

Algebraic Specialization of Generic Functions for Recursive Types

Alcino Cunha¹ Hugo Pacheco²

*DI-CCTC
Universidade do Minho
Braga, Portugal*

Abstract

Defining functions over large, possibly recursive, data structures usually involves a lot of boilerplate. This code simply traverses non-interesting parts of the data, and rapidly becomes a maintainability problem. Many generic programming libraries have been proposed to address this issue. Most of them allow the user to specify the behavior just for the interesting bits of the structure, and provide traversal combinators to “scrap the boilerplate”. The expressive power of these libraries usually comes at the cost of efficiency, since runtime checks are used to detect where to apply the type-specific behavior.

In previous work we have developed an effective rewrite system for specialization and optimization of generic programs. In this paper we extend it to also cover recursive data types. The key idea is to specialize traversal combinators using well-known recursion patterns, such as folds or paramorphisms. These are ruled by a rich set of algebraic laws that enable aggressive optimizations. We present a type-safe encoding of this rewrite system in Haskell, based on recent language extensions such as type-indexed type families.

Keywords: Generic programming, Recursion patterns, Rewrite systems, Specialization, Type families

1 Introduction

Modeling real-world problems in a functional language typically leads to a large set of recursive data types, each with a lot of different constructors. That is the case, for example, when developing language processing tools, where grammars are represented using a different data type for each non-terminal and a different constructor for each production rule. Similarly, schema-aware XML processing usually involves mapping a huge schema to an equivalent data type, with each of the many elements mapped to a different type. Such proliferation of data types makes it hard to implement even conceptually simple functions, that manipulate a very small subset of the data constructors.

A classic (but still benign) example used to illustrate this problem is the so called “paradise benchmark” [LJ03]. Suppose one has a XML schema to model a company with several departments, each having a name, a manager and a collection

¹ Email: alcino@di.uminho.pt

² Email: hpacheco@di.uminho.pt

of employees or sub-departments. This schema could be represented by the following Haskell data type.

```

data Company = C [Dept]
data Dept = D Name Manager [Either Employee Dept]
data Employee = E Person Salary
data Person = P Name Address
data Salary = S Int
type Manager = Employee
type Name = String
type Address = String

```

Suppose one also wants to define a function to increase all salaries by a fixed amount k . A possible definition of this function could be

```

increase :: Int → Company → Company
increase k (C ds) = C (map (incD k) ds)
  where incD k (D nm mgr us) = D nm (incE k mgr) (map (incU k) us)
        incU k (Left e)   = Left (incE k e)
        incU k (Right d)  = Right (incD k d)
        incE k (E p s)    = E p (incS k s)
        incS k (S s)      = S (s + k)

```

Even this rather simple definition is filled with boilerplate code, whose only purpose is to perform a standard traversal of the *Company* data type to find salaries to increase. Apart from aesthetical reasons, this kind of boilerplate has some major drawbacks: (1) it makes code understanding rather difficult, since the only interesting functions (in this case *incS*) are lost amid bucketloads of uninteresting code; and (2) it rapidly becomes a maintainability problem, since each model evolution implies changes to many functions that are only concerned with parts of the model not affected by the evolution.

Many generic programming libraries have been proposed to address this issue. Most of them allow the user to specify the behavior just for the interesting bits of the structure, and provide traversal combinators to encode the remaining boilerplate. One of the most successful libraries is the conveniently named “Scrap you Boilerplate” (SYB), first introduced in [LJ03] and subsequently extended with additional functionalities [LJ05]. Using this library, the *increase* function could be redefined as follows.

```

increase :: Int → Company → Company
increase k = everywhere (mkT (incS k))
  where incS k (S s) = S (s + k)

```

The *everywhere* combinator traverses a data structure in bottom-up fashion, applying the given generic transformation to all its nodes. The *mkT* combinator builds a generic transformation from a type specific one: given a function $f :: a \rightarrow a$, *mkT f* behaves like f for all values of type a and like the identity function otherwise. Be-

sides being much easier to understand what *increase* does, its definition will stay the same even if the *Company* data type changes.

Unfortunately, the obvious advantages provided by this style of generic programming come at a price: the performance of generic functions is much worse than analogous non-generic ones. In [MR07], the SYB implementation of a standard set of benchmark functions was reported to run 7 times slower in average than the non-generic implementation. Part of this performance loss is due to the run-time checks needed to determine at each node whether to apply specific or generic behavior. The remaining is due to structural reasons inherent to this style of generic programming: the traversal combinators must blindly traverse the whole data structure, even if a certain branch does not mention types where the specific behavior applies.

Some new SYB-like generic programming libraries have been proposed to address this efficiency problem. According to a recent survey [RJJ⁺08], two of the most efficient are Uniplate [MR07] and Smash [Kis06]. The former outperforms SYB by restricting the power of the traversal combinators. The latter offsets some of the run-time checks to compile-time, but needs extra work from the programmer in order to support new data types.

In previous work [CV07b], we have taken a different approach to tackle this problem: we developed a rewrite system that specializes generic functions for specific types. This specialization proceeds in two phases: (1) the generic functions are specialized to non-optimized point-free definitions; (2) these definitions are then optimized using standard algebraic laws for point-free combinators. The major drawback of that approach was the lack of support for user defined recursive types, such as the *Company* type above.

The major contribution of this paper is to extend that specialization mechanism to also cover recursive data types. More specifically, we will focus on inductive data types that can be defined as fixpoints of functors. The key idea is to specialize traversal combinators using well-known recursion patterns for inductive types, such as folds or paramorphisms. Likewise to standard point-free combinators, these recursion patterns are also characterized by a rich set of algebraic laws that enable further optimizations after specialization. Most of the definitions that result from the new rewrite system have runtimes close to the hand-written non-generic ones. Another contribution is the Haskell encoding of these new laws: thanks to recent language extensions such as type-indexed type families, we managed to get an implementation that closely mimics the theory.

