# Derivable Type classes

RALF HINZE

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


SIMON PEYTON JONES

*Microsoft Research Ltd*
*St George House, 1 Guildhall St, Cambridge CB2 3NH, England*
(*e-mail:* `simonpj@microsoft.com`)

## Abstract

Generic programming allows you to write a function once, and use it many times at different types. A lot of good foundational work on generic programming has been done. The goal of this paper is to propose a practical way of supporting generic programming within the Haskell language, while going "with the grain" of the language.

On the way, we came across a separate issue, concerning type-class overloading where higher kinds are involved. We propose a simple type-class system extension to allow the programmer to write richer contexts than is currently possible.

## 1 Introduction

A *generic*, or *polytypic*, function is one that the programmer writes once, but which works over many different data types. The standard examples are parsing and printing, serialising, taking equality, ordering, hashing, and so on. There is lots of work on generic programming (Bird *et al.*, 1996; Jansson & Jeuring, 1997; Backhouse *et al.*, 1999; Hinze, 2000b).

In this paper we present the design and implementation of an extension to Haskell that supports generic programming. At first sight it might seem that Haskell's type classes are in competition with generic programming — after all, both concern functions that work over many data types. But we have found that they can be combined very gracefully, offering a smooth upward extension to Haskell.

On the way we describe an orthogonal, but complementary idea. Haskell allows one to define higher-order kinded data types for which it is impossible to define, for example, an equality instance. This seems unfortunate: one part of the language is more powerful than another. We describe a modest extension of Haskell's type-class system that removes this difficulty.

More specifically, our contributions are these:

- We present the language design for an extension to Haskell that supports generic programming (Sections 2–6). Generic functions appear solely in class declarations.
- We describe an entirely separate extension that lets one write certain instance declarations for higher-order kinded data types that are simply inexpressible in Haskell 98 (Sections 7–8).
- We discuss the implementation of both parts.

The first part of this paper is a follow-up to (Hinze, 1999); the new achievements are detailed in Section 9.

## 2 Setting the scene

In this section we set the context for our proposal. We begin by reviewing Haskell's type-class overloading mechanism.

### 2.1 A brief review of type class overloading

Haskell supports overloading, based on type classes. For example, the Prelude defines the class *Eq*:

$$\textbf{class } \textit{Eq } t \textbf{ where}$$
$$(==), (/=) \quad :: \quad t \to t \to \textit{Bool} \ .$$

This defines two overloaded top-level functions, $(==)$ and $(/=)$, whose types are

$$(==), (/=) \quad :: \quad (\textit{Eq } t) \Rightarrow t \to t \to \textit{Bool} \ .$$

Before we can use $(==)$ on values of, say *Int*, we must explain how to take equality over *Int* values:

$$\textbf{instance } \textit{Eq Int } \textbf{where}$$
$$(==) \quad = \quad \textit{eqInt}$$
$$(/=) \quad = \quad \textit{neInt} \ .$$

Here we suppose that $\textit{eqInt} :: \textit{Int} \to \textit{Int} \to \textit{Bool}$, and similarly *neInt* are provided from somewhere. The **instance** declarations says "the $(==)$ function at type *Int* is implemented by *eqInt*".

How can we take equality of pairs of values? Presumably by comparing their components; but that requires equality over the component types:

$$\textbf{instance } (\textit{Eq } a, \textit{Eq } b) \Rightarrow \textit{Eq } (a, b) \textbf{ where}$$
$$(x_1, y_1) == (x_2, y_2) \quad = \quad (x_1 == x_2) \wedge (y_1 == y_2)$$
$$(x_1, y_1) /= (x_2, y_2) \quad = \quad (x_1 /= x_2) \vee (y_1 /= y_2) \ .$$

It is a bit annoying to keep having to write code for the $(/=)$ method, since it is

simply the negation of the code for the $(==)$ method, so Haskell allows you to write a *default method* in the class declaration:

$$
\begin{array}{lll}
\textbf{class } Eq\ t\ \textbf{where} & & \\
(==), (/=) & :: & t \to t \to Bool \\
x_1\ /= x_2 & = & not\ (x_1 == x_2)\ .
\end{array}
$$

Now, if you give an instance declaration for $Eq$ that lacks a definition for $(/=)$, Haskell "fills in" the missing method definition with code copied from the class declaration. So we can write:

$$
\begin{array}{lll}
\textbf{instance } (Eq\ a, Eq\ b) \Rightarrow Eq\ (a, b)\ \textbf{where} & & \\
(x_1, y_1) == (x_2, y_2) & = & (x_1 == x_2) \wedge (y_1 == y_2)
\end{array}
$$

and get just the same effect as before. You can even specify a default method for both methods:

$$
\begin{array}{lll}
\textbf{class } Eq\ t\ \textbf{where} & & \\
(==), (/=) & :: & t \to t \to Bool \\
x_1 == x_2 & = & not\ (x_1\ /= x_2) \\
x_1\ /= x_2 & = & not\ (x_1 == x_2)\ .
\end{array}
$$

In an instance declaration, you can now *either* give a definition for $(==)$, *or* a definition for $(/=)$, or both. If you specify neither, then you will get an infinite loop, unfortunately!

If you give an instance declaration without specifying code for method *op*, and the class has no default method for *op*, then invoking the method will halt the program with an error message. It is not a compile-time error; sometimes a method just doesn't make sense for a particular instance type.

## 2.2 Overloading is not generic programming

Haskell as it stands does not support generic, or polytypic, programming. In particular, suppose you define a new data type:

$$
\textbf{data } Wibble \quad = \quad Wag\ Int \mid Wug\ Bool\ .
$$

It is "obvious" how to take equality over *Wibble*, and support for generic programming would allow us to specify this "obvious" precisely. In Haskell, however, you have to give an explicit instance declaration, containing the code for equality:

$$
\begin{array}{lll}
\textbf{instance } Eq\ Wibble\ \textbf{where} & & \\
(Wag\ i_1) == (Wag\ i_2) & = & i_1 == i_2 \\
(Wug\ b_1) == (Wug\ b_2) & = & b_1 == b_2 \\
w_1 == w_2 & = & False\ .
\end{array}
$$

Cranking out this sort of boilerplate code is so tiresome that Haskell provides special support — the so-called "**deriving**" mechanism — for a handful of built-in classes.

In particular, for *Eq* you say

<div align="center">

**data** *Wibble*  =  *Wag Int* | *Wug Bool* **deriving** (*Eq*) .

</div>

The "**deriving** (*Eq*)" part tells Haskell to generate the "obvious" code for equality. What "obvious" means is specified informally in an Appendix of the language definition (Peyton Jones & Hughes, 1999). This is all rather *ad hoc*, and in particular *it only works for a fixed set of built-in classes* (*Eq*, *Ord*, *Enum*, *Bounded*, *Read*, *Show*, and *Ix*).

### *2.3  Generalising default methods*

What we seek, then, is an automatic mechanism that "fills in" a suitable implementation for the methods of an instance declaration. But wait a minute! That's what a default method does! Indeed so but, as we have already remarked, default methods as they stand are too weak. If we write merely:

<div align="center">

**instance** *Eq Wibble*

</div>

we would, as remarked earlier, just get an infinite loop. We have to provide some real code somewhere! What we want is a richer language in which to write default methods. That is what we turn our attention to now.

## 3  Generic definitions

From a language-design point of view, our story is this: *providing a richer language for default method definitions in a class declaration gives an elegant way to extend Haskell with the power of generic programming.* We will justify this statement more fully in Section 4.4, but first we must present our design.

### *3.1  Two examples*

We adopt with minor changes the proposal in (Hinze, 1999). Two examples will serve to give the idea. First, here is the *Eq* class augmented with generic equality:

```
class Eq t where
  (==), (/=)                        ::  t → t → Bool
      -- generic default method
  (==){1} Unit Unit                 =   True
  (==){ a + b }(Inl x₁) (Inl x₂)    =   x₁ == x₂
  (==){ a + b }(Inr y₁) (Inr y₂)    =   y₁ == y₂
  (==){ a + b }_ _                  =   False
  (==){ a × b }(x₁ :*: y₁) (x₂ :*: y₂) =   x₁ == x₂ ∧ y₁ == y₂
      -- vanilla, non-generic default method
  (/=) x₁ x₂                        =   not (x₁ == x₂) .
```

This new class declaration contains an ordinary, default declaration for $(/=)$, just as before. The new feature is a *generic definition* for equality, distinguished by the curly braces on the left hand side, which enclose a *type* argument. We will study such generic definitions in more detail in Section 4.3. For now, we simply observe that a generic definition works by induction over the structure of the type (written in curly braces) at which the class is instantiated.

Here is another example. The class *Binary* has methods *showBin* and *readBin* that respectively convert a value to a list of bits and vice versa:

$$
\begin{aligned}
&\textbf{data } Bit & &= & &0 \mid 1 \\
&\textbf{type } Bin & &= & &[\,Bit\,] \\
&\textbf{class } Binary\ t\ \textbf{where} \\
&\quad showBin & &:: & &t \to Bin \\
&\quad readBin & &:: & &Bin \to (t, Bin)\ .
\end{aligned}
$$

A real implementation might have a more sophisticated representation for *Bin* but that is a separate matter. We can give generic definitions for *showBin* and *readBin* like this:

$$
\begin{aligned}
&showBin\{1\}\,Unit & &= & &[\,] \\
&showBin\{a+b\}(Inl\ x) & &= & &0 : showBin\ x \\
&showBin\{a+b\}(Inr\ y) & &= & &1 : showBin\ y \\
&showBin\{a \times b\}(x :\!*\!: y) & &= & &showBin\ x \,{+\!\!+}\, showBin\ y \\
\\
&readBin\{1\}\,bs & &= & &(Unit, bs) \\
&readBin\{a+b\}[\,] & &= & &error\ \texttt{"readBin:}_{\sqcup}\texttt{unexpected}_{\sqcup}\texttt{end}_{\sqcup}\texttt{of}_{\sqcup}\texttt{input"} \\
&readBin\{a+b\}(0 : bs) & &= & &\textbf{let } (x, bs') = readBin\ bs\ \textbf{in}\ (Inl\ x, bs') \\
&readBin\{a+b\}(1 : bs) & &= & &\textbf{let } (y, bs') = readBin\ bs\ \textbf{in}\ (Inr\ y, bs') \\
&readBin\{a \times b\}\,bs & &= & &\textbf{let } (x, bs_1) = readBin\ bs \\
& & & & &\qquad\ (y, bs_2) = readBin\ bs_1 \\
& & & & &\textbf{in } ((x :\!*\!: y), bs_2)\ .
\end{aligned}
$$

Notice that *readBin produces* a value of the unknown type $t$, whereas *showBin* and $(==)$ both *consume* such values.

