# Generic and Indexed Programming

*Jeremy Gibbons*

*University of Oxford*

*WG2.1#62, December 2006*

# 1. Background

- generic programming: *parametrization*

- datatype-generic programming: parametrization *by datatype*

- special-purpose languages and constructs: *GH*, *SyB*...

- *lightweight embeddings* in general-purpose languages

## 1.1. Capturing properties

Linguistic approaches to modelling:
find new ways to express properties within programs.

Narrowing the *semantic gap* between
the programmer's head and the program.

- type systems

- assertions and testing frameworks

(There are extra-linguistic modelling approaches too,
but we won't discuss them here.)

## 1.2. Dependently-typed programming

- types classify values

- dependent types classify values more precisely:
  in particular, the way in which values depend on other values

- eg $Vector_n\ \mathbb{Z}$, the type of $n$-vectors of integers

- more generally, *dependent product* type $\Sigma\ n :: \mathbb{N}.\ f\ (n)$ of pairs $(n, x)$
  with $n :: \mathbb{N}$ and $x :: f\ (n)$

- play a central role in constructive logics
  ('propositions as types', 'Curry–Howard isomorphism')

## 1.3. Generalized algebraic datatypes

Slight generalization of algebraic datatypes, allowing result type of constructor to be *strictly more specific* than declared datatype.

Allows use of types as indices, capturing program properties. A kind of lightweight dependently-typed programming, by lifting some index values to the type level.

Also known as *first-class phantom types*, *guarded recursive datatypes*, *indexed types*, *equality-constrained types*…apparently a good idea!

# 2. Generalizing algebraic datatypes

Standard algebraic datatypes, as in Haskell:

> **data** *Expr = N Int*
>       *| Add Expr Expr*
>       *| B Bool*
>       *| IsZ Expr*
>       *| If Expr Expr Expr*

They can be polymorphic too, with type parameters:

> **data** *List a = Nil*
>       *| Cons a (List a)*

## 2.1. Definitions by pattern-matching

**data** *Result = NR Int | BR Bool*

*eval :: Expr → Result*

*eval* (*N n*)        = *NR n*

*eval* (*Add x y*) = **case** (*eval x, eval y*) **of**

$\qquad\qquad$ (*NR m, NR n*) → *NR* (*m + n*)

*eval* (*B b*)        = *BR b*

*eval* (*IsZ x*)      = **case** (*eval x*) **of**

$\qquad\qquad$ *NR n* → *NB* (0 ≡ *n*)

*eval* (*If x y z*)  = **case** (*eval x*) **of**

$\qquad\qquad$ *NB b* → **if** *b* **then** *eval y* **else** *eval z*

Note the explicit tagging and untagging
(and the lack of error-checking for ill-formed expressions!).

## 2.2. Extended syntax

New syntax, explicitly stating constructor types (and datatype kind):

> **data** *Expr* :: * **where**
>   *N*   :: *Int* →                          *Expr*
>   *Add* :: *Expr* → *Expr* →        *Expr*
>   *B*   :: *Bool* →                 *Expr*
>   *IsZ*  :: *Expr* →                *Expr*
>   *If*   :: *Expr* → *Expr* → *Expr* → *Expr*
>
> **data** *List* :: * → * **where**
>   *Nil*   ::                 *List a*
>   *Cons* :: *a* → *List a* → *List a*

Note that for ordinary polymorphic algebraic datatypes, all constructors have the same (most general) result type.

## 2.3. GADT declaration

Make the datatype polymorphic, with a type parameter
(in this case, expressing the represented type).

> **data** *Expr* :: $* \to *$ **where**
>   *N*   :: *Int* $\to$                                      *Expr Int*
>   *Add* :: *Expr Int* $\to$ *Expr Int* $\to$                *Expr Int*
>   *B*   :: *Bool* $\to$                                     *Expr Bool*
>   *IsZ*  :: *Expr Int* $\to$                                *Expr Bool*
>   *If*   :: *Expr Bool* $\to$ *Expr Int* $\to$ *Expr Int* $\to$ *Expr Int*

Now constructors may have more specialized return types.

Note that the type parameter is a *phantom type*:
a value of type *Expr a* need not contain elements of type *a*.

## 2.4. GADT use

Specialized return types of constructors induce type constraints, which are exploited in type-checking definitions.

$$eval :: Expr\ a \rightarrow a$$
$$eval\ (N\ n)\qquad = n$$
$$eval\ (Add\ x\ y) = eval\ x + eval\ y$$
$$eval\ (B\ b)\qquad = b$$
$$eval\ (IsZ\ x)\qquad = 0 \equiv eval\ x$$
$$eval\ (If\ x\ y\ z)\ = \textbf{if}\ eval\ x\ \textbf{then}\ eval\ y\ \textbf{else}\ eval\ z$$

Note that all the tagging and untagging has gone,
and with it the possibility of run-time errors.