Section 2 briefly surveys the SYB approach to generic programming and recaps our previous work on the specialization of generic functions for non-recursive types. Section 3 extends this work with new algebraic rules for the specialization for recursive types. Section 4 discusses how these new rules can be accommodated in a type-safe rewrite system implemented in Haskell. Section 5 presents specialization examples and compares the respective speedups. Section 6 presents related work, and Section 7 makes some concluding remarks, including some future work ideas.

2 Specialization for Non-recursive Types

In this paper we will focus on a limited set of combinators that capture the essence of strategic generic programming libraries like SYB. Generic functions come in two flavors: type-preserving (transformations) and type-unifying (queries). For defining type-preserving functions we have the following combinators:

$$\begin{aligned}
 \text{nop} &:: T \\
 (\triangleright) &:: T \rightarrow T \rightarrow T \\
 \text{gmap}T &:: T \rightarrow T \\
 \text{everywhere} &:: T \rightarrow T \\
 \text{mk}T_A &:: (A \rightarrow A) \rightarrow T \\
 \text{ap}T_A &:: T \rightarrow (A \rightarrow A)
 \end{aligned}$$

In the SYB library, the type T of type-preserving generic functions is defined as $\forall a. \text{Data } a \Rightarrow a \rightarrow a$. Type classes like Data are extensively used in SYB to infer type representations for data types. Among others, these are necessary in the definition of $\text{mk}T$ to determine where the type-specific transformation should be applied. To simplify the presentation, instead of using type classes, we will parameterize $\text{mk}T$ with an explicit type representation. Besides everywhere , we have combinators to map a transformation over all direct children of node ($\text{gmap}T$), to sequence transformations (\triangleright), and to denote the identity transformation (nop). We also have an explicit combinator to apply a generic transformation to a particular type ($\text{ap}T$). In SYB this is done implicitly via type-classes. The transformation to increase all salaries would now be written as follows.

$$\begin{aligned}
 \text{increase} &:: \text{Int} \rightarrow \text{Company} \rightarrow \text{Company} \\
 \text{increase } k &= \text{ap}T_{\text{Company}} \text{ everywhere } (\text{mk}T_{\text{Salary}} (\text{inc}S \ k)) \\
 \text{where } \text{inc}S \ k &(S \ s) = S \ (s + k)
 \end{aligned}$$

For defining type-unifying functions we have the following combinators:

$$\begin{aligned}
 \emptyset &:: Q \ R \\
 (\cup) &:: Q \ R \rightarrow Q \ R \rightarrow Q \ R \\
 \text{gmap}Q &:: Q \ R \rightarrow Q \ R \\
 \text{everything} &:: Q \ R \rightarrow Q \ R \\
 \text{mk}Q_A &:: (A \rightarrow R) \rightarrow Q \ R \\
 \text{ap}Q_A &:: Q \ R \rightarrow (A \rightarrow R)
 \end{aligned}$$

In this case, $Q \ R$ represents the type of generic queries with result type R . Once more, to simplify the presentation of the specialization laws we will assume that R is a monoid, with a *zero* element and an associative *plus* operator. In practice, this makes little difference since most typical results, namely lists and integers, are indeed monoids. $\text{mk}Q_A$ creates a generic query out of a type-specific one, returning *zero* for types other than A . everything collects all results in a bottom-up traversal using the *plus* operator. $\text{gmap}Q$ collects the results of applying a query to all direct children, \cup sums the results of two queries, and \emptyset denotes the query that always returns *zero*. To apply a generic query to a particular type we have $\text{ap}Q$.

id	$:: A \rightarrow A$
(\circ)	$:: (B \rightarrow C) \rightarrow (A \rightarrow B) \rightarrow (A \rightarrow C)$
π_1	$:: A \times B \rightarrow A$
π_2	$:: A \times B \rightarrow B$
(Δ)	$:: (A \rightarrow B) \rightarrow (A \rightarrow C) \rightarrow (A \rightarrow B \times C)$
\times	$:: (A \rightarrow B) \rightarrow (C \rightarrow D) \rightarrow (A \times C \rightarrow B \times D)$
i_1	$:: A \rightarrow A + B$
i_2	$:: B \rightarrow A + B$
(∇)	$:: (A \rightarrow C) \rightarrow (B \rightarrow C) \rightarrow (A + B \rightarrow C)$
$(+)$	$:: (A \rightarrow B) \rightarrow (C \rightarrow D) \rightarrow (A + C \rightarrow B + D)$

Fig. 1. Point-free combinators.

For example, to compute the total salary bill of a company we could define the following generic query:

```

salaries :: Company → Int
salaries = apQCompany everything (mkQSalary bills)
where bills (S s) = s
  
```

In previous work [CV07b] we have presented a rewrite system to specialize generic functions to type-specific point-free definitions. In the point-free style of programming, functions are composed using a standard set of higher-order combinators, avoiding the need to explicitly mention the domain points as variables. This variable-free style (popularized by John Backus in his 1977 Turing award lecture [Bac78]) is particularly amenable for program calculation since its combinators are characterized by a rich set of algebraic laws. We use a rather standard set of point-free combinators for handling products and sums (see Figure 1). Their behavior should be clear from the type signatures. For more information about the laws ruling this combinators and point-free program calculation in general see [Gib02, CPP06].

Most user defined data types can be defined as the fixpoint of a regular functor. The *base functor* that captures the signature of a data type A will be denoted F_A (when the type A is clear from the context we will often omit it from the subscript). A regular functor is either the identity functor Id , the constant functor \underline{A} (that always returns A), the lifting of the sum \oplus and product bifunctors \otimes , or the composition \odot of functors. For example, for lists we have $F_{[A]} = \underline{1} \oplus \underline{A} \otimes Id$. If the type is not recursive, its base functor will not have any identity. For example, $F_{(Maybe\ A)} = \underline{1} \oplus \underline{A}$. Associated with each data type A we also have two unique functions $in_A : F_A\ A \rightarrow A$ and $out_A :: A \rightarrow F_A\ A$, that are each other's inverse. They allow us to encode and inspect values of the given type, respectively.