### 3.2 Instances and deriving

Suppose we have given a class declaration of *Eq* with generic default methods, as we did in the previous section. Now we can give an instance declaration like this:

**instance** *Eq Wibble*

*without giving code for either method.* Both methods will be "filled in" from the class declaration. The ordinary, non-generic default method, $(/=)$, is filled in verbatim. The generic default method, $(==)$, is specialised in a way we will describe, to give essentially the code in Section 2.2. That is, the effect of the **instance** declaration

is exactly as if we had written

$$
\begin{array}{lll}
\textbf{instance } Eq \; Wibble \textbf{ where} \\
(Wag \; i_1) == (Wag \; i_2) & = & i_1 == i_2 \\
(Wug \; b_1) == (Wug \; b_2) & = & b_1 == b_2 \\
w_1 == w_2 & = & False \\
w_1 \; /= w_2 & = & not \; (w_1 == w_2) \; .
\end{array}
$$

And that is all. Though this sounds simple enough, it has interesting and important consequences:

- Though the instance declaration is now rather brief, it must still be given to make *Wibble* an instance of *Eq*. It is *not* the case that all types become instances of *Eq* when one gives a generic default method.
- It is not necessary for the instance declaration to appear in the same module as the data type declaration, or the class declaration. By contrast, in Haskell 98 a **deriving** clause must be attached to the data type declaration. This separation is useful, because the class might not even be defined in the scope where the type is declared.
- The compiler only fills in a method definition if the programmer omits it. For example, if we said

$$
\begin{array}{lll}
\textbf{instance } Eq \; Wibble \textbf{ where} \\
(Wag \; \_) == (Wag \; \_) & = & True \\
w_2 == w_2 & = & False
\end{array}
$$

then this programmer-supplied code is, of course, used. This is one way in which our proposal differs from others. In most generic-programming systems, a generic function works generically over all types. In our design, *the programmer can override the generic definition on a type-by-type basis.*

In particular, this is how we specify instances for primitive types such as *Int*, *Double*, "→": we just use ordinary instance declarations.

- A **deriving** clause can now be seen as shorthand (albeit now not so much shorter) for an instance declaration. There is a difference, though. Consider

$$
\textbf{data } Tree \; a \;\; = \;\; Node \; a \; [Tree \; a] \; \textbf{deriving} \; (Eq) \; .
$$

In our design, the **deriving** clause is shorthand for

$$
\textbf{instance } (Eq \; a) \Rightarrow Eq \; (Tree \; a)
$$

Note that in an instance declaration we must explicitly specify the context $(Eq \; a)$, which is inferred automatically by the **deriving** mechanism.

There are two quite separate questions here:

— Inferring the right context for the instance declaration.

— Filling in the code for the methods.

This paper addresses the second of these questions, and not the first. The first is certainly an interesting question in its own right, and not just for instance

declarations. For example, it might be desirable to allow one to say

$$f \quad :: \quad (...) \Rightarrow a \rightarrow a$$

where the "(...)" means "please infer the appropriate context". The ability to write such partial type signatures, with the ellipsis filled in by type inference, has been discussed on the Haskell mailing list. Such an ellipsis makes sense wherever a context is explicitly written in Haskell. In the case of instance declarations, we might like to be able to write

$$\textbf{instance } (...) \Rightarrow Eq \ (Tree \ a)$$

and have Haskell fill in the ellipsis, as well as the method declarations. This is quite feasible (albeit involving a fixpoint iteration) for first-order kinded types, but we believe that it is infeasible for higher-order types (see Section 7.3).

### 3.3 Generic representation types

The arguments in braces on the left hand side of a generic definition are *types*. The idea is, of course, that these generic definitions can be specialised for any particular type. Suppose, for example, we have a data type *List*, and we make *List* an instance of *Binary*:

$$\textbf{data } List \ a \quad = \quad Cons \ a \ (List \ a) \mid Nil$$
$$\textbf{instance } (Binary \ a) \Rightarrow Binary \ (List \ a) \ .$$

How is the compiler to fill in the missing method definitions?

First, we define the *generic representation type* for *List*, which we will call $List^\circ$:

$$\textbf{type } List^\circ \ a \quad = \quad (a \times List \ a) + 1 \ .$$

We will have more to say about representation types in Section 6.2, but for now we can just think of $List^\circ$ as a type that is more-or-less isomorphic to *List*, but one that uses only a small, fixed set of type constructors, namely sums, products, and unit, defined like this:

$$\textbf{data } a \times b \quad = \quad a \ :*: \ b$$
$$\textbf{data } a + b \quad = \quad Inl \ a \mid Inr \ b$$
$$\textbf{data } 1 \qquad = \quad Unit \ .$$

Of course, infix $(\times)$ is not a legal Haskell type constructor, and nor are $(+)$ and 1. We give them special syntax to distinguish them from their "normal" counterparts, $(a, b)$, *Either a b* and (), and extend the syntax of types to accommodate them.

In our example, a *List* is a sum $(+)$ of two types: the unit type (1), and a product $(\times)$ of the element type $a$ and a *List*. To make the isomorphism explicit, let us write functions that convert to and fro:

$$to\text{-}List \qquad\qquad\qquad :: \quad \forall a \, . \, List^\circ \ a \rightarrow List \ a$$
$$to\text{-}List \ (Inl \ (x :*: xs)) \quad = \quad Cons \ x \ xs$$
$$to\text{-}List \ (Inr \ Unit) \qquad = \quad Nil$$

$$
\begin{array}{lll}
\textit{from-List} & :: & \forall a \,.\, \textit{List } a \to \textit{List}^{\circ}\, a \\
\textit{from-List } (\textit{Cons } x\ xs) & = & \textit{Inl } (x \mathbin{:\!*\!:} xs) \\
\textit{from-List Nil} & = & \textit{Inr Unit} \ .
\end{array}
$$

### 3.4 The generic instances

The idea is that *by regarding a List as a List$^{\circ}$, the generic code explains what to do.* The generic method for *showBin*, for example, says what to do if the argument is a sum, what to do if it is a product, and what to do if it is a unit type.

It's useful to imagine re-expressing these default methods as three ordinary instance declarations:

> **instance** *Binary* 1 **where**
> $\quad$ *showBin Unit* $\quad = \quad$ [ ]
> $\quad$ *readBin bs* $\quad\quad = \quad$ (*Unit*, *bs*)
> **instance** (*Binary a*, *Binary b*) $\Rightarrow$ *Binary* (*a* + *b*) **where**
> $\quad$ *showBin* (*Inl x*) $\quad = \quad$ 0 : *showBin x*
> $\quad$ *showBin* (*Inr y*) $\quad = \quad$ 1 : *showBin y*
>
> $\quad$ *readBin* [ ] $\quad\quad\quad = \quad$ *error* `"readBin:␣unexpected␣end␣of␣input"`
> $\quad$ *readBin* (0 : *bs*) $\quad = \quad$ **let** (*x*, *bs'*) = *readBin bs* **in** (*Inl x*, *bs'*)
> $\quad$ *readBin* (1 : *bs*) $\quad = \quad$ **let** (*y*, *bs'*) = *readBin bs* **in** (*Inr y*, *bs'*)
> **instance** (*Binary a*, *Binary b*) $\Rightarrow$ *Binary* (*a* × *b*) **where**
> $\quad$ *showBin* (*x* :*:* *y*) $\quad = \quad$ *showBin x* ++ *showBin y*
> $\quad$ *readBin bs* $\quad\quad\quad = \quad$ **let** (*x*, *bs*$_1$) = *readBin bs*
> $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ (*y*, *bs*$_2$) = *readBin bs*$_1$
> $\quad\quad\quad\quad\quad\quad\quad\quad$ **in** ((*x* :*:* *y*), *bs*$_2$) .

We describe these instance declarations for generic representation types as *generic instance declarations.* They are not written explicitly by the programmer, but instead are derived by the compiler from a class declaration that has generic default methods. We discuss generic instance declarations further in Section 4.2.

### 3.5 Filling in the missing methods

We are now ready to say more precisely how the compiler fills in the missing methods. In this section we sketch the idea using an example, while Section 6 deals with the general case.

When the programmer writes

$$\textbf{instance } (\textit{Binary } a) \Rightarrow \textit{Binary } (\textit{List } a)$$

the compiler will fill in the method declarations as follows:

> $\quad$ *showBin xs* $\quad = \quad$ *showBin* (*from-List xs*)
> $\quad$ *readBin bs* $\quad = \quad$ **case** *readBin bs* **of** (*xs*, *bs'*) $\to$ (*to-List xs*, *bs'*) .

Let us focus on the definition for *showBin*. It works in two stages:

1. First, *from-List* :: ∀*a* . *List a* → *List*° *a* converts the input list of type *List a* into a value of type *List*° *a*.
2. Second, we call the overloaded *showBin* function to complete the job, using the methods of the generic instance declarations.

