### Embedding Generic Programming into Haskell

#### RALF HINZE

### Institut für Informatik III, Universität Bonn Römerstraße 164, 53117 Bonn Email: ralf@informatik.uni-bonn.de Homepage: http://www.informatik.uni-bonn.de/~ralf

July, 2006

Joint work with Andres Löh (Pick up the slides at .../~ralf/talks.html#T51.)



### **1** Introduction



Many of us (most of us?) will probably agree that type systems, especially, polymorphic type systems are a good thing.

A type system is like a suit of armour:

- it shields against the modern dangers of illegal instructions and memory violations, but
- ▶ it also restricts flexibility.

In Haskell 98, for instance, it is not possible to define an equality test that works for all types.

Equality, comparison functions, pretty printers (Haskell's *show*), parsers (Haskell's *read*) have to become known as data-generic or polytypic functions.

Broadly speaking, generic programming is about defining functions that work for all types but that also exhibit type-specific behaviour.



- ▶ In this talk, we show how to embed generic programming into Haskell.
- The embedding builds upon recent advances in type theory: generalised algebraic data types and open data types.
- Put differently, we propose and employ language features that are useful for generic programming.
- We will identify the basic building blocks of generic programming and we will provide an overview of the overall design space.
- We hope to convince you that generic programming is useful and that you can use generic programming techniques today!



- ▶ 1. Introduction
- ▶ 2. Preliminaries
- ► 3. Overloaded functions
- ▶ 4. Generic functions
- ▶ 5. Conclusion



# 2. Preliminaries



- > 2.1 Values, types and kinds
- 2.2 Generalised algebraic data types
- ▶ 2.3 Open data types and open functions



### 2.1 Values, types and kinds



Haskell has the three level structure depicted below.

kinds: \*, \*  $\rightarrow$  \* types: Bool, List  $\alpha$ ,  $\forall \alpha.\alpha \rightarrow \alpha$ values: False, Nil,  $\lambda f \ x \rightarrow f \ (f \ x)$ 

The lowest level — where computations take place — consists of values.

- The second level, which imposes structure on the value level, is inhabited by types.
- On the third level, which imposes structure on the type level, we have so-called kinds. A kind is simply the 'type' of a type constructor.



#### In Haskell, new data types are declared using the data construct.

data Bool = False | True

data  $[\alpha] = Nil | Cons \alpha [\alpha]$ 

data Pair  $\alpha \beta = (\alpha, \beta)$ 

 $\bigcirc$  Data constructors are written in *blue*; type constructors in *red*.

### A data type declaration of the schematic form

data 
$$T \alpha_1 \ldots \alpha_s = C_1 \tau_{1,1} \ldots \tau_{1,m_1} | \cdots | C_n \tau_{n,1} \ldots \tau_{n,m_n}$$

introduces data constructors  $C_1, \ldots, C_n$  with signatures

$$C_i :: \forall \alpha_1 \ldots \alpha_s. \tau_{i,1} \to \cdots \to \tau_{i,m_i} \to T \alpha_1 \ldots \alpha_s$$

rightarrow The data construct is a beast; it combines no less than four features: type abstraction, *n*-ary disjoint sums, *n*-ary cartesian products and type recursion.



### 2.2 Generalised algebraic data types



# Using a recent version of GHC, there is an alternative way of defining data types:

data [] :: 
$$* \to *$$
 where  
 $Nil$  ::  $\forall \alpha.[\alpha]$   
 $Cons$  ::  $\forall \alpha.\alpha \to [\alpha] \to [\alpha]$ 

The first line declares the kind of the new data type.

The type is then inhabited by listing the signatures of the data constructors.



Generalised algebraic data types (GADTs) lift the Haskell 98 restriction that the result type of the constructors must be of the form  $T \alpha_1 \ldots \alpha_s$ .

data  $Expr :: * \to *$  where  $Num :: Int \to Expr Int$   $Plus :: Expr Int \to Expr Int \to Expr Int$   $Eq :: Expr Int \to Expr Int \to Expr Bool$  $If :: \forall \alpha. Expr Bool \to Expr \alpha \to Expr \alpha$ 

rightarrow The data type *Expr* represents typed expressions.