Figure 2 presents the laws used to specialize type-preserving combinators into point-free. Specialization proceeds by pushing down the *apT* combinator until it gets consumed by the *mkT-APPLY* law. Similar laws exist for the type-unifying combinators. Although not generic, the definitions produced by the specialization phase are very inefficient because they still traverse the whole data structure. However, using point-free program calculation laws they can be optimized in order to eliminate redundant traversals. Notice that the *everywhere-APPLY* law uses the

$apT_A \text{ nop} = id$	<i>nop</i> -APPLY
$apT_A (f \triangleright g) = apT_A f \circ apT_A g$	\triangleright -APPLY
$apT_A (gmapT f) = id, \text{ if } A \text{ base type}$	$gmapT$ -APPLY
$apT_A (gmapT f) = in_A \circ apT_{F A} f \circ out_A, \text{ if } A \text{ datatype}$	
$apT_A (everywhere f) = apT_A (gmapT (everywhere f) \triangleright f)$	<i>everywhere</i> -APPLY
$apT_A (mkT_A f) = f$	mkT -APPLY
$apT_A (mkT_B f) = id, \text{ if } A \neq B$	
$apT_{A \times B} f = apT_A f \times apT_B f$	\times -APPLY
$apT_{A+B} f = apT_A f + apT_B f$	$+$ -APPLY
$apT_1 f = id$	1-APPLY

Fig. 2. Laws for specializing generic transformations.

recursive definition of this traversal combinator using $gmapT$ and \triangleright . Since the previous rewrite system only handled non-recursive user defined data types, this law did not pose any termination problems. However, for recursive types it cannot be used since it would lead to an infinite expansion of the definition (due to successive expansions of *everywhere* in recursive occurrences of the type).

3 Specialization for Recursive Types

The key to avoid infinite expansions is to specialize traversal combinators using an alternative definition based on standard recursion patterns such as folds. Likewise to point-free combinators, these recursion patterns are characterized by powerful algebraic laws, that will enable us to optimize the specialized definitions. For a comprehensive presentation of most standard recursion patterns and the respective laws see [MFP91].

The standard recursion pattern of iteration, usually known as fold or *catamorphism*, consumes an inductive type by replacing its constructors with a given argument function. For an inductive type A , given a function $g :: F B \rightarrow B$, $\langle g \rangle_A :: A \rightarrow B$ denotes a fold over that type that produces values of type B . Its recursive definition can be clearly depicted in the following diagram.

$$\begin{array}{ccc}
 A & \xrightarrow{out_A} & F A \\
 \langle g \rangle_A \downarrow & & \downarrow F \langle g \rangle_A \\
 B & \xleftarrow{g} & F B
 \end{array}$$

While folds can express functions defined by iteration, *paramorphisms* can express all functions that can be defined by primitive recursion [Mee90]. In practice, this means that the result can depend not only on the recursive result, but also on the recursive occurrence of the type. For an inductive type A , given a function $g :: F (B \times A) \rightarrow B$, $\langle g \rangle_A :: A \rightarrow B$ denotes a paramorphism over that type that produces values of type B . Again, its recursive definition can be expressed by a

$\langle in_A \rangle_A = id$	<i>reflex-CATA</i>
$\langle g \rangle_A \circ in_A = g \circ F \langle g \rangle_A$	<i>cancel-CATA</i>
$f \circ \langle g \rangle_A = \langle h \rangle_A \Leftarrow f \circ g = h \circ F f$	<i>fusion-CATA</i>
$\langle in_{[A]} \circ (id + f \times id) \rangle_{[A]} = map f$	<i>map-CATA</i>
$\langle in_A \circ F \pi_1 \rangle_A = id$	<i>reflex-PARA</i>
$\langle g \rangle_A \circ in_A = g \circ F (\langle g \rangle_A \triangle id)$	<i>cancel-PARA</i>
$f \circ \langle g \rangle_A = \langle h \rangle_A \Leftarrow f \circ g = h \circ F (f \times id)$	<i>fusion-PARA</i>
$\langle f \circ F \pi_1 \rangle_A = \langle f \rangle_A$	<i>cata-PARA</i>

Fig. 3. Some laws for folds and paramorphisms.

diagram.

$$\begin{array}{ccc}
 A & \xrightarrow{out_A} & F A \xrightarrow{F(id \triangle id)} F(A \times A) \\
 \langle g \rangle_A \downarrow & & \downarrow F(\langle g \rangle_A \times id) \\
 B & \xleftarrow{g} & F(B \times A)
 \end{array}$$

Notice how a copy of the recursive occurrence is made before the recursive invocation. For optimization of functions defined as folds and paramorphisms we will use the laws presented in Figure 3.

When applied to an inductive type, the bottom-up traversal *everywhere* will be specialized into a fold over that type. The *everywhere-APPLY* law will now be defined as follows.

$$apT_A (everywhere f) = \langle apT_A f \circ in_A \circ apT_{F \overline{A}} (everywhere f) \rangle_A$$

The behavior of this fold is better understood with the help of the following diagram:

$$\begin{array}{ccc}
 A & \xrightarrow{out_A} & F A \\
 \langle \cdot \rangle_A \downarrow & & \downarrow F \langle \cdot \rangle_A \\
 A & \xleftarrow{apT_A f} A \xleftarrow{in_A} F A \xleftarrow{apT_{F \overline{A}} (everywhere f)} F A &
 \end{array}$$

The intent of the function $apT_{F \overline{A}} (everywhere f)$ is to apply the transformation to all content of the type, apart from its recursive occurrences (which were already processed recursively by the fold itself). This behavior is achieved by adding the following law to the set presented in Figure 2:

$$apT_{\overline{A}} f = id \quad \text{rec-APPLY}$$

This law guarantees that a type marked with an overline is ignored by the apT combinator. For example, for lists $apT_{F \overline{A}} (everywhere f)$ would be instantiated as $apT_{1+A \times \overline{[A]}} (everywhere f)$, which is equivalent to $id + apT_A (everywhere f) \times id$.

