datatypes (extended abstract). 3rd latin-american conference on functional programming (CLaPF'99).
- Jones, Mark P. (1995). Functional programming with overloading and higher-order polymorphism. Pages 97-136 of: First international spring school on advanced functional programming techniques. Lecture Notes in Computer Science, vol. 925. Springer Verlag.
- Okasaki, Chris. (1998). Purely functional data structures. Cambridge University Press.
- Okasaki, Chris. (1999). From fast exponentiation to square matrices: An adventure in types. Submitted for publication.
- Peyton Jones, Simon, & Hughes, John (eds). 1998 (December). Haskell 98 A non-strict, purely functional language.
- Snyder, Lawrence. (1977). On uniquely represented data structures (extended abstract). Pages 142-146 of: 18th annual symposium on foundations of computer science, Providence. Long Beach, Ca., USA: IEEE Computer Society Press.