At first this looks utterly bizarre. We are defining *showBin* in terms of *showBin*. But look at the definition one would write by hand:

$$\textbf{instance } (Binary\ a) \Rightarrow Binary\ (List\ a)\ \textbf{where}$$
$$showBin\ Nil \quad\quad\quad\quad = \quad 0 : [\,]$$
$$showBin\ (Cons\ x\ xs) \quad = \quad 1 : showBin\ x \mathbin{+\!\!+} showBin\ xs \ .$$

The first call is to *showBin* at the list element type; the second is a recursive call to the same *showBin* at the list type.

Something similar happens with the generic definition. Here *showBin* is called on an argument of type *List*° *a*. This is a sum type, so the sum instance of *Binary* kicks in (Section 3.4). It in turn will call *showBin*, once at type 1, and once at type *a* × *List a*. The former has an instance declaration, while the latter uses the product instance and makes calls to *showBin* at type *a* and *List a*. But the former is just like the *showBin x* in the hand-written instance, while the latter is like the *showBin xs*. So everything works out.

Let us return briefly to the first step above. In the case of *showBin* it was fairly simple to convert the argument to its generic representation type. On the other hand *readBin* was a bit more complicated because it returned a pair, only one component of which had to be converted. How, in general, does the compiler perform this conversion? We devote the whole of Section 5 to this topic. First, though, we elaborate on the programmer-visible aspects of our design.

## 4 Discussion and elaboration

We have now sketched the bones of our design. In this section we elaborate on some of the details.

### *4.1 Constructor names and record labels*

Annoyingly, sum, product, and unit are not quite enough. Consider, for example, the standard Haskell class *Show*. To be able to give generic definitions for *showsPrec*, the *names* of the constructors, and their fixities, must be made accessible.

To this end we provide an additional generic representation type, of the form *c* **of** *a* where *c* is a *value* variable of type *ConDescr* and *a* is a *type* variable. The type

*ConDescr* is a new primitive type that comprises all constructor names. To manipulate constructor names the following operations among others can be used — for an exhaustive list see (Hinze, 1999).

$$
\begin{array}{lll}
\textbf{data } ConDescr & \quad \text{-- abstract} \\
\textbf{data } Fixity & = & Nofix \mid Infix\ Int \mid Infixl\ Int \mid Infixr\ Int \\
conName & :: & ConDescr \rightarrow String \quad \text{-- primitive} \\
conArity & :: & ConDescr \rightarrow Int \quad\quad \text{-- primitive} \\
conFixity & :: & ConDescr \rightarrow Fixity \quad \text{-- primitive} \\
\end{array}
$$

**instance** *Show ConDescr* **where**
$$show \quad = \quad conName$$

Using *conName* and *conArity* we can implement a simple variant of Haskell's *Show* class — for the full-fledged version see (Hinze, 1999).

**class** *Show t* **where**
$$
\begin{array}{lll}
show & :: & t \rightarrow String \\
showsPrec & :: & Int \rightarrow t \rightarrow String \\
show\ x & = & showsPrec\ 0\ x \\
showsPrec\{\,a+b\,\}d\ (Inl\ x) & = & showsPrec\ d\ x \\
showsPrec\{\,a+b\,\}d\ (Inr\ y) & = & showsPrec\ d\ y \\
showsPrec\{\,c\ \textbf{of}\ a\,\}d\ (Con\ c\ x) \\
\quad \mid conArity\ c == 0 & = & show\ c \\
\quad \mid otherwise & = & showParen\ (d \geqslant 10) \\
& & \quad (show\ c \mathbin{+\!\!+} \text{"}\sqcup\text{"} \mathbin{+\!\!+} showsPrec\ 10\ x) \\
showsPrec\{\,a \times b\,\}d\ (x :\!\!*\!\!: y) & = & showsPrec\ d\ x \mathbin{+\!\!+} \text{"}\sqcup\text{"} \mathbin{+\!\!+} showsPrec\ d\ y
\end{array}
$$

The representation type $c\ \textbf{of}\ a$ is defined by the following pseudo-Haskell declaration:

$$\textbf{newtype } c\ \textbf{of}\ a \quad = \quad Con\ c\ a \ .$$

Uniquely for Haskell, $c$ is a *value* that is carried by a *type*. It is best to think of the above declaration as defining a family of types: for each $c$ there is a type constructor "$c\ \textbf{of}$" of kind $\star \rightarrow \star$ with a value constructor "$Con\ c$" of type $a \rightarrow (c\ \textbf{of}\ a)$. Now, why does the type $c\ \textbf{of}\ a$ incorporate information about $c$? One might suspect that it is sufficient to supply this information on the value level. Doing so would work for *show*, but would fail for *read*:

**class** *Read t* **where**
$$
\begin{array}{lll}
read & :: & String \rightarrow [(t, String)] \\
read\{\,c\ \textbf{of}\ a\,\}s & = & [(Con\ c\ x, s_3) \mid (s_1, s_2) \leftarrow lex\ s, s_1 == conName\ c, \\
& & \quad\quad\quad\quad (x, s_3) \leftarrow read\ s_2] \ .
\end{array}
$$

The important point is that *read produces* (not consumes) the value, and yet it requires access to the constructor name.

Haskell allows the programmer to assign labels to the components of a constructor, and these, too, are needed by *read* and *show*. For the purpose of presentation,

however, we choose to ignore field names. In fact, they can be handled completely analogously to constructor names.

## 4.2 Generic instance declarations

In Section 3.4 we said that the generic definitions in a class declaration are re-expressed by the compiler as a set of instance declarations, one for each generic representation type. One might ask: why not get the programmer to write these instance declarations directly?

Our answer is stylistic rather than technical. We want to present generic programming in Haskell as a richer language in which to write default method declarations, and scattering them over several instance declarations does not convey that message. The question about whether a generic-default declaration was available to use would become more diffuse, because some parts, but not others, might be available. Writing the declaration all at once, in the class declaration, seems to be the simplest and most direct thing to do, even though it does involve a little new syntax.

Another stylistic reason for our decision is that it is rather easy to confuse the generic instance declaration for, say, products $a \times b$ with "ordinary" instance declarations for the corresponding "ordinary" product type $(a, b)$. For example, in the case of *Show*, the ordinary instance declaration for products might look like this:

**instance** $(Show\ a, Show\ b) \Rightarrow Show\ (a, b)$ **where**
$\quad showsPrec\ d\ (x, y) \quad = \quad$ "(" $+\!\!+ showsPrec\ 0\ x +\!\!+$ "," $+\!\!+ showsPrec\ 0\ y +\!\!+$ ")".

Because tuples are typically shown using distfix notation, we choose to over-ride the generic definition. Nevertheless, the class declaration for *Show* will have given rise to the generic instance declaration

$\quad$ **instance** $(Show\ a, Show\ b) \Rightarrow Show\ (a \times b)$ **where**
$\qquad showsPrec\ d\ (x :\!*: y) \quad = \quad showsPrec\ d\ x +\!\!+$ "␣" $+\!\!+ showsPrec\ d\ y$ .

Recall that products $a \times b$ are used to represent the arguments of a constructor. Consequently, the generic instance declaration specifies the layout of constructor arguments: they are shown consecutively separated by spaces.

## 4.3 Generic default declarations

In general, a class declaration consists of a signature, which specifies the types of the class methods, and an implementation part, which gives default definitions for the class methods. The type signature has the general form:

$\quad$ **class** $ctx \Rightarrow C\ a$ **where**
$\qquad op_1 \quad :: \quad Op_1\ a$
$\qquad \cdots$
$\qquad op_n \quad :: \quad Op_n\ a$ .

The implementation part consists of generic and non-generic default definitions. A non-generic definition is an ordinary Haskell definition

$$op \quad = \quad \dots \quad .$$

A generic definition can be recognised by the *type patterns* on the left hand side, enclosed in curly braces. It has the schematic form

$$
\begin{aligned}
op\{\,()\,\} &= \dots \\
op\{\,a + b\,\} &= \dots \\
op\{\,a \times b\,\} &= \dots \\
op\{\,c \textbf{ of } a\,\} &= \dots \quad .
\end{aligned}
$$

The type patterns are mandatory, so that the equations can be correctly grouped. For example, consider the generic definition of $(==)$ given earlier:

$$
\begin{aligned}
(==)\{\,1\,\}\,Unit\ Unit &= \quad True \\
(==)\{\,a + b\,\}\,(Inl\ x_1)\ (Inl\ x_2) &= \quad x_1 == x_2 \\
(==)\{\,a + b\,\}\,(Inr\ y_1)\ (Inr\ y_2) &= \quad y_1 == y_2 \\
(==)\{\,a + b\,\}\,\_\ \_ &= \quad False \\
(==)\{\,a \times b\,\}\,(x_1 :*: y_1)\ (x_2 :*: y_2) &= \quad x_1 == x_2 \wedge y_1 == y_2 \quad .
\end{aligned}
$$

Without the type patterns there is no way to decide whether the second but last equation belongs to the $(+)$ or to the $(\times)$ case.

Apart from the type patterns, an generic definition has exactly the form of a normal Haskell definition.

We note the following points:

- A class declaration may specify an arbitrary mixture of generic and non-generic default-method declarations. In the case of *Eq* above, $(==)$ is defined by induction over the argument type, while $(/=)$ is non-generic. The generic and non-generic methods may be mutually recursive.

- Class declarations are the *only* place that generic definitions appear in our design. There are no "free-standing" generic definitions.

- At the moment, generic default declarations may only be given for *type* classes, that is, for classes whose type parameter ranges over types of kind $\star$. For example, we cannot specify a generic default method for the *Functor* class:

$$
\begin{aligned}
&\textbf{class } Functor\ f\ \textbf{where} \\
&\quad fmap \quad :: \quad (a \rightarrow b) \rightarrow (f\ a \rightarrow f\ b) \quad .
\end{aligned}
$$

  This is the first extension we plan to add in the future.