Since *everywhere f* is a bottom-up traversal, after transforming both the recursive occurrences and the remaining content, f still needs to be applied to the resulting value. To do so, in_A is first used to reconstruct a value of type A , followed by an application of $apT_A f$.

To exemplify the specialization of a generic transformation to a recursive type, consider the following example, where $f = mkT_{Int} succ$:

$$\begin{aligned}
 & apT_{[Int]} (\text{everywhere } f) \\
 = & \{ \text{everywhere-APPLY} \} \\
 & (\llbracket apT_{[Int]} (mkT_{Int} \text{ succ}) \circ in_{[Int]} \circ apT_{F \overline{[Int]}} (\text{everywhere } f) \rrbracket_{[Int]}) \\
 = & \{ mkT\text{-APPLY}; +\text{-APPLY}; \times\text{-APPLY}; \text{rec-APPLY} \} \\
 & (\llbracket id \circ in_{[Int]} \circ (id + apT_{Int} (\text{everywhere } f) \times id) \rrbracket_{[Int]}) \\
 = & \{ \text{everywhere-APPLY} \} \\
 & (\llbracket in_{[Int]} \circ (id + apT_{Int} (gmapT (\text{everywhere } f) \triangleright f) \times id) \rrbracket_{[Int]}) \\
 = & \{ \triangleright\text{-APPLY}; gmapT\text{-APPLY} \} \\
 & (\llbracket in_{[Int]} \circ (id + (id \circ apT_{Int} (mkT_{Int} \text{ succ})) \times id) \rrbracket_{[Int]}) \\
 = & \{ mkT\text{-APPLY} \} \\
 & (\llbracket in_{[Int]} \circ (id + \text{succ} \times id) \rrbracket_{[Int]}) \\
 = & \{ \text{map-CATA} \} \\
 & \text{map succ}
 \end{aligned}$$

As expected, if f is applied to a type that does not contain integers the result is the identity function:

$$\begin{aligned}
 & apT_{[Char]} (\text{everywhere } f) \\
 = & \{ \dots \} \\
 & (\llbracket in_{[Char]} \circ (id + (id \circ apT_{Char} (mkT_{Int} \text{ succ})) \times id) \rrbracket_{[Char]}) \\
 = & \{ mkT\text{-APPLY} \} \\
 & (\llbracket in_{[Char]} \circ (id + id \times id) \rrbracket_{[Char]}) \\
 = & \{ id \times id = id; id + id = id \} \\
 & (\llbracket in_{[Char]} \rrbracket_{[Char]}) \\
 = & \{ \text{reflex-CATA} \} \\
 & id
 \end{aligned}$$

The bottom-up *everything* combinator will be specialized into a paramorphism:

$$\begin{aligned}
 & apQ_A (\text{everything } f) = \\
 & \llbracket plus \circ (apQ_{F\overline{R}} (\text{everything } f) \times apQ_A f) \circ (F \pi_1 \Delta in_A \circ F \pi_2) \rrbracket_A
 \end{aligned}$$

Again, this paramorphism is easier to understand with a diagram:

$$\begin{array}{ccccc}
 A & \xrightarrow{\text{out}_A} & F A & \xrightarrow{F (id \Delta id)} & F (A \times A) \\
 \Downarrow \llbracket \cdot \rrbracket_A & & & & \downarrow F (\llbracket \cdot \rrbracket_A \times id) \\
 R & \xleftarrow{\text{plus}} & R \times R & \xleftarrow{apQ_{F\overline{R}} (\text{everything } f) \times apQ_A f} & F R \times A & \xleftarrow{F \pi_1 \Delta in_A \circ F \pi_2} & F (R \times A)
 \end{array}$$

After recursion, the input value is reconstructed using in_A in order to feed it to the generic query. Simultaneously, the query is applied to the non-recursive type contents, and finally both results are put together with the monoid *plus* operator.

4 Encoding in Haskell

In order to harness the above algebraic laws into a type-safe rewrite system for the specialization of generic functions, we must provide type-safe representations for

both functions and types. For functions we will use the same representation presented in [CV07b], based on a *generalized algebraic data type* (GADT) [JVWW06]:

data $PF\ a\ \mathbf{where}$

$Id \quad :: PF\ (a \rightarrow a)$
 $(\Delta) \quad :: PF\ (a \rightarrow b) \rightarrow PF\ (a \rightarrow c) \rightarrow PF\ (a \rightarrow (b, c))$
 $(\times) \quad :: PF\ (a \rightarrow b) \rightarrow PF\ (c \rightarrow d) \rightarrow PF\ ((a, c) \rightarrow (b, d))$
 $(\nabla) \quad :: PF\ (a \rightarrow c) \rightarrow PF\ (b \rightarrow c) \rightarrow PF\ (Either\ a\ b \rightarrow c)$
 $(+)$ $:: PF\ (a \rightarrow b) \rightarrow PF\ (c \rightarrow d) \rightarrow PF\ (Either\ a\ c \rightarrow Either\ b\ d)$
 $mkT \quad :: Type\ a \rightarrow PF\ (a \rightarrow a) \rightarrow PF\ T$
 $apT \quad :: Type\ a \rightarrow PF\ T \rightarrow PF\ (a \rightarrow a)$
 \dots

This type contains both point-free combinators and SYB combinators. Generic transformations and queries have the following types:

type $T = \forall a. Type\ a \rightarrow a \rightarrow a$
type $Q\ r = \forall a. Type\ a \rightarrow a \rightarrow r$

Instead of using type classes to infer type representations, these are explicitly passed to generic functions. Type representations are also defined using a GADT. For base types, sums and products, $Type\ a$ can be defined as follows.