**ξ**2.

An evaluator for the *Expr* data type:

 $eval :: Expr \ \alpha \to \alpha$   $eval \ (Num \ i) = i$   $eval \ (Plus \ e_1 \ e_2) = eval \ e_1 + eval \ e_2$   $eval \ (Eq \ e_1 \ e_2) = eval \ e_1 == eval \ e_2$  $eval \ (If \ e_1 \ e_2 \ e_3) = \mathbf{if} \ eval \ e_1 \mathbf{then} \ eval \ e_2 \mathbf{else} \ eval \ e_3$ 

For functions on GADTs, type signatures are mandatory.

Each equation has a more specific type: the first equation has type  $Expr Int \rightarrow Int$  as Num constrains  $\alpha$  to Int.

The interpreter is quite noticeable in that it is tag free.



### 2.3 Open data types and open functions



For embedding generic function into Haskell, we make use of open data types and open functions, data types and functions that can be freely extended.

An open data type is declared as follows:

open data  $Expr :: * \rightarrow *$ 

Constructors can then be introduced just by providing their type signatures, at any point in the program.



§2.:

§2.3

An open function is declared as follows:

**open**  $eval :: Expr \alpha \to \alpha$ 

The definition of an open function needs not be contiguous; the defining equations may be scattered around the program.

 $\begin{array}{ll} eval \; (Str \; s) &= s \\ eval \; (Show \; e) &= show_{Int} \; (eval \; e) \\ eval \; (Cat \; e_1 \; e_2) &= eval \; e_1 \; \# \; eval \; e_2 \end{array}$ 

The semantics of open data types and open functions is the same as if data types and functions had been defined closed, in a single place.



Open data types and open functions provide two dimensions of extensibility:

- ▶ we can add additional sorts of data, by providing new constructors,
- ▶ we can add additional operations, by defining new functions:

**open** string :: Expr  $\alpha \rightarrow String$ string (Num i) = "(Num"  $\diamond$  show<sub>Int</sub> i + ")" string (Plus  $e_1 e_2$ ) = "(Plus"  $\diamond$  string  $e_1 \diamond$  string  $e_2 + "$ )" string (Eq  $e_1 e_2$ ) = "(Eq"  $\diamond$  string  $e_1 \diamond$  string  $e_2 + "$ )" string (If  $e_1 e_2 e_3$ ) = "(If"  $\diamond$  string  $e_1 \diamond$  string  $e_2 \diamond$  string  $e_3 + "$ )" string (Str s) = "(Str"  $\diamond$  show<sub>String</sub> s + ")" string (Show e) = "(Show"  $\diamond$  string e + ")" string (Cat  $e_1 e_2$ ) = "(Cat"  $\diamond$  string  $e_1 \diamond$  string  $e_2 + "$ )"

#### ◀ ◀ ▶ ▶ □

§2.3

For open functions, first-fit pattern matching is not suitable.

string \_ = ""

Using first-fit pattern matching, this equation effectively closes the definition of *string*.

Instead we use **best-fit left-to-right** pattern matching: the most specific match rather than the first match wins.



### 3. Overloaded functions



In Haskell, showing values of a data type is particularly easy:

data Tree  $\alpha = Empty \mid Node \ (Tree \ \alpha) \ \alpha \ (Tree \ \alpha)$ deriving (Show)

The compiler automatically generates a suitable *show* function.

This function is used, for instance, in interactive sessions:

Now tree [0..3] Node (Node (Node Empty 0 Empty) 1 Empty) 2 (Node Empty 3 Empty)

Here  $tree :: [\alpha] \rightarrow Tree \alpha$  transforms a list into a balanced tree.

The function *show* can be seen as a pretty printer.

The display of larger structures, however, is not especially pretty, due to lack of indentation.

Now tree [0..9] Node (Node (Node Empty 0 Empty) 1 Empty) 2 (Node (Node Em pty 3 Empty) 4 Empty)) 5 (Node (Node (Node Empty 6 Empty) 7 Empt y) 8 (Node Empty 9 Empty))

In the sequel we shall develop a replacement for show, a generic prettier printer.



We use a basic pretty printing library, which just offers support for indentation.