- For a multi-parameter type class there would be multiple type arguments. We do not consider this complication in this paper.

### *4.4 Why our scheme is so beautiful*

We have described a design for adding generic programming to Haskell. We are now in a position to briefly enumerate its merits:

- It fits with the "spirit of Haskell". At first sight, generic programming and Haskell type classes are in competition, but we use generic programming to smoothly extend the power of type classes.
- We are able to explain what "**deriving**" means in a systematic way. The *ad hoc* nature of **deriving** has long been considered a wart, and programmers often want to add new "derivable" classes — that is, classes for which you can say "**deriving** $(C)$". Now they can.
- Generic definitions can be over-ridden at particular types by programmer-supplied instance declarations. This sets our approach apart from other generic programming schemes. Not only is this useful for primitive types, but generic methods are often inapplicable for abstract types — consider equality on sets represented as unordered lists, for example.
- There is no run-time passing or case-analysis of types, beyond Haskell's existing dictionary passing. Of course, dictionary-passing *is* a sort of type passing, but it already exists in Haskell, and it would be extremely tiresome to introduce another, overlapping mechanism.

  Nor are there any new requirements to inspect the run-time representation of a value, a feature of some proposals. Our proposal is a 100% compile-time transformation.
- Like Haskell's type classes, static specialisation is possible to eliminate run-time overhead (see Section 6.6).
- Our scheme deals successfully with constructor names and labels. We have to admit, though, that this is one of the trickier corners of the design.
- We cunningly re-use Haskell's type-class mechanism to instantiate the generic methods for particular types, by expressing the generic methods as generic instance declarations (Section 4.2). This approach means that we do not need to explain or implement exactly how this code instantiation takes place (e.g. how much is done at compile time). Instead we just piggy-back on an existing piece of implementation technology. (This is really a point about the implementation, not about the design.)

## 5 Mapping functions

We now turn our attention to the implementation of our design. As a first step, we discuss so-called *bidirectional maps*. Recall from Section 3.5 our general plan for filling in the generic method of an instance declaration. Suppose we have the following class declaration:

$$\textbf{class } C \ a \ \textbf{where}$$
$$op \quad :: \quad Op \ a \ .$$

We will deal only with single-parameter type classes, but see Section 10. We also assume, for notational clarity, that the type of method *op* is given simply by *Op a*. We can always introduce a type synonym to make this so[1]. Now suppose that the programmer writes the instance declaration

$$\textbf{instance } ctx \Rightarrow C\ (T\ \bar{a})\ .$$

where *ctx* is a context, and $\bar{a}$ is a sequence of type variables. How is the compiler to fill in the definition of method *op*? Following Section 3.5 it can fill in thus:

$$\textbf{instance } ctx \Rightarrow C\ (T\ \bar{a})\ \textbf{where}$$
$$op\ \ =\ \ adapt\text{-}Op\ (op :: Op\ (T^{\circ}\ \bar{a}))\ .$$

That is, we call *op* at type $T^{\circ}\ \bar{a}$, to produce a value of type $Op\ (T^{\circ}\ \bar{a})$, and then convert the value to $Op\ (T\ \bar{a})$. The function *adapt-Op* does this impedance-matching by converting a function that works on values of type $T^{\circ}\ \bar{a}$ to one that works on $T\ \bar{a}$.

$$adapt\text{-}Op\quad::\quad \forall \bar{a}\,.\, Op\ (T^{\circ}\ \bar{a}) \to Op\ (T\ \bar{a})$$

Clearly, how *adapt-Op* works depends on the form of *Op*, the type of the method. Here are some examples:

$$
\begin{array}{lll}
\textbf{type } In\ a & = & a \to String \\
adapt\text{-}In & :: & \forall \bar{a}\,.\, In\ (T^{\circ}\ \bar{a}) \to In\ (T\ \bar{a}) \\
adapt\text{-}In & = & \lambda f \to f \cdot from\text{-}T \\[4pt]
\textbf{type } Out\ a & = & String \to a \\
adapt\text{-}Out & :: & \forall \bar{a}\,.\, Out\ (T^{\circ}\ \bar{a}) \to Out\ (T\ \bar{a}) \\
adapt\text{-}Out & = & \lambda f \to to\text{-}T \cdot f \\[4pt]
\textbf{type } InOut\ a & = & a \to a \\
adapt\text{-}InOut & :: & \forall \bar{a}\,.\, InOut\ (T^{\circ}\ \bar{a}) \to InOut\ (T\ \bar{a}) \\
adapt\text{-}InOut & = & \lambda f \to to\text{-}T \cdot f \cdot from\text{-}T\ .
\end{array}
$$

These *adapt* functions use the functions *to-T* and *from-T*, that convert between $T\ \bar{a}$ and $T^{\circ}\ \bar{a}$; they were introduced in Section 3.3. Notice that both *to-T* and *from-T* are needed; one by itself will not do. Furthermore, while we define the class, and hence the *Op* types, once, we may write instances of that class at many different types, *T*. So we want to abstract out the *to-T* and *from-T* functions from the *adapt* functions (note that $a^{\circ}$ is a type variable below):

$$
\begin{array}{lll}
adapt\text{-}InOut' & :: & \forall a\ a^{\circ}\,.\, (a^{\circ} \to a) \to (a \to a^{\circ}) \to (InOut\ a^{\circ} \to InOut\ a) \\
adapt\text{-}InOut'\ to\ from & = & \lambda f \to to \cdot f \cdot from \\
adapt\text{-}InOut & = & adapt\text{-}InOut'\ to\text{-}T\ from\text{-}T
\end{array}
$$

---

[1] Technically, Haskell type synonyms are not powerful enough to do such an abbreviation for a method like *properFraction*:

$$\textbf{class } RealFrac\ a\ \textbf{where}$$
$$properFraction \quad :: \quad (Integral\ b) \Rightarrow a \to (b, a)\ .$$

since its type has a context at the beginning. But we will pretend that such an abbreviation is possible.

It turns out to be convenient to package up the *to-T* and *from-T* functions into an *embedding-projection pair*:

$$\textbf{data } EP\ a\ a^\circ \quad = \quad EP\{\,to :: a^\circ \to a, from :: a \to a^\circ\,\}\ .$$

Here $EP\ a\ a^\circ$ is just a pair of functions, one to convert in one direction and one to convert back. Now we can write

$$
\begin{array}{lll}
conv\text{-}T & :: & \forall\bar{a}\,.\,EP\ (T\ \bar{a})\ (T^\circ\ \bar{a}) \\
conv\text{-}T & = & EP\{\,to = to\text{-}T, from = from\text{-}T\,\} \\[4pt]
adapt\text{-}InOut'' & :: & \forall a\ a^\circ\,.\,EP\ a\ a^\circ \to InOut\ a^\circ \to InOut\ a \\
adapt\text{-}InOut''\ ep\text{-}a & = & \lambda f \to to\ ep\text{-}a \cdot f \cdot from\ ep\text{-}a \\[4pt]
adapt\text{-}InOut & = & adapt\text{-}InOut''\ conv\text{-}T
\end{array}
$$

The last step is to make the *adapt* function itself return an embedding-projection pair, rather than just the "to" function; and at this stage we adopt the name *bimap* — for "bidirectional mapping":

$$
\begin{array}{lll}
bimap\text{-}InOut & :: & \forall a\ a^\circ\,.\,EP\ a\ a^\circ \to EP\ (InOut\ a)\ (InOut\ a^\circ) \\
bimap\text{-}InOut\ ep\text{-}a & = & EP\{\,to = \lambda f \to from\ ep\text{-}a \cdot f \cdot to\ ep\text{-}a, \\
& & \qquad\quad from = \lambda f \to to\ ep\text{-}a \cdot f \cdot from\ ep\text{-}a\,\} \\[4pt]
adapt\text{-}InOut & = & to\ (bimap\text{-}InOut\ conv\text{-}T)\ .
\end{array}
$$

It is not at all obvious why we construct mappings in both directions, only to discard one of them when we use it, but it turns out to be essential when constructing *bimap* for arbitrary types, as we will see in the next section.

### 5.1  Generating bidirectional mapping functions

In the last section we generated *bimap-InOut* for a particular method type *InOut a*. We also observed that appropriate code depends on the structure of the type of the method, so the million-dollar question is: how do we generate the bidirectional maps for *arbitrary* method types? We do it simply by induction over the type of the method, thus:

$$
\begin{array}{lll}
bimap\text{-}Op & :: & \forall a\ a^\circ\,.\,EP\ a\ a^\circ \to EP\ (Op\ a)\ (Op\ a^\circ) \\
bimap\text{-}Op\ ep\text{-}a & = & bimap\{\,Op\ a\,\}[a := ep\text{-}a]\ .
\end{array}
$$

This definition is not proper Haskell; *bimap* should be thought of as a meta-function, evaluated at compile time, that returns a Haskell expression. It takes as arguments: a type (written in curly braces), and an environment $\varrho$ mapping type variables to expressions. The syntax $[a := ep\text{-}a]$ means an environment that binds $a$ to *ep-a*.