data $Type\ a\ \mathbf{where}$

$Int \quad :: Type\ Int$
 $Char \quad :: Type\ Char$
 \dots
 $Sum \quad :: Type\ a \rightarrow Type\ b \rightarrow Type\ (Either\ a\ b)$
 $Prod \quad :: Type\ a \rightarrow Type\ b \rightarrow Type\ (a, b)$
 $Func \quad :: Type\ a \rightarrow Type\ b \rightarrow Type\ (a \rightarrow b)$

One consequence of using a GADT to encode the combinators is the ability to define an evaluation function $eval :: PF\ a \rightarrow a$. For example, the evaluation of mkT follows closely the SYB semantics:

$$eval\ (mkT\ a\ f) = \lambda b\ x \rightarrow \mathbf{case}\ teq\ a\ b\ \mathbf{of}\ Just\ Eq \rightarrow eval\ f\ x$$

$$Nothing \rightarrow x$$

The resulting generic function behaves as f if applied to a value of type a , or as the identity function otherwise. Function teq tests equality of type representations, and is nowadays a classical example of the usefulness of GADTs [JVWW06]:

data $Equal\ a\ b\ \mathbf{where}\ Eq \quad :: Equal\ a\ a$

$teq \quad :: Type\ a \rightarrow Type\ b \rightarrow Maybe\ (Equal\ a\ b)$
 $teq\ Int\ Int = return\ Eq$
 $teq\ (Sum\ a\ b)\ (Sum\ c\ d) = \mathbf{do}\ Eq \leftarrow teq\ a\ c$
 $Eq \leftarrow teq\ b\ d$
 $return\ Eq$

...
 $teq _ _ = Nothing$

The constructor Eq of the *Equal* GADT can be seen as a proof that types a and b are indeed equal.

The representation of user defined recursive types follows directly from the theoretical definition given in Section 2. For each type A we need to represent its base functor $F A$. When applying a functor to another type we want to get a sum of products as result, capable of being processed with point-free combinators. If functors are defined with normal Haskell polymorphic data types, it is impossible to obtain this behavior, since type-equivalence in Haskell is not structural but name-based.

In order to overcome this problem, we decided to represent functors using type-indexed type families [SSJC07,CKJ05], a new extension to the Haskell type system already supported in GHC. Developed with type-level programming in mind, type families are type constructors that represent sets of types. Set members are aggregated according to the type parameters passed to the type family constructor, called type indices: family constructors can have different representation types for different type indices. A type family to represent functors can be defined as follows.

type family $F a x :: *$

In this definition a is the type index that stands for the type whose functor is being defined, and x is the type argument of the functor itself. For example, for lists we have the following instance:

type instance $F [a] x = Either One (a, x)$

The GADT that represents functions can now be extended with constructors for the recursion patterns, together with the *in* and the *out* functions.

data $PF a$ **where**

...
 $In \quad :: PF (F a a \rightarrow a)$
 $Out \quad :: PF (a \rightarrow F a a)$
 $Cata :: PF (F a c \rightarrow c) \rightarrow PF (a \rightarrow c)$
 $Para :: PF (F a (c, a) \rightarrow c) \rightarrow PF (a \rightarrow c)$

In our rewrite system we will need to apply a functor both to type and function representations. Using the above type family, we can capture this behavior in the following data type, that represents the functor of an inductive type a .

data $Functor a = Functor\{ mapT :: \forall b. Type b \rightarrow Type (F a b),$
 $mapF :: \forall x y. PF (x \rightarrow y) \rightarrow PF (F a x \rightarrow F a y)\}$

Our type representation can now be extended with a new constructor to represent user defined recursive types.

```
data Type a where
...
Data :: String → Functor a → Type a
```

Given a ground type a , it is possible to use the Haskell type system to infer its representation. We can define a class with all representable types:

```
class Typeable a where
  typeof :: Type a
```

For example, for lists the *Typeable* instance can be defined as follows:

```
instance Typeable a ⇒ Typeable [a] where
  typeof = Data name functor
  where name = "[" ++ show (typeof :: Type a) ++ "]"
        functor = Functor{ mapT = λb → Sum One (Prod typeof b),
                          mapF = λf → Id + (Id × f)}
```

In order to guarantee that the rewrite system is type-safe, rewrite rules are represented by a monadic function that takes a function representation and returns a representation of the same type.

```
type Rule = ∀a. Type a → PF a → RewriteM (PF a)
```

RewriteM is a stateful monad that keeps a trace of the applied rules and is an instance of *MonadPlus*, thus modeling partiality in rule application. The extra type representation passed as argument allows the rule to make type-based rewriting decisions.

Both the specialization laws of Figure 2, and the point-free optimization laws (such as the ones presented in Figure 3 for folds and paramorphisms) are encoded as rewrite rules. For example, the reflexivity rule for folds can be defined as follows.

```
reflex_Cata :: Rule
reflex_Cata (Func a b) (Cata In) = do Eq ← teq a b
                                     success "reflex-Cata" Id
reflex_Cata _ _ = mzero
```

This rule uses *teq* to guarantee that is only applied to functions of type $a \rightarrow a$. The monadic function *success* updates the *RewriteM* monad to keep trace of the successful reduction.

Rewrite systems are built from basic rules using a standard set of strategic combinators. There are two main top-level strategies: *optimize_syb* for specialization of type-preserving and type-unifying generic programs into point-free expressions; and *optimize_pf* for simplification and optimization of point-free definitions. The latter also applies some “beautifying” rules to produce more concise results.

5 Examples

We will now present some specialization examples, and compare the performance between the resulting definitions and the original SYB functions.