By explicitly stating a property formerly implicit in the code,
we have gained both in safety and in efficiency.

# 3. Application: indexing by size

Empty datatypes as indices (so $S\ (S\ Z)$ is a type).

> **data** $Z$
> **data** $S\ n$

Size-indexed type of vectors:

> **data** $Vector :: * \to * \to *$ **where**
>   $VNil$  $::$                              $Vector\ a\ Z$
>   $VCons :: a \to Vector\ a\ n \to Vector\ a\ (S\ n)$

Size constraint on $vzip$ is captured in the type:

> $vzip :: Vector\ a\ n \to Vector\ b\ n \to Vector\ (a, b)\ n$
> $vzip\ VNil\ VNil$                          $=\ VNil$
> $vzip\ (VCons\ a\ x)\ (VCons\ b\ y) = VCons\ (a, b)\ (vzip\ x\ y)$

# 4. Application: indexing by shape

*Red-black trees* are binary search trees in which:

- every node is coloured either red or black

- every red node has a black parent

- every path from the root to a leaf contains
  the same number of black nodes (enforcing approximate balance)

In *RBTree a c n*, type *a* is the element type, *c* the root colour, and *n* the black height.

```
data R
data B
data RBTree :: * → * → * → * where
   Empty ::                                              RBTree a B Z
   Red    :: RBTree a B n → a → RBTree a B n  → RBTree a R n
   Black  :: RBTree a c n → a → RBTree a c' n → RBTree a B (S n)
```

# 5. Application: indexing by unit

Suppose dimensions of non-negative powers of metres and seconds:

> **data** *Dim* :: $* \to * \to *$ **where**
>   *D* :: *Float* $\to$ *Dim m s*
>
> *distance* :: *Dim* (*S Z*) *Z*
> *distance* = *D* 3.0
>
> *time* :: *Dim Z* (*S Z*)
> *time* = *D* 2.0

A dimensioned value is a *Float* with two type-level tags.

> *dadd* :: *Dim m s* $\to$ *Dim m s* $\to$ *Dim m s*
> *dadd* (*D x*) (*D y*) = *D* (*x* + *y*)

Now *dadd time time* is well-typed, but *dadd distance time* is ill-typed.

(More interesting to allow negative powers too, but for brevity...)

## 5.1. Type-level functions

Proofs of properties about indices:

> **data** *Add* :: $* \rightarrow * \rightarrow * \rightarrow *$ **where**
>   *AddZ* ::                          *Add Z n n*
>   *AddS* :: *Add m n p* $\rightarrow$ *Add (S m) n (S p)*

Used to constrain the type of dimensioned multiplication:

> *dmult* :: (*Add* $m_1$ $m_2$ *m*, *Add* $s_1$ $s_2$ *s*) $\rightarrow$
>              *Dim* $m_1$ $s_1$ $\rightarrow$ *Dim* $m_2$ $s_2$ $\rightarrow$ *Dim m s*
> *dmult* (_, _) (*D x*) (*D y*) = *D* (*x* $\times$ *y*)

Thus, type-index of product is computed from indices of arguments.

## 5.2. Inferring proofs of properties

Capture the proof as a type class (multi-parameter, with functional dependency; essentially a function on types).

$$\textbf{class } Add\ m\ n\ p\ |\ m\ n \to p$$

$$\textbf{instance} \qquad\qquad Add\ Z\ n\ n$$
$$\textbf{instance } Add\ m\ n\ p \Rightarrow Add\ (S\ m)\ n\ (S\ p)$$

Now the proof can be (type-)inferred rather than passed explicitly.

$$dmult :: (Add\ m_1\ m_2\ m, Add\ s_1\ s_2\ s) \Rightarrow$$
$$Dim\ m_1\ s_1 \to Dim\ m_2\ s_2 \to Dim\ m\ s$$
$$dmult\ (D\ x)\ (D\ y) = D\ (x \times y)$$

Note that the type class has no methods, so corresponds to an empty dictionary; it can be optimized away.