We define *bimap* by induction on the structure of type expressions:

$$
\begin{array}{lll}
bimap\{\,a\,\}\varrho & = & \varrho(a) \\
bimap\{(\to)\}\varrho & = & bimap\text{-}Arrow \\
bimap\{\,T\,\}\varrho & = & bimap\text{-}T \\
bimap\{\,t\ u\,\}\varrho & = & (bimap\{\,t\,\}\varrho)\ (bimap\{\,u\,\}\varrho)
\end{array}
$$

where

$$bimap\text{-}Arrow \quad :: \quad \forall a\ a^\circ\ b\ b^\circ.\ EP\ a\ a^\circ \to EP\ b\ b^\circ \to EP\ (a \to b)\ (a^\circ \to b^\circ)$$

$$bimap\text{-}Arrow\ ep\text{-}a\ ep\text{-}b$$
$$= \quad EP\{\,to = \lambda f \to to\ ep\text{-}b \cdot f \cdot from\ ep\text{-}a,$$
$$from = \lambda f \to from\ ep\text{-}b \cdot f \cdot to\ ep\text{-}a\,\}\ .$$

The $bimap\{\,T\,\}$ case deals with type constructors other than $(\to)$, which we discuss in Section 5.2. Let us take an example. Recall our $InOut$ type:

$$\textbf{type}\ InOut\ a \quad = \quad a \to a\ .$$

Setting $\varrho = [\,a := ep\text{-}a\,]$ we have

$$\begin{aligned}
bimap\text{-}InOut\ ep\text{-}a \quad &= \quad bimap\{\,a \to a\,\}\varrho \\
&= \quad (bimap\{\,(\to)\ a\,\}\varrho)\ (bimap\{\,a\,\}\varrho) \\
&= \quad ((bimap\ (\to)\ \varrho)\ (bimap\{\,a\,\}\varrho))\ (bimap\{\,a\,\}\varrho) \\
&= \quad bimap\text{-}Arrow\ ep\text{-}a\ ep\text{-}a \\
&= \quad EP\{\,to = \lambda f \to from\ ep\text{-}a \cdot f \cdot to\ ep\text{-}a, \\
&\qquad\qquad from = \lambda f \to to\ ep\text{-}a \cdot f \cdot from\ ep\text{-}a\,\}\ .
\end{aligned}$$

### 5.2 Mapping over data types

What if there is a data type involved? For example in the type of $readBin$, there is a pair in the result type:

$$\textbf{type}\ ReadBin\ a \quad = \quad Bin \to (a, Bin)\ .$$

If we just try our current scheme we get stuck:

$$\begin{aligned}
bimap\text{-}ReadBin\ ep \quad &= \quad bimap\{\,Bin \to (a, Bin)\,\}\varrho \\
&= \quad bimap\text{-}Arrow\ (bimap\{\,Bin\,\}\varrho)\ (bimap\{\,(a, Bin)\,\}\varrho)\ .
\end{aligned}$$

Since $Bin$ is not a parameterised type, there is nothing to do,

$$\begin{aligned}
bimap\text{-}Bin \quad &:: \quad EP\ Bin\ Bin \\
bimap\text{-}Bin \quad &= \quad id\text{-}EP \\
\\
id\text{-}EP \quad &:: \quad \forall a.\ EP\ a\ a \\
id\text{-}EP \quad &= \quad EP\{\,to = \lambda x \to x, from = \lambda x \to x\,\}
\end{aligned}$$

whereas pairs are parameterised over two types, so we must push the mapping functions inside:

$$bimap\text{-}Pair \quad :: \quad \forall a\ a^\circ\ b\ b^\circ.\ EP\ a\ a^\circ \to EP\ b\ b^\circ \to EP\ (a, b)\ (a^\circ, b^\circ)$$

$$bimap\text{-}Pair\ ep\text{-}a\ ep\text{-}b$$
$$= \quad EP\{\,to = \lambda(x^\circ, y^\circ) \to (to\ ep\text{-}a\ x^\circ, to\ ep\text{-}b\ y^\circ),$$
$$from = \lambda(x, y) \to (from\ ep\text{-}a\ x, from\ ep\text{-}b\ y)\,\}\ .$$

In general, we can define *bimap-T* by induction on the structure of datatype declarations. The mapping function for the data type

$$\textbf{data } T\ a_1\ \ldots\ a_k\quad =\quad K_1\ t_{11}\ \ldots\ t_{1\,m_1}\ |\ \cdots\ |\ K_n\ t_{n1}\ \ldots\ t_{nm_n}$$

is given by

$$
\begin{array}{ll}
\textit{bimap-T bimap-a}_1\ \ldots\ \textit{bimap-a}_k\quad =\quad EP\{\,to = \textit{to-T}, \textit{from} = \textit{from-T}\,\} \\
\quad\textbf{where} \\
\quad\textit{to-T}\ (K_1\ x_{11}\ \ldots\ x_{1\,m_1}) & =\quad K_1\ (to\{\,t_{11}\,\}x_{11})\ \ldots\ (to\{\,t_{1\,m_1}\,\}x_{1\,m_1}) \\
\quad\ldots \\
\quad\textit{to-T}\ (K_n\ x_{n1}\ \ldots\ x_{nm_n}) & =\quad K_n\ (to\{\,t_{n1}\,\}x_{n1})\ \ldots\ (to\{\,t_{1\,m_1}\,\}x_{nm_n}) \\
\quad\textit{from-T}\ (K_1\ x_{11}\ \ldots\ x_{1\,m_1}) & =\quad K_1\ (from\{\,t_{11}\,\}x_{11})\ \ldots\ (from\{\,t_{1\,m_1}\,\}x_{1\,m_1}) \\
\quad\ldots \\
\quad\textit{from-T}\ (K_n\ x_{n1}\ \ldots\ x_{nm_n}) & =\quad K_n\ (from\{\,t_{n1}\,\}x_{n1})\ \ldots\ (from\{\,t_{1\,m_1}\,\}x_{nm_n})
\end{array}
$$

where

$$
\begin{array}{rcl}
to\{\,t\,\} & = & to\ (bimap\{\,t\,\}[a_1 := \textit{bimap-a}_1, \ldots, a_k := \textit{bimap-a}_k]) \\
from\{\,t\,\} & = & from\ (bimap\{\,t\,\}[a_1 := \textit{bimap-a}_1, \ldots, a_k := \textit{bimap-a}_k])\ .
\end{array}
$$

Now, what is type of this bidirectional map? The answer is simple yet mind-boggling: the type of *bimap-T* depends on the kind of *T*. Assume that *T* has kind $\star$ as, for instance, *Bin*. Then the bidirectional map simply has type *EP T T* (and, in fact, *bimap-T = id-EP*). If *T* has kind $\star \to \star$ as all the *Op*'s have, then *bimap-T*'s type is close to that of an "ordinary" mapping function:

$$\textit{bimap-T}\quad ::\quad \forall a\ a^\circ .\ EP\ a\ a^\circ \to EP\ (T\ a)\ (T\ a^\circ)\ .$$

A more involved kind, say $(\star \to \star) \to (\star \to \star)$, gives rise to a more complicated type:

$$
\begin{array}{rcl}
\textit{bimap-T}\quad :: & \forall f\ f^\circ .\ (\forall b\ b^\circ .\ EP\ b\ b^\circ \to EP\ (f\ b)\ (f^\circ\ b^\circ)) \\
& \to (\forall a\ a^\circ .\ EP\ a\ a^\circ \to EP\ (T\ f\ a)\ (T\ f^\circ\ a^\circ))\ .
\end{array}
$$

Now *bimap-T* has a so-called rank-2 type signature (Leivant, 1983). Roughly speaking, *bimap-T* takes bidirectional maps to bidirectional maps (this is why the arguments of *bimap-T* are called *bimap-a_i*). In general, if *T* has kind $\kappa$, then *bimap-T* has type $Bimap\{\kappa\}\ T\ T$ where *Bimap* is defined by induction on the structure of kinds:

$$
\begin{array}{rcl}
Bimap\{\star\}t\ t^\circ & = & EP\ t\ t^\circ \\
Bimap\{\kappa_1 \to \kappa_2\}t\ t^\circ & = & \forall a\ a^\circ .\ Bimap\{\kappa_1\}a\ a^\circ \to Bimap\{\kappa_2\}(t\ a)\ (t^\circ\ a^\circ)\ .
\end{array}
$$

We will say a bit more about higher-order kinded types in Section 7. For further information on kind-indexed types such as *Bimap* the reader if referred to (Hinze, 2000c).

# 6 Implementing generic default methods

In this section we discuss the implementation of our design.

We describe our implementation as a Haskell source-to-source translation, performed (at least notionally) prior to type checking. Why? The type checker already does a lot of what we require. Also we probably have a better chance that generic default methods will work smoothly with complications such as multi-parameter type classes (Peyton Jones *et al.*, 1997), implicit parameters (Lewis *et al.*, 2000), and functional dependencies (Jones, 2000).

The source-to-source translation goes as follows. For each *datatype declaration*, $T$, we generate the following:

- For each constructor $K$ a value *con-K* of type *ConDescr* that describes the properties of the constructor (Section 6.1).
- A type synonym, $T^\circ$, for $T$'s generic representation type (Section 6.2).
- An embedding-projection pair *conv-T* :: $\forall \bar{a} . EP$ $(T$ $\bar{a})$ $(T^\circ$ $\bar{a})$, that converts between $T$ and its generic representation $T^\circ$ (Section 6.3).

For each *class declaration*, for class $C$, we generate the following (see Section 6.4):

- An emaciated class declaration for $C$, generated simply by omitting the generically-defined methods.
- For each generic method *op* :: $Op$ $a$ in the class declaration, a bidirectional map *bimap-Op* :: $\forall a$ $a^\circ . EP$ $a$ $a^\circ \to EP$ $(Op$ $a)$ $(Op$ $a^\circ)$ (see Section 5).
- Instance declarations for $C$ 1, $C$ $(a+b)$, $C$ $(a \times b)$ and $C$ $(c$ **of** $a)$, all obtained by selecting the eponymous equations from the original class declaration (see Section 3.4).

For each *instance declaration* we generate (see Section 6.5):

- An extended instance declaration, obtained by adding definitions for the generic methods that are not specified explicitly in the instance declaration.