The first example is the generic transformation to increase all salaries. In order to increase the readability of the specialized point-free definitions, we will consider that the type-specific behavior is for the *Employee* type instead of *Salary*. In SYB we have the following definition.

$$\begin{aligned} \text{increase} &:: \text{Int} \rightarrow \text{Company} \rightarrow \text{Company} \\ \text{increase } k &= \text{everywhere } (\text{mkT } (\text{incE } k)) \\ \textbf{where } \text{incE } k &(\text{E } p \text{ (S } s)) = \text{E } p \text{ (S } (s + k)) \end{aligned}$$

After encoding this definition using type $PF\ a$, and applying the specialization and optimization strategies, we get the following definition, where C stands for *Company* and D for *Dept*.

$$\text{inc}_C \circ \text{map } (\text{in}_D \circ (\text{id} \times (\text{incE } k \times \text{map } (\text{incE } k + \text{id})))) \text{)}_D \circ \text{out}_C$$

This definition is a pretty-print of the respective representation in $PF\ a$. It approximates the hand-written presented in Section 1: the fold will be recursively applied to each department of a company; at each department both the manager and all direct employees will have their salaries increased by the function $\text{incE } k$.

The second example, presented before in Section 2, addresses the specialization of a generic query to compute the total salary bill of a company.

$$\begin{aligned} \text{salaries} &:: \text{Company} \rightarrow \text{Int} \\ \text{salaries } k &= \text{everything } (\text{mkQ } \text{billE}) \\ \textbf{where } \text{billE } &(\text{E } p \text{ (S } s)) = s \end{aligned}$$

In this case we get the following definition.

$$\text{sum} \circ \text{map } (\text{plus} \circ (\text{billE} \times (\text{sum} \circ \text{map } (\text{billE } \nabla \text{id}))) \circ \pi_2)) \text{)}_D \circ \text{out}_C$$

Note how the paramorphism was simplified as a fold, since the query does not mention the recursive type *Dept*. The expression $\text{sum} \circ \text{map } (\text{billE } \nabla \text{id})$ collects all salaries from the direct employees of a department, and all recursively computed salaries from sub-departments. This result is then summed with the salary of the manager to compute the total salary bill. The expression $\text{sum} \circ \text{map } f$ was used in this example just to make the result more readable. In fact, they are fused together as a single fold by the rewrite system.

The last example combines the two previous examples into a single function:

$$\begin{aligned} \text{higher_salaries} &:: \text{Int} \rightarrow \text{Company} \rightarrow \text{Int} \\ \text{higher_salaries } k &= \text{salaries} \circ \text{increase } k \end{aligned}$$

Although the two operations are performed in sequence in the original query, after specialization we get a result very similar to the previous one, with a single traversal over the type.

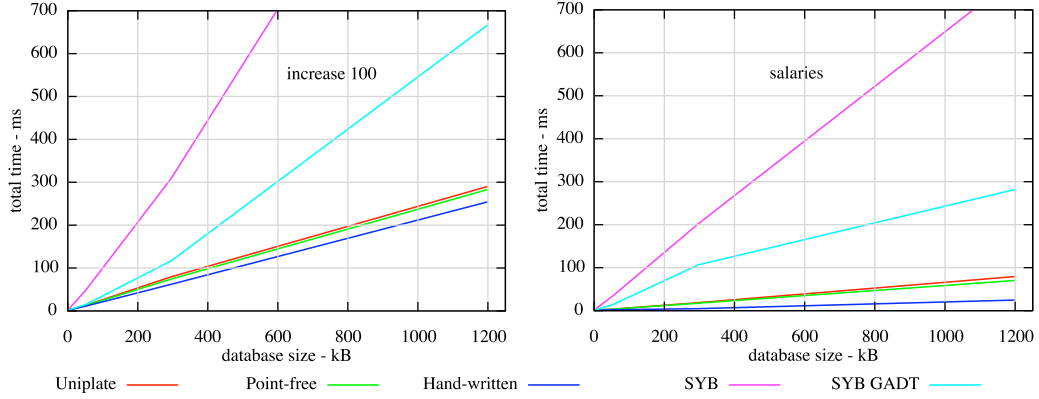


Fig. 4. Timing results.

```

let aux = bille ◦ incE k
in sum ◦ map (plus ◦ (aux × (sum ◦ map (aux ∇ id)) ◦ π2)) ◦ outC
    
```

Unlike systems specially designed to implement fusion (such as [SdM03]), our rewrite system cannot implement the full power of the fusion laws. However it covers most of the particular instances that occur during the specialization of generic functions. For example, the above optimization was possible due to the following instance of *fusion-CATA*.

$$(f)_A \circ (in_A \circ g)_A = (f \circ g)_A \Leftarrow F (f)_A \circ g = g \circ F (f)_A$$

To verify the side-condition of this law, we first apply the rewrite system *optimize_pf* to $F (f)_A \circ g$ and $g \circ F (f)_A$, and then check for syntactic equality.

Performance analysis

We have compared the runtimes of the first two examples for hand-written, specialized, and generic definitions written in SYB and Uniplate. The results are presented in Figure 5. A large part of SYB’s inefficiency is due to the heavy use of type-classes to infer type representations. To factor out this penalty, and better quantify the speedup achieved by our specialization mechanism, we also include the runtimes of both generic functions obtained by evaluating their representation using the *eval* function presented in section 4 (denoted in the graphic as SYB GADT). We compiled each function using GHC 6.8.2 with optimization flag `02`. Each example was tested with *Company* values of increasing size (measured in kBytes needed to store the Haskell definition of each value).

As expected, for both examples, the SYB generic definitions perform much worse than the hand-written, and the loss factor grows with the database size. The SYB GADT variant is at least twice as fast, but still much slower than the hand-written. The specialized point-free definitions perform closer to the hand-written, with loss factors of 1.11 (*increase 100*) and 2.85 (*salaries*) for the biggest sample. This performance loss is mainly due to the use of *in* and *out* to convert between user defined types and they structural representation as a sum of products. For these particular examples, the performance of the specialized point-free code is tangentially better than Uniplate. As discussed in the next section, Uniplate also has some mechanisms to avoid traversing unnecessary branches, which justify the proximity in the results.