# 6. Application: indexing by state

The 'ketchup problem':

**data** *O*

**data** *C*

**data** *Edge* :: $* \rightarrow * \rightarrow *$ **where**

  *Open*  :: *Edge O C*

  *Close*  :: *Edge C O*
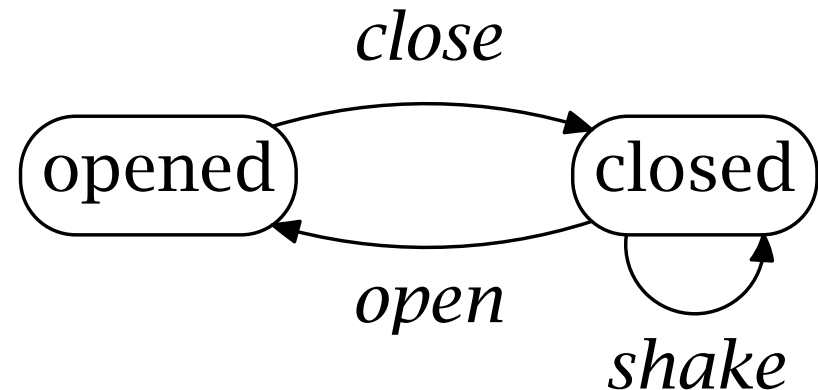
  *Shake* :: *Edge C C*

**data** *Path* :: $* \rightarrow * \rightarrow *$ **where**

  *Empty* :: *Path s s*

  *PCons* :: *Edge x y* $\rightarrow$ *Path y z* $\rightarrow$ *Path x z*

*scenario* :: *Path O O*

*scenario = PCons Open* (*PCons Shake* (*PCons Close Empty*))

# 7. Application: indexing by type

*Generic programming* is about writing programs parametrized by datatypes; for example, a generic data marshaller.

One implementation of generic programming manifests the parameters as some family of *type representations*.

For example, C's *sprintf* is generic over a family of *format specifiers*.

> **data** *Format* :: $* \to *$ **where**
>     *I* ::            *Format a → Format* (*Int → a*)
>     *B* ::            *Format a → Format* (*Bool → a*)
>     *S* :: *String → Format a → Format a*
>     *F* ::                        *Format String*

A term of type *Format a* is a representation of the type *a*, for various types *a* appropriate for *sprintf*, such as *Int → Bool → String*.

## 7.1. Type-indexed dispatching

The function *sprintf interprets* that representation.

$$sprintf :: Format\ a \to a$$
$$sprintf\ fmt = aux\ id\ fmt\ \textbf{where}$$
$$aux :: (String \to String) \to Format\ a \to a$$
$$aux\ f\ (I\ fmt)\ \ = \lambda n \to aux\ (f \circ (show\ n +\!\!+))\ fmt$$
$$aux\ f\ (B\ fmt)\ \ = \lambda b \to aux\ (f \circ (show\ b +\!\!+))\ fmt$$
$$aux\ f\ (S\ s\ fmt) = aux\ (f \circ (s +\!\!+))\ fmt$$
$$aux\ f\ (F)\ \ \ \ \ \ \ \ = f\ \texttt{""}$$

For example, *sprintf f* 13 *True* = `"Int is 13, bool is True."`, where

$$f :: Format\ (Int \to Bool \to String)$$
$$f = S\ \texttt{"Int is "}\ (I\ (S\ \texttt{", bool is "}\ (B\ (S\ \texttt{"."}\ F))))$$

# 8. Adding weight

We have shown some examples in Haskell with small extensions.

This is a very lightweight approach to dependently-typed programming.

Lightweight approaches have low entry cost, but relatively high continued cost: encoding via type classes etc is a bit painful.

Tim Sheard's *Ωmega* is a cut-down version of Haskell with explicit support for GADTs:

- kind declarations

- type-level functions

- statically-generated witnesses

Xi and Pfenning's *Dependent ML* provides natural-number indices, and incorporates decision procedures for discharging proof obligations.

These are more heavyweight approaches (such as McBride's *Epigram*).

## 8.1. Transfer to the mainstream

Kennedy and Russo (OOPSLA 2005) showed that Java and C#

> 'can express GADT definitions, and a large class of GADT-manipulating programs, through the use of generics, subclassing and virtual dispatch'

(with a few casts, that they propose ways around).

# 9. The GIP project at Oxford

- EPSRC-funded, about £500k

- three and a half years, from November 2006

- Jeremy Gibbons, principal investigator

- Bruno Oliveira, postdoctoral researcher

- Meng Wang, doctoral student

## 9.1. Workpackages

- *capturing properties*

    case studies as benchmarks

- *generics for GADTs*

    GADTs no longer sums of products: spine views, idioms

- *extensible generic functions*

    expression problem, combining structural and nominal views

- *design patterns as a library*

    GADTs in Scala, Java and C#

- *type classes and GADTs*

    inferring proof objects

- *impedence transformers*

    statically-checked metaprogramming; multi-tier