▶ *Text* is type of documents with indentation.

- text converts a string to a text.
- The string passed to *text* must not contain newline characters. The constant *nl* has to be used for that purpose.
- ▶ *indent i* adds *i* spaces after each newline.
- $\blacktriangleright$  ' $\diamondsuit$ ' concatenates two pieces of text.



#### It is a simple exercise to write a prettier printer for trees of integers.



While the program does the job, it is not very general: we can print trees of integers, but not, say, trees of characters.



Of course, it is easy to add another two ad-hoc definitions.

 $\begin{array}{l} pretty_{Char} :: Char \to Text\\ pretty_{Char} c = text \ (show_{Char} \ c) \end{array}$   $\begin{array}{l} pretty_{TreeChar} :: Tree \ Char \to Text\\ pretty_{TreeChar} \ Empty &= text \ "Empty"\\ pretty_{TreeChar} \ (Node \ l \ x \ r) = align \ "(Node \ " \ (pretty_{TreeChar} \ l \ \diamond \ nl \ \diamond \\ pretty_{Char} \ x \ \diamond \ nl \ \diamond \\ pretty_{TreeChar} \ r \ \diamond \ text \ ")") \end{array}$ 

The code of  $pretty_{Tree Char}$  is almost identical to that of  $pretty_{Tree Int}$ .

We actually need a family of pretty printers. A typical case for type classes?

In the sequel we explore a different route!



#### Towards a generic prettier printer — continued

We could define a single function that receives the type as an additional argument and suitably dispatches on this type argument.

ξ3

 $pretty :: (\alpha :: *) \to \alpha \to Text$ 

Haskell doesn't permit the explicit passing of types.

We could pass the pretty printer an additional argument that represents the type.

 $pretty :: Type \rightarrow \alpha \rightarrow Text$ 

This is too simple-minded: a function of this type must necessarily ignore its second parameter (parametricity, "free theorem").

pretty :: Type  $\alpha \to \alpha \to Text$ 

An element of type  $Type \tau$  is a representation of the type  $\tau$ .

### Using a generalised algebraic data type, we can define Type directly in Haskell.

**open data**  $Type :: * \to *$  **where**  *Char* :: Type *Char Int* :: Type *Int Pair* ::  $Type \ \alpha \to Type \ \beta \to Type \ (\alpha, \beta)$  *List* ::  $Type \ \alpha \to Type \ [\alpha]$ *Tree* ::  $Type \ \alpha \to Type \ (Tree \ \alpha)$ 



Each type has a unique representation:

- $\blacktriangleright$  the type *Int* is represented by the constructor *Int*,
- ▶ the type (*Char*, *Int*) is represented by *Pair Char Int*,
- ▶ the type [*Tree Char*] is represented by List (*Tree Char*).

Recall: type constructors are written in red; data constructors in blue.

For any given  $\tau$  the type  $Type \tau$  comprises exactly one element:  $Type \tau$  is a so-called singleton type.



We shall often need to annotate an expression with its type representation.

data Typed  $\alpha = (:) \{ val :: \alpha, type :: Type \alpha \}$ 

The definition, which makes use of Haskell's record syntax, introduces the colon ':' as an infix data constructor.

- ▶ 4711 : *Int* is an element of *Typed Int*.
- ► (47, 'R'): *Pair Int Char* is an element of *Typed* (*Int*, *Char*).

Solution Note the difference between x : t and  $x :: \tau$ .

- $\blacktriangleright x: t$  is a pair consisting of a value x and a representation t of its type.
- $\blacktriangleright x :: \tau$  is Haskell syntax for 'x has type  $\tau$ '.



#### An almost generic prettier printer

Given these prerequisites, we can finally define the desired pretty printer.

We declare pretty to be open so that we can later extend it.

 $rightarrow Typed \ \alpha \to Text$  is an uncurried version of  $Type \ \alpha \to \alpha \to Text$ .



The pretty printer produces output in the following style.

```
\begin{array}{c} Now \rangle \ pretty \ (tree \ [0 \, . \, 3] : Tree \ Int) \\ (Node \ (Node \ (Node \ Empty) \\ 0 \\ Empty) \\ 2 \\ (Node \ Empty \\ 3 \\ Empty)) \end{array}
```



The type of type representations is, of course, by no means special to pretty printing.