Although quite standard when comparing generic programming libraries, these example are not particularly flattering to our optimization mechanism: in fact, there are no large branches of data that can be avoided in the traversals. For example, if the *Company* data type had any other information besides departments (not containing the type *Salary*), the runtime would remain the same, further widening the gap to SYB. We also achieve a significant advantage when optimizing compositions of generic functions: for example, in the *higher_salaries* example our specialized point-free definition was already 1.35 times faster than Uniplate for the biggest sample.

6 Related Work

Uniplate

Unlike SYB, some generic programming libraries have been designed with performance issues in mind, usually at the cost of expressiveness. One such library is Uniplate [MR07], that is among the fastest libraries currently available for generic programming in Haskell [RJJ⁺08]. That fact, together with the SYB-like flavor of its combinators, motivated an obvious inclusion in the comparative performance analysis of the previous section. The key idea behind Uniplate is that most generic traversals have value-specific behavior for just one type. Building on this insight, this library provides two key combinators to specify bottom-up generic transformations:

$$\begin{aligned} \text{transform} &:: \text{Uniplate } a \Rightarrow (a \rightarrow a) \rightarrow a \rightarrow a \\ \text{transformBi} &:: \text{Biplate } b \ a \Rightarrow (a \rightarrow a) \rightarrow b \rightarrow b \end{aligned}$$

The *transform* combinator applies its argument to every *a* occurring inside a value of type *a*, while *transformBi* applies its argument to every *a* occurring inside a value of a different type *b*. Recalling our examples, the *increase* transformation can be defined using *transformBi*, since it looks for all salaries inside a company. The *Uniplate* and *Biplate* classes contain primitive methods to find the substructures of type *a* inside values of type *a* and *b*, respectively. Instances of these classes can be defined using a variety of methods, ranging from more generic and less efficient to more verbose and more efficient. The most efficient method (used in the comparison of the previous section) is to define the instances by hand, which for *Biplate* requires defining n^2 instances to support *n* types. When defining instances for *Biplate b a* it is possible to avoid traversing down branches of *a* that do not contain the target type *b*, thus optimizing generic traversals.

The main advantage of *Uniplate* is that generic functions execute fast out of the box, without the need of an explicit optimization phase. On the other hand, likewise to SYB, our optimization technique can handle more powerful combinators, that target different types in a single traversal. Using fusion techniques, our approach can also further optimize combinations of traversals, while Uniplate speedups are constrained to individual traversals.

Other techniques to optimize generic programs

Another very efficient SYB-like generic programming library is Smash [Kis06]. Instead of using run-time checks to find the target types, it offsets them to compile-time by using heterogeneous collections [KLS04] to encode the type-specific cases of generic functions. Unfortunately, the speedup obtained with this technique comes at the cost of extra work from the programmer: in order to support a new data type, all different traversal combinators must be defined from scratch, while in SYB they can all be generically implemented using just two primitive methods.

A different approach has been followed in [AS04], where a technique named symbolic evaluation was developed to optimize Generic Haskell programs [LJC⁺05]. It focus on the specialization of fully applied functions and tries to eliminate conversions between types and their structural representations. Symbolic evaluation guarantees that the intermediate structures are completely removed from the optimized code. A similar technique could be used in our framework to further optimize the point-free definitions, via an additional translation step to explicitly recursive point-wise code.

Application scenarios

As previously mentioned, our main goal was to extend the specialization mechanism presented in [CV07b] to also cover inductive types. In [CV07b] we already described how it could be used to optimize the structure-shy XPath query language. This technique was harnessed into the prototype schema-aware XPath compiler XPTO [FP07]. Query compilation in XPTO proceeds as follows: the XML Schema is parsed into a sum of products representation using *Type a*; the XPath query is parsed into a type-safe representation of type *PF a*; the rewrite system is used to specialize the query to the given schema; the specialized point-free definition is output into a new Haskell program to be compiled and linked with an XML parser and point-free execution library; the resulting program can then be used to execute the original query against XML files conforming to the given schema. We are currently deploying the new technique presented here into the XPTO compiler in order to handle some recursive XML Schemas.

A similar type-safe rewriting system was also used in [CV07a] to optimize two-level data transformations [COV06]. A two-level data transformation consists of a type-level transformation coupled with value-level transformations of the respective inhabitants. More specifically, we developed a framework that allows us to specify data type refinements $A \leq B$ using strategic combinators, and get for free the migration functions between values of type *A* and *B*, and vice-versa. Both the types and the migration functions are again encoded using *Type a* and *PF a*, allowing us to use the rewrite system to optimize them, and migrate queries/producers from the abstract type *A* to the concrete type *B*. The inclusion of inductive types in the rewrite system will allow us to extend the applicability of this framework.

Template meta-programming

We believe that our algebraic approach could be instructed at a lower level to provide compile-time specialization of generic functions, through template meta-programming [SJ02]. This rewriting process would encompass transformation of

Haskell generic programs through direct manipulation of their abstract syntax trees. However, since the current implementation of template meta-programming in Haskell is completely untyped, we would lose the guarantee that the rewrite system is type-safe. Template meta-programming could also be used to infer automatically the recursive types' sum of products representation.

7 Concluding Remarks

We have extended an existent mechanism to specialize SYB-like generic functions to also cover user defined recursive data types. By focusing on inductive types (fixpoints of functors) we were able to use recursion patterns such as folds and paramorphisms to encode generic traversals. These recursion patterns are characterized by nice algebraic laws, that were incorporated in a type-safe rewrite system to further optimize the specialized code. The definitions produced by our specialization mechanism perform close to hand-written non-generic ones. Thanks to recent extensions of the Haskell type-system, such as type-indexed type families or generalized algebraic data types, our implementation of the rewrite system closely mimics the theoretical presentation.

The major limitation of the current approach is that it only supports single-recursive inductive types. We are currently investigating how to extend it to cover more general forms of recursion, such as mutually-inductive data types or nested data types. Particularly relevant to this endeavor is the work described in [JG07], showing that higher-order functors can be used to give an initial algebra semantics to nested data types (likewise to standard inductive types).