### 6.1 Constructors

For each constructor, $K$, in a datatype declaration, we produce a value of type *ConDescr* that gives information about the constructor (in fact, the type *ConDescr* used in the compiler is slightly more elaborate):

**data** $ConDescr$ $=$ $ConDescr\{name :: String, arity :: Int, fixity :: Fixity\}$ .

As an example, for the *List* datatype we generate:

$$con\text{-}Cons, con\text{-}Nil \quad :: \quad ConDescr$$
$$con\text{-}Cons \quad = \quad ConDescr \text{ "Cons" } 2 \text{ } NoFixity$$
$$con\text{-}Nil \quad = \quad ConDescr \text{ "Nil" } 0 \text{ } NoFixity \text{ .}$$

## *6.2  Generic representation types*

For each data type, $T$, we produce a type synonym $T°$, for its *generic representation type*. For example, for the data type

$$\textbf{data } \textit{List a} \quad = \quad \textit{Cons a (List a)} \mid \textit{Nil}$$

we generate the representation type

$$\textbf{type } \textit{List}° \textit{ a} \quad = \quad (\textit{con-Cons } \textbf{of } (a \times \textit{List a})) + (\textit{con-Nil } \textbf{of } 1)$$

Our generic representation type constructors are just sum, product, unit and "$c$ **of**". In particular, there is no recursion operator. Thus, we observe that $\textit{List}°$ is just a non-recursive type synonym: $\textit{List}$ (not $\textit{List}°$) appears on the right-hand side. So $\textit{List}°$ is not a recursive type; rather, it expresses just the top "layer" of a list structure, leaving the original $\textit{List}$ to do the rest. But as we have seen, this is enough: a recursive function just does one "layer" of recursion at a time.

This is unusual compared to other approaches. In PolyP (Jansson & Jeuring, 1997), for instance, there is an additional type pattern for type recursion (at kind $\star \rightarrow \star$). A very significant advantage here is that there is no problem with mutually-recursive data types, nor with data types with many parameters, both of which make explicit recursion operators extremely clumsy and hard to use in practice.

Our design makes do with just *binary* sum and product. Algebraic data types with many constructors, each of which has many fields, are encoded as nested uses of sum and product. The exact way in which the nesting is done is unimportant to our method. For example:

$$
\begin{array}{lll}
\textbf{data } \textit{Color} & = & \textit{Red} \mid \textit{Blue} \mid \textit{Green} \\
\textbf{type } \textit{Color}° & = & \textit{con-Red } \textbf{of } 1 + (\textit{con-Blue } \textbf{of } 1 + \textit{con-Green } \textbf{of } 1) \\
\textbf{data } \textit{Tree a b} & = & \textit{Leaf a} \mid \textit{Node (Tree a b) b (Tree a b)} \\
\textbf{type } \textit{Tree}° \textit{ a b} & = & \textit{con-Leaf } \textbf{of } a + (\textit{con-Node } \textbf{of } (\textit{Tree a b} \times (b \times \textit{Tree a b})))
\end{array}
$$

One may wonder about the efficiency of translating a user-defined data type into a generic form before operating on it, especially if everything is encoded with only binary sums and products. However, sufficiently vigorous inlining means that the generic data representations never exist at run-time (see Section 6.6). But, in fact, we might want to explore space-time trade-offs, by getting much more compact code in exchange for some data translation. Our design allows this trade-off to be made on a case-by-case basis.

Whether the encoding into sums and products is done in a linear or binary-sub-division fashion may or may not affect efficiency, depending on how vigorous the inlining is.

### *6.3 Embedding-projection pairs*

For each datatype $T$, we also generate functions to convert between $T$ and $T^\circ$. We saw the conversion functions for *List* in Section 3.3. The process is entirely straightforward, driven by the encoding. For example:

$$
\begin{array}{lll}
\textit{from-Color} & :: & \textit{Color} \to \textit{Color}^\circ \\
\textit{from-Color Red} & = & \textit{Inl (Con con-Red Unit)} \\
\textit{from-Color Blue} & = & \textit{Inr (Inl (Con con-Blue Unit))} \\
\textit{from-Color Green} & = & \textit{Inr (Inr (Con con-Green Unit))}
\end{array}
$$

$$
\begin{array}{lll}
\textit{to-Color} & :: & \textit{Color}^\circ \to \textit{Color} \\
\textit{to-Color (Inl (Con con-Red Unit))} & = & \textit{Red} \\
\textit{to-Color (Inr (Inl (Con con-Blue Unit)))} & = & \textit{Blue} \\
\textit{to-Color (Inr (Inr (Con con-Green Unit)))} & = & \textit{Green} \ .
\end{array}
$$

For *bimap* we have to package the two conversion functions into a single value:

$$
\begin{array}{lll}
\textit{conv-List} & :: & \forall a . \textit{EP (List a) (List}^\circ\textit{ a)} \\
\textit{conv-List} & = & \textit{EP}\{\textit{to} = \textit{to-List}, \textit{from} = \textit{from-List}\} \\
\\
\textit{conv-Color} & :: & \textit{EP Color Color}^\circ \\
\textit{conv-Color} & = & \textit{EP}\{\textit{to} = \textit{to-Color}, \textit{from} = \textit{from-Color}\} \ .
\end{array}
$$

### *6.4 Translating class declarations*

For each generic method $op :: Op\ a$ contained in a class declaration we generate a bidirectional map

$$
\textit{bimap-Op} \quad :: \quad \forall a\ a^\circ . \textit{EP a a}^\circ \to \textit{EP (Op a) (Op a}^\circ)
$$

that allows us to convert between types and representation types (the definition of *bimap-Op* is given in Section 5).

Furthermore, we produce instance declarations

$$
\begin{array}{l}
\textbf{instance } C\ 1 \\
\textbf{instance } (C\ a, C\ b) \Rightarrow C\ (a + b) \\
\textbf{instance } (C\ a, C\ b) \Rightarrow C\ (a \times b) \\
\textbf{instance } (C\ a) \Rightarrow C\ (c\ \textbf{of}\ a)
\end{array}
$$

whose bodies are filled with the generic methods from the original class declaration (see Section 3.4). If an equation for a type pattern is missing, the method of the corresponding instance is undefined. There is, however, one important exception to this rule: if no equation is given for the type pattern $c\ \textbf{of}\ a$ as, for example, in the classes *Eq* and *Binary*, we define the generic methods of the $C\ (c\ \textbf{of}\ a)$ instance by:

$$
op\{\,c\ \textbf{of}\ a\,\} \quad = \quad to\ (\textit{bimap-Op (con-EP c)})\ (op :: Op\ a)
$$

where *con-EP* is given by the following pseudo-Haskell code (which defines a family

of functions):

$$con\text{-}EP \ c \quad :: \quad \forall a \, . \, EP \ a \ (c \ \textbf{of} \ a)$$
$$con\text{-}EP \ c \quad = \quad EP\{ to = \lambda x \to Con \ c \ x, from = \lambda(Con \ c \ x) \to x \} \ .$$

Again, we employ the bidirectional map to convert between two isomorphic types.

### 6.5 Translating instance declarations

An instance declaration for type $T$ is extended by filling in implementations for the methods. More specifically, if the method *op* is not specified and if it enjoys a generic default definition, then the following equation is supplemented:

$$op \quad = \quad to \ (bimap\text{-}Op \ conv\text{-}T) \ (op :: Op \ T^\circ) \ .$$

That's it.

### 6.6 Inlining

It does not sound very efficient to translate a value from $T \ \bar{a}$ to $T^\circ \ \bar{a}$ and then to operate on it, but we believe that a bit of judicious inlining can yield more or less the code one would have written by hand. Let us consider, for example, *showBin* at type *List*. The *showBinList* method will look something like this:

$$showBinList \qquad :: \quad (Binary \ a) \Rightarrow List \ a \to Bin$$
$$showBinList \ xs \quad = \quad showBin \ (from\text{-}List \ xs)$$
$$\textbf{type} \ List^\circ \ a \qquad = \quad (a \times List \ a) + 1$$
$$from\text{-}List \qquad \qquad :: \quad List \ a \to List^\circ \ a \ .$$

The call to *showBin* is at type $List^\circ \ a$, so the overloading can be resolved statically. Assuming that the method bodies (given in Section 3.1) are inlined, we get:

$$showBinList \ xs \quad = \quad \textbf{case} \ from\text{-}List \ xs \ \textbf{of}$$
$$Inl \ z \to 0 : \textbf{case} \ z \ \textbf{of}$$
$$(x :*: y) \to showBin \ x \ {+\!\!+}\ showBin \ y$$
$$Inr \ z \to 1 : \textbf{case} \ z \ \textbf{of} \ Unit \to [\,] \ .$$

But remember that *from-List* also has a simple, non-recursive definition:

$$from\text{-}List \ (Cons \ x \ xs) \quad = \quad Inl \ (x :*: xs)$$
$$from\text{-}List \ Nil \qquad \qquad = \quad Inr \ Unit \ .$$

If we inline this definition in *showBinList* and simplify using standard transformations, we get

$$showBinList \ xs \quad = \quad \textbf{case} \ xs \ \textbf{of}$$
$$Cons \ x \ y \to 0 : showBin \ x \ {+\!\!+}\ showBin \ y$$
$$Nil \to 1 : [\,]$$

which is about as good as we can hope for.

## 7 Higher-order kinded types

Functional programmers love abstraction. In Haskell we can, for instance, abstract over the *List* datatype in

$$\textbf{data } Rose\ a\ \ =\ \ Branch\ a\ (List\ (Rose\ a))$$

obtaining the more general type:

$$\textbf{data } GRose\ f\ a\ \ =\ \ GBranch\ a\ (f\ (GRose\ f\ a))\ .$$

Here, the type variable "*f*" ranges over type constructors, rather than types. Formally, *GRose* has kind $(\star \to \star) \to \star \to \star$. There are numerous examples of such type definitions in (Okasaki, 1998; Hinze, 2000a). Alas, it is impossible to define many instance declarations for *GRose* in Haskell at all. In this section we describe the problem and a solution. This section is quite independent of the rest of the paper. Though we became aware of the issue when working on generic programming, we propose an extension to Haskell that is completely orthogonal to generic programming.