A second example: collecting integers.



We have declared Type to be open so that we can freely add new constructors to the Type data type and that we can freely add new equations to existing open functions on Type.

Whenenver we define a new data type

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

we extend Type by a new constructor.

Perfect :: Type  $\alpha \rightarrow$  Type (Perfect  $\alpha$ )

Perfect is a so-called nested data type.



§3

Then we extend pretty by suitable equations.

 $\begin{array}{l} pretty \; (Zero \; x : Perfect \; a) \\ = a lign \; "(\texttt{Zero } \; " \; (pretty \; (x : a) \diamondsuit text \; ")") \\ pretty \; (Succ \; x : Perfect \; a) \\ = a lign \; "(\texttt{Succ } \; " \; (pretty \; (x : Perfect \; (Pair \; a \; a)) \diamondsuit text \; ")") \end{array}$ 



#### An example session





Whenever we define a new data type,

- ▶ we add a constructor of the same name to the type of type representations,
- ▶ we add corresponding equations to all generic functions.

Observations:

- ▶ the extension of *Type* is cheap and easy (a compiler could do this for us),
- the extension of all functions on *Type* is laborious and difficult (can you imagine a compiler doing that?).

In the next section we shall develop a scheme so that it suffices to extend one or two particular overloaded functions. The remaining functions adapt themselves.

#### Terminology: overloaded and generic functions

**overloaded and generic functions**. An overloaded function works for a fixed family of types. By contrast, a generic function works for all types, including types that the programmer is yet to define.



# **3.3 Generic functions**



We need to find a way to treat elements of a data type in a uniform way. Consider an arbitrary element of some data type:

 $C e_1 \cdots e_n$ 

The idea is to make this applicative structure visible and accessible: we mark the constructor using Con and each function application using ' $\diamond$ '.

- ► *Empty* becomes *Con empty*,
- ▶ Node  $l \ a \ r$  becomes Con node  $\diamond$  (l: Tree Int)  $\diamond$  (a: Int)  $\diamond$  (r: Tree Int).

The arguments are additionally annotated with their types and the constructor itself with information on its syntax.



**§**3.3

The functions Con and ' $\diamond$ ' are constructors of a data type called *Spine*.

data Spine ::  $* \to *$  where Con :: Constr  $\alpha \to Spine \alpha$ ( $\diamond$ ) :: Spine ( $\alpha \to \beta$ )  $\to$  Typed  $\alpha \to Spine \beta$ 

The type is called *Spine* because its elements represent the possibly partial spine of a constructor application.



#### §3.3

#### The following table illustrates the stepwise construction of a spine.

 $\begin{array}{l} \textit{node} :: \textit{Constr} (\textit{Tree Int} \rightarrow \textit{Int} \rightarrow \textit{Tree Int} \rightarrow \textit{Tree Int}) \\ \textit{Con node} :: \textit{Spine} (\textit{Tree Int} \rightarrow \textit{Int} \rightarrow \textit{Tree Int} \rightarrow \textit{Tree Int}) \\ \textit{Con node} \diamond (l:\textit{Tree Int}) :: \textit{Spine} (\textit{Int} \rightarrow \textit{Tree Int} \rightarrow \textit{Tree Int}) \\ \textit{Con node} \diamond (l:\textit{Tree Int}) \diamond (a:\textit{Int}) :: \textit{Spine} (\textit{Tree Int} \rightarrow \textit{Tree Int}) \\ \textit{Con node} \diamond (l:\textit{Tree Int}) \diamond (a:\textit{Int}) :: \textit{Spine} (\textit{Tree Int} \rightarrow \textit{Tree Int}) \\ \textit{Con node} \diamond (l:\textit{Tree Int}) \diamond (a:\textit{Int}) :: \textit{Spine} (\textit{Tree Int}) \\ \end{array}$ 



Elements of type  $Constr \alpha$  comprise an element of type  $\alpha$ , namely the original data constructor, plus additional information about its syntax.

data Constr  $\alpha = Descr\{constr :: \alpha$ , name :: String}



Given a value of type  $Spine \alpha$ , we can easily recover the original value of type  $\alpha$ :

from Spine :: Spine  $\alpha \to \alpha$ from Spine (Con c) = constr c from Spine ( $f \diamond x$ ) = (from Spine f) (val x)

*fromSpine* is parametrically polymorphic.