References

- [AS04] Artem Alimarine and Sjaak Smetsers. Optimizing generic functions. In *Proceedings of the 7th International Conference on the Mathematics of Program Construction*, volume 3125 of *LNCS*, pages 16–31. Springer-Verlag, 2004.
- [Bac78] John Backus. Can programming be liberated from the von Neumann style? a functional style and its algebra of programs. *Communications of the ACM*, 21(8):613–641, 1978.
- [CKJ05] Manuel M. T. Chakravarty, Gabriele Keller, and Simon Peyton Jones. Associated type synonyms. In *Proceedings of the 10th ACM SIGPLAN international conference on Functional programming*, pages 241–253. ACM Press, 2005.
- [COV06] Alcino Cunha, José Nuno Oliveira, and Joost Visser. Type-safe two-level data transformation. In J. Misra, T. Nipkow, and E. Sekerinski, editors, *Proceedings of the 14th International Symposium on Formal Methods*, volume 4085 of *LNCS*, pages 284–299. Springer-Verlag, 2006.
- [CPP06] Alcino Cunha, Jorge Sousa Pinto, and José Proença. A framework for point-free program transformation. In A. Butterfield, C. Grellck, and F. Huch, editors, *Selected Papers of the 17th International Workshop on Implementation and Application of Functional Languages*, volume 4015 of *LNCS*, pages 1–18. Springer-Verlag, 2006.
- [CV07a] Alcino Cunha and Joost Visser. Strongly typed rewriting for coupled software transformation. *ENTCS*, 174(1):17–34, 2007. Proceedings of the 7th International Workshop on Rule-Based Programming.
- [CV07b] Alcino Cunha and Joost Visser. Transformation of structure-shy programs - applied to XPath queries and strategic functions. In *Proceedings of the ACM SIGPLAN 2007 Workshop on Partial Evaluation and Program Manipulation*, pages 11–20. ACM Press, 2007.
- [FP07] Flávio Ferreira and Hugo Pacheco. XPTO - an Xpath preprocessor with type-aware optimization. In *Proceedings of CORTA'07: Compilers, Related Technologies and Applications*, February 2007.

- [Gib02] Jeremy Gibbons. Calculating functional programs. In R. Backhouse, R. Crole, and J. Gibbons, editors, *Algebraic and Coalgebraic Methods in the Mathematics of Program Construction*, volume 2297 of *LNCS*, chapter 5, pages 148–203. Springer-Verlag, 2002.
- [JG07] Patricia Johann and Neil Ghani. Initial algebra semantics is enough! In *Proceedings of the 8th International Conference on Typed Lambda Calculi and Applications*, volume 4583 of *LNCS*, pages 207–222. Springer-Verlag, 2007.
- [JVWW06] Simon Peyton Jones, Dimitrios Vytiniotis, Stephanie Weirich, and Geoffrey Washburn. Simple unification-based type inference for GADTs. *SIGPLAN Notices*, 41(9):50–61, 2006. Proceedings of the 11th ACM SIGPLAN International Conference on Functional Programming.
- [Kis06] Oleg Kiselyov. Smash your boilerplate without class and typeable. <http://article.gmane.org/gmane.comp.lang.haskell.general/14086>, 2006.
- [KLS04] Oleg Kiselyov, Ralf Lämmel, and Kean Schupke. Strongly typed heterogeneous collections. In *Proceedings of the 2004 ACM SIGPLAN Haskell Workshop*, pages 96–107. ACM Press, 2004.
- [LJ03] Ralf Lämmel and Simon Peyton Jones. Scrap your boilerplate: a practical design pattern for generic programming. *SIGPLAN Notices*, 38(3):26–37, 2003. Proceedings of the ACM SIGPLAN Workshop on Types in Language Design and Implementation.
- [LJ05] Ralf Lämmel and Simon Peyton Jones. Scrap your boilerplate with class: extensible generic functions. *SIGPLAN Notices*, 40(9):204–215, 2005. Proceedings of the 10th ACM SIGPLAN International Conference on Functional Programming.
- [LJC⁺05] Andres Löh, Johan Jeuring, Dave Clarke, Ralf Hinze, Alexey Rodriguez, and Jan de Wit. The generic haskell user’s guide. Technical Report UU-CS-2005-004, Utrecht University, 2005.
- [Mee90] Lambert Meertens. Paramorphisms. Technical Report RUU-CS-90-4, Utrecht University, Department of Computer Science, January 1990.
- [MFP91] Erik Meijer, Maarten Fokkinga, and Ross Paterson. Functional programming with bananas, lenses, envelopes and barbed wire. In J. Hughes, editor, *Proceedings of the 5th ACM Conference on Functional Programming Languages and Computer Architecture*, volume 523 of *LNCS*, pages 124–144. Springer-Verlag, 1991.
- [MR07] Neil Mitchell and Colin Runciman. Uniform boilerplate and list processing. In *Proceedings of the 2007 ACM SIGPLAN Haskell Workshop*, pages 49–60. ACM Press, 2007.
- [RJJ⁺08] Alexey Rodriguez, Johan Jeuring, Patrik Jansson, Alex Gerdes, Oleg Kiselyov, and Bruno C. d. S. Oliveira. Comparing libraries for generic programming in Haskell. Technical Report UU-CS-2008-010, Utrecht University, 2008.
- [Sdm03] Ganesh Sittampalam and Oege de Moor. Mechanising fusion. In J. Gibbons and O. de Moor, editors, *The Fun of Programming*, chapter 5, pages 79–104. Palgrave Macmillan, 2003.
- [SJ02] Tim Sheard and Simon Peyton Jones. Template meta-programming for Haskell. In *Proceedings of the 2002 ACM SIGPLAN Haskell Workshop*, pages 1–16. ACM Press, 2002.
- [SSJC07] Tom Schrijvers, Martin Sulzmann, Simon Peyton Jones, and Manuel M. T. Chakravarty. Towards open type functions for haskell. In *Proceedings of the 19th International Symposium on Implementation and Application of Functional Languages*, pages 233–251, 2007.