### 7.1 What's the problem?

Consider first defining an instance for *Binary* (*Rose* *a*) by hand — we ignore *readBin* here:

$$\textbf{instance } (Binary\ a) \Rightarrow Binary\ (Rose\ a)\ \textbf{where}$$
$$showBin\ (Branch\ x\ ts)\ \ =\ \ showBin\ x \mathbin{+\!\!+} showBin\ ts\ .$$

The first call to *showBin* on the right hand side requires that *Binary* *a* should hold; the context, (*Binary* *a*), takes care of that. The second call is at type *List* (*Rose* *a*). Assuming we have an instance elsewhere of the form

$$\textbf{instance } (Binary\ t) \Rightarrow Binary\ (List\ t)$$

the second call requires *Binary* (*Rose* *a*), and there is an instance declaration for that too — it gives rise to a recursive call to *showBin*.

But matters are not so simple when we want to write an instance *Binary* (*GRose* *f* *a*). We might try

$$\textbf{instance } (???) \Rightarrow Binary\ (GRose\ f\ a)\ \textbf{where}$$
$$showBin\ (GBranch\ x\ ts)\ \ =\ \ showBin\ x \mathbin{+\!\!+} showBin\ ts\ .$$

The context (???) must account for the calls to *showBin* on the right-hand side. The first one is fine: it requires *Binary* *a* as before. But the latter is bad news: it requires *Binary* (*f* (*GRose* *f* *a*)), and we certainly cannot write

$$\textbf{instance } (Binary\ a, Binary\ (f\ (GRose\ f\ a))) \Rightarrow Binary\ (GRose\ f\ a)\ \textbf{where}$$
$$showBin\ (GBranch\ x\ ts)\ \ =\ \ showBin\ x \mathbin{+\!\!+} showBin\ ts\ .$$

This is not legal Haskell and, even if it were, the typechecker would diverge. Indeed, no ordinary Haskell context will do.

## 7.2  A solution

What we need is a way to simplify the predicate $f$ (*GRose f a*). The trick is to take the "constant" instance declaration that we assumed for *Binary* (*List a*) above, and abstract over it:

**instance** (*Binary a*, $\forall b$ . (*Binary b*) $\Rightarrow$ *Binary* (*f b*)) $\Rightarrow$ *Binary* (*GRose f a*) **where**
    *showBin* (*GBranch x ts*)   =   *showBin x* ++ *showBin ts* .

Now, as well as (*Binary a*), the context also contains a *polymorphic predicate*. This predicate can be used to reduce the predicate *Binary* (*f* (*GRose f a*)) to just *Binary* (*GRose f a*), and we have an instance declaration for that.

Viewed in operational terms, the predicate (*Binary a*) in a context corresponds to passing a *dictionary* for class *Binary*. A predicate $\forall b$ . *Binary b* $\Rightarrow$ *Binary* (*f b*) corresponds to passing a *dictionary transformer* to the function.

## 7.3  Deriving instance declarations

Of course, since *Binary* is a derivable class by virtue of the generic default definitions, we need not define *showBin* at all. We can simply write

    **instance** (*Binary a*, $\forall b$ . (*Binary b*) $\Rightarrow$ *Binary* (*f b*)) $\Rightarrow$ *Binary* (*GRose f a*)

and get just the same effect as before. In other words, the deriving mechanism works happily for types of arbitrary kinds.

Can we use a **deriving** clause to get the same effect? Unfortunately, the answer is in the negative. The problem is that the above instance declaration is not the most general one and, in fact, there is no "most general instance declaration". To illustrate the point consider the following instance declaration for the abstract type *Set*:

        **instance** (*Binary a*, *Ord a*) $\Rightarrow$ *Binary* (*Set a*) .

Note that we additionally require that $a$ is an instance of *Ord*. Now, given the instance declaration for *GRose* above, we cannot infer *Binary* (*GRose Set Int*) since *Set* does not satisfy $\forall b$ . (*Binary b*) $\Rightarrow$ *Binary* (*Set b*). If we require such an instance, we must *generalize* the *GRose* instance:

**instance** (*Binary a*, $\forall b$ . (*Binary b*, *Ord b*) $\Rightarrow$ *Binary* (*f b*)) $\Rightarrow$ *Binary* (*GRose f a*).

By adding further class constraints to $f$'s context, we can generalize the instance declaration even more. Sadly, this implies that there is no "most general" instance which **deriving** could infer. Note that this problem does not crop up for first-order kinded types.

### 7.4 Formalising the extension

Here is the grammar for generalized instance declarations:

$$
\begin{array}{llll}
\text{instance head} & ::= & \textbf{instance } (ctx_1, \ldots, ctx_n) \Rightarrow C\ t \\
\text{context} & ctx & ::= & \forall \bar{a} . (ctx_1, \ldots, ctx_n) \Rightarrow C\ t \ .
\end{array}
$$

A context of the form $\forall \bar{a} . (ctx_1, \ldots, ctx_n) \Rightarrow C\ t$ with $n \geqslant 1$ is called a *polymorphic predicate*. Note that for $n = 0$ we have "ordinary" Haskell 98 predicates.

## 8 Implementing generalized instance declarations

How do we translate a method call $op :: Op\ T$? We must create a $C$-dictionary for $T$ if $op$ is a method of class $C$. In the higher-order kinded situation, we may need to create a dictionary *transformer* to pass to $op$. Fortunately, it turns out that the now-standard machinery to construct the correct dictionary to pass can easily be extended to construct dictionary transformers too.

At a call site we have to solve the following problem: we have a set of *assumptions* $\mathbb{H}$ and a single clause $H$, the dictionary (transformer) required, and we want to know whether $H$ is a logical consequence of $\mathbb{H}$. Additionally we return an expression for the dictionary (transformer) for $H$. We use the following notation: $\mathbb{H} \vdash H \longmapsto d$ means that $d$ is a dictionary (transformer) expression that shows how $H$ can be deduced from $\mathbb{H}$.

The assumptions $\mathbb{H}$ embody:

- Any instance declarations in scope. For example:

$$
\begin{array}{c}
Eq\ Int \longmapsto dict\text{-}Eq\text{-}Int \\
\forall a . Eq\ a \Rightarrow Eq\ (List\ a) \longmapsto dict\text{-}Eq\text{-}List \quad .
\end{array}
$$

- Information about superclasses. For example:

$$
\forall a . Ord\ a \Rightarrow Eq\ a \longmapsto dict\text{-}Eq\text{-}Ord \ .
$$

  This says that if we have $Ord\ a$ we can deduce $Eq\ a$; in concrete terms we witness this fact by the selector function $dict\text{-}Eq\text{-}Ord$ which selects the $Eq$ dictionary from the $Ord$ one.
- Constraints from the type signature. For example, if we are checking types for the function

$$
\begin{array}{lll}
f & :: & \bar{H} \Rightarrow T \\
f\ x & = & \ldots
\end{array}
$$

  then we put the assumptions $\bar{H}$ in our assumption set, and try to deduce all the dictionaries that are needed by calls in the body of $f$.

We use the following inference rules:

$$
\frac{(H \longmapsto d) \in \mathbb{H}}{\mathbb{H} \vdash H \longmapsto d} \qquad (\textsc{Assumption})
$$

$$\frac{\mathbb{H} \vdash H_1 \longmapsto d_1 \qquad \cdots \qquad \mathbb{H} \vdash H_n \longmapsto d_n}{\mathbb{H} \vdash (H_1, \ldots, H_n) \longmapsto (d_1, \ldots, d_n)} \quad (\text{CONJUNCTION})$$

$$\frac{\mathbb{H} \vdash (\forall \bar{a} \,.\, \bar{H} \Rightarrow Q) \longmapsto f \qquad \mathbb{H} \vdash \bar{H}\varrho\theta \longmapsto d}{\mathbb{H} \vdash P \longmapsto (f\ d)} \quad (\text{MODUS PONENS})$$

where $\varrho = [\bar{a} := \bar{x}]$ is a renaming substitution (the $x_i$ are fresh variables) and $\theta = match(Q\varrho, P)$ is the result of matching $Q\varrho$ against $P$ (note that only the variables in $Q\varrho$ are bound).

So far, these rules are entirely standard, see, for instance, (Jones, 1994). To these we add one new rule.

$$\frac{\mathbb{H}, (\bar{H}\varrho \longmapsto v) \vdash Q\varrho \longmapsto d}{\mathbb{H} \vdash (\forall \bar{a} \,.\, \bar{H} \Rightarrow Q) \longmapsto (\lambda v \to d)} \quad (\text{DEDUCTION RULE})$$

where $\varrho = [\bar{a} := \bar{c}]$ is a Skolem substitution, that is, the $c_i$ are Skolem constants. Thus, to deduce the polymorphic predicate $\forall \bar{a} \,.\, \bar{H} \Rightarrow Q$ we add the body $\bar{H}$ to the set of assumptions and try to deduce $Q$. The Skolem substitution ensures that this derivation works for all $\bar{a}$.

The new rule is called DEDUCTION RULE because it resembles the deduction theorem of first-order logic. It is also reminiscent of the usual typing rule for $\lambda$-abstraction while MODUS PONENS corresponds to the typing rule for application. In fact, these two rules capture dictionary abstraction and dictionary application.