The inverse of from Spine is an overloaded function of type  $Typed \alpha \rightarrow Spine \alpha$ .

Its definition, however, follows a trivial pattern: if the data type comprises a constructor  ${\cal C}$ 

 $C::\tau_1\to\cdots\to\tau_n\to\tau_0$ 

then the equation for toSpine takes the form

$$toSpine (C x_1 \ldots x_n : t_0) = Con c \diamond (x_1 : t_1) \diamond \cdots \diamond (x_n : t_n)$$

where c is the annotated version of C and  $t_i$  is the type representation of  $\tau_i$ .

§3.3

As an example, here is the definition of toSpine for binary trees.

```
open toSpine :: Typed \alpha \rightarrow Spine \alpha

toSpine (Empty : Tree a) = Con empty

toSpine (Node l x r : Tree a)

= Con node \diamond (l : Tree a) \diamond (x : a) \diamond (r : Tree a)

empty :: Constr (Tree \alpha)

empty = Descr{ constr = Empty,

name = "Empty"}

node :: Constr (Tree \alpha \rightarrow \alpha \rightarrow Tree \alpha \rightarrow Tree \alpha)

node = Descr{ constr = Node,

name = "Node"}
```



**ξ**3.3

With all the machinery in place we can now turn *pretty* and *ints* into truly generic functions.

The idea is to add a catch-all case that takes care of all the remaining type cases in a uniform manner.

 $\begin{array}{l} ints :: \textit{Typed } \alpha \rightarrow [\textit{Int}] \\ ints \; (i:\textit{Int}) = [i] \\ ints \; x \qquad = intsSpine \; (toSpine \; x) \\ intsSpine :: \textit{Spine } \alpha \rightarrow [\textit{Int}] \\ intsSpine \; (Con \; c) = [] \\ intsSpine \; (f \diamond x) \qquad = intsSpine \; f \; + \; ints \; x \end{array}$ 



§3.3

The catch-all case for pretty is almost as easy. We only have to take care that we do not parenthesize nullary constructors.

 $pretty \ x = prettySpine \ (toSpine \ x)$   $prettySpine :: Spine \ \alpha \to Text$   $prettySpine \ (Con \ c) = text \ (name \ c)$   $prettySpine \ (f \diamond x) = prettySpine1 \ f \ (pretty \ x)$   $prettySpine1 :: Spine \ \alpha \to Text \to Text$   $prettySpine1 \ (Con \ c) \ d = align \ ("(" + name \ c + " ")) \ (d \ \diamond text ")")$   $prettySpine1 \ (f \diamond x) \ d = prettySpine1 \ f \ (pretty \ x \ \diamond nl \ \diamond d)$ 



**§**3.3

§3.3

Now, why are we in a better situation than before?

- ▶ When we introduce a new data type we still have to extend *Type* and provide cases for the data constructors in the *toSpine* function.
- ▶ This has to be done only once per data type.
- The code for the generic functions (of which there can be many) is completely unaffected by the addition of a new data type.
- As a further plus, the generic functions are unaffected by changes to a given data type. Only the function *toSpine* must be adapted to the new definition.

# 4. Conclusion



Using reflected types we can program overloaded functions. Using a uniform view on data we can generalise overloaded functions to generic ones.

Support for generic programming consists of three essential ingredients:

- ▶ a type reflection mechanism,
- ► a type representation, and
- ▶ a generic view on data.

For each dimension there are several choices: instead of the data type Type we could use type classes or a type-safe cast.



view(s)	representation of overloaded functions			
	type reflection	type classes	type-safe cast	specialisation
none	ITA	-	_	_
fixed point	Reloaded	PolyP	-	PolyP
sum-of-products	LIGD	DTC, GC, GM	-	GH
spine	Reloaded, Revolutions	SYB, Reloaded	SYB	-