Here is an example of a deduction using these rules. Later lines are deduced from earlier ones using the specified rule (A for ASSUMPTION, etc).

$$\begin{aligned}
\mathbb{H} \quad = \quad &\{ \text{Ord Int} \longmapsto \text{d-O-I}, \\
&\ \ \text{Binary Int} \longmapsto \text{d-B-I}, \\
&\ \ (\forall a \,.\, (\text{Binary } a, \text{Ord } a) \Rightarrow \text{Binary } (\text{Set } a)) \longmapsto \text{d-B-S} \}
\end{aligned}$$

| | | |
|---|---|---|
| (5) | $\mathbb{H} \vdash$ *Ord Int* $\longmapsto$ *d-O-I* | A |
| (4) | $\mathbb{H} \vdash$ *Binary Int* $\longmapsto$ *d-B-I* | A |
| (3) | $\mathbb{H} \vdash$ *(Binary Int, Ord Int)* $\longmapsto$ *(d-B-I, d-O-I)* | C(4,5) |
| (2) | $\mathbb{H} \vdash$ *($\forall a \,.\, $(Binary a, Ord a) $\Rightarrow$ Binary (Set a))* $\longmapsto$ *d-B-S* | A |
| (1) | $\mathbb{H} \vdash$ *Binary (Set Int)* $\longmapsto$ *d-B-S (d-B-I, d-O-I)* | MP(2,3) |

Here is another example, this time of a higher-order case:

$$\begin{aligned}
\mathbb{H} \quad = \quad &\{ \text{Binary Int} \longmapsto \text{d-B-I}, \\
&\ \ (\forall a \,.\, (\text{Binary } a) \Rightarrow \text{Binary } (\text{List } a)) \longmapsto \text{d-B-L}, \\
&\ \ (\forall f\ a \,.\, (\text{Binary } a, \forall b \,.\, (\text{Binary } b) \Rightarrow \text{Binary } (f\ b)) \\
&\qquad\qquad \Rightarrow \text{Binary } (\text{GRose } f\ a)) \longmapsto \text{d-B-G} \} \ .
\end{aligned}$$

In the following we abbreviate *Binary* by *B*.

(10)  $\mathbb{H}, (B\ c \longmapsto v) \vdash (\forall b\,.\,(B\ b) \Rightarrow B\ (List\ b)) \longmapsto d\text{-}B\text{-}L$          A

(9)   $\mathbb{H}, (B\ c \longmapsto v) \vdash B\ c \longmapsto v$          A

(8)   $\mathbb{H}, (B\ c \longmapsto v) \vdash (B\ c, \forall b\,.\,(B\ b) \Rightarrow B\ (List\ b)) \longmapsto (v, d\text{-}B\text{-}L)$   C(9,10)

(7)   $\mathbb{H}, (B\ c \longmapsto v) \vdash (\forall f\ a\,.\,(\ldots) \Rightarrow B\ (GRose\ f\ a)) \longmapsto d\text{-}B\text{-}G$   A

(6)   $\mathbb{H}, (B\ c \longmapsto v) \vdash B\ (GRose\ List\ c) \longmapsto d\text{-}B\text{-}G\ (v, d\text{-}B\text{-}L)$   MP(7,8)

(5)   $\mathbb{H} \vdash (\forall b\,.\,(B\ b) \Rightarrow B\ (GRose\ List\ b))$
       $\longmapsto (\lambda v \to d\text{-}B\text{-}G\ (v, d\text{-}B\text{-}L))$          DR(6)

(4)   $\mathbb{H} \vdash B\ Int \longmapsto d\text{-}B\text{-}I$          A

(3)   $\mathbb{H} \vdash (B\ Int, \forall b\,.\,(B\ b) \Rightarrow B\ (GRose\ List\ b))$
       $\longmapsto (d\text{-}B\text{-}I, \lambda v \to d\text{-}B\text{-}G\ (v, d\text{-}B\text{-}L))$          C(4,5)

(2)   $\mathbb{H} \vdash (\forall f\ a\,.\,(\ldots) \Rightarrow B\ (GRose\ f\ a)) \longmapsto d\text{-}B\text{-}G$   A

(1)   $\mathbb{H} \vdash B\ (GRose\ (GRose\ List)\ Int)$
       $\longmapsto d\text{-}B\text{-}G\ (d\text{-}B\text{-}I, \lambda v \to d\text{-}B\text{-}G\ (v, d\text{-}B\text{-}L))$          MP(2,3)

The new inference rule kicks in at line (5) and introduces a new assumption, $B\ c \longmapsto v$, that is used in line (2).

# 9  Related work

This paper improves on our earlier work (Hinze, 1999) in several respects. First, generic definitions now appear solely in class declarations as generic default methods. In the previous design generic definitions and classes were two competing features. We feel that the new proposal fits better with "the spirit of Haskell". Second, we have spelled out the implementation in considerable detail. Especially, the notion of generic representation types and the conversion between types and representation types has been made precise. Third, we have described a separate extension that allows the programmer to define instance declarations for higher-order kinded types. The need for this extension was noted in (Hinze, 1999) but no solution was given.

Though there is a considerable amount of work on generic programming (Ruehr, 1992; Cockett & Fukushima, 1992; Jay *et al.*, 1998) this is the first paper we are aware of — apart from PolyP (Jansson & Jeuring, 1997) — that aims at adding generic features to an *existing* functional language. The PolyP extension offers a special construct (essentially, a type case) for defining generic functions. The resulting definitions are similar to ours (modulo notation) except that the generic programmer must additionally consider cases for type composition and for type recursion. Furthermore, PolyP is restricted to regular datatypes of kind $\star \to \star$, whereas our proposal works for all types of all kinds. This is quite a significant advantage. In particular, our proposal deals gracefully with mutually-recursive data types and with data types with many parameters, both of which make explicit recursion operators clumsy and hard to use in practice.

## 10 Conclusions and further work

This paper describes two separate extensions to Haskell. The first extension supports generic programming through a new form of default method declaration. The second extension allows one to define instance declarations for higher-order kinded types through the notion of polymorphic predicates. Though both extensions are orthogonal to each other, the second ensures that one gets the most out of the first one (surely, one wants to derive instances for higher-order kinded types).

We are currently implementing the proposal and we hope to make the new features available in the next release of the Glasgow Haskell Compiler.

There are several directions we plan to explore in the future:

- Currently, generic default declarations may be given only for type classes. However, the theory (Hinze, 2000b) also deals with constructor classes whose type parameter range over types of first-order kind. Consequently, we plan to lift this restriction.
- In Haskell 98 instance heads must have the general form $C$ $(T$ $\bar{a})$ where $\bar{a}$ is a sequence of distinct variables. The Glasgow Haskell Compiler, however, allows for non-general instance heads such as $C$ (*List Char*). We are confident that the implementation scheme for generic methods can be extended to deal with this extra complication.
- Multi-parameter type classes are on the wish list of many Haskell programmers. So it would be a shame if the generic extension failed to support them. Now, multi-parameter classes correspond to generic definitions with multiple type arguments, which are theoretically well understood. So we are confident that we can also deal with this generalization.

## References

Backhouse, R., Jansson, P., Jeuring, J. and Meertens, L. (1999) Generic Programming — An Introduction —. Swierstra, S. D., Henriques, P. R. and Oliveira, J. N. (eds), *3rd International Summer School on Advanced Functional Programming, Braga, Portugal.* Lecture Notes in Computer Science 1608, pp. 28–115. Springer-Verlag.

Bird, R., de Moor, O. and Hoogendijk, P. (1996) Generic functional programming with types and relations. *Journal of Functional Programming* **6**(1):1–28.

Cockett, R. and Fukushima, T. (1992) *About Charity*. Yellow Series Report 92/480/18. Dept. of Computer Science, Univ. of Calgary.

Hinze, R. (1999) A generic programming extension for Haskell. Meijer, E. (ed), *Proceedings of the 3rd Haskell Workshop, Paris, France.* The proceedings appear as a technical report of Universiteit Utrecht, UU-CS-1999-28.

Hinze, R. (2000a) Manufacturing datatypes. *Journal of Functional Programming, Special Issue on Algorithmic Aspects of Functional Programming Languages.* To appear.

Hinze, R. (2000b) A new approach to generic functional programming. Reps, T. W. (ed), *Proceedings of the 27th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Boston, Massachusetts, January 19-21* pp. 119–132.

Hinze, R. (2000c) Polytypic values possess polykinded types. Wadler, P. (ed), *Proceedings of the Fifth International Conference on Mathematics of Program Construction (MPC 2000), July 3-5, 2000.*

Jansson, P. and Jeuring, J. (1997) PolyP—a polytypic programming language extension. *Conference Record 24th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL'97, Paris, France* pp. 470–482. ACM-Press.

Jay, C., Bellè, G. and Moggi, E. (1998) Functorial ML. *Journal of Functional Programming* **8**(6):573–619.

Jones, M. P. (1994) *Qualified Types: Theory and Practice.* Cambridge University Press.

Jones, M. P. (2000) Type classes with functional dependencies. Smolka, G. (ed), *Proceedings of the 9th European Symposium on Programming, ESOP 2000, Berlin, Germany.* Lecture Notes in Computer Science 1782, pp. 230–244. Springer-Verlag.

Leivant, D. (1983) Polymorphic type inference. *Proc. 10th Symposium on Principles of Programming Languages.*

Lewis, J. R., Shields, M. B., Meijer, E. and Launchbury, J. (2000) Implicit parameters: Dynamic scoping with static types. Reps, T. W. (ed), *Proceedings of the 27th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Boston, Massachusetts, January 19-21* pp. 108–118.

Okasaki, C. (1998) *Purely Functional Data Structures.* Cambridge University Press.

Peyton Jones, S. and Hughes, J. (eds). (1999) *Haskell 98 — A Non-strict, Purely Functional Language.* Available from `http://www.haskell.org/definition/`.

Peyton Jones, S., Jones, M. and Meijer, E. (1997) Type classes: Exploring the design space. *Proceedings of the Haskell Workshop.*

Ruehr, K. F. (1992) *Analytical and Structural Polymorphism Expressed using Patterns over Types.* PhD thesis, University of Michigan.