

Optimizing SYB Is Easy!

Michael D. Adams^{a,*}, Andrew Farmer^b, José Pedro Magalhães^c

^a*Department of Computer Science, University of Illinois at Urbana-Champaign*

^b*Information and Telecommunication Technology Center, University of Kansas*

^c*Department of Computer Science, University of Oxford*

Abstract

The most widely used generic-programming system in the Haskell community, Scrap Your Boilerplate (SYB), also happens to be one of the slowest. Generic traversals in SYB are often an order of magnitude slower than equivalent handwritten, non-generic traversals. Thus while SYB allows the concise expression of many traversals, its use incurs a significant runtime cost. Existing techniques for optimizing other generic-programming systems are not able to eliminate this overhead.

This paper presents an optimization that completely eliminates this cost. Essentially, it is a partial evaluation that takes advantage of domain-specific knowledge about the structure of SYB. It optimizes SYB-style traversals to be as fast as handwritten, non-generic code, and benchmarks show that this optimization improves the speed of SYB-style code by an order of magnitude or more.

Keywords: optimization, partial evaluation, datatype-generic programming, Haskell, Scrap Your Boilerplate (SYB), performance

2010 MSC: 68N20

1. Introduction

1 Scrap Your Boilerplate (SYB) (Lämmel and Peyton Jones, 2003, 2004) is
2 one of the oldest and most widely used systems for generic programming in
3 Haskell. It is the most downloaded package for generic programming in the
4 Hackage archive (Industrial Haskell Group, 2013). It is easy to use and has
5 strong support from the Glasgow Haskell Compiler (GHC) (GHC Team, 2013).
6

7 While SYB allows the easy and concise expression of traversals that oth-
8 erwise require large amounts of handwritten code, it has a serious drawback,
9 namely, poor runtime performance. Our own benchmarks show it to be an or-
10 der of magnitude slower than handwritten, non-generic code, and this fact is

*Corresponding author

Email addresses: afarmer@itc.ku.edu (Andrew Farmer), jpm@cs.ox.ac.uk (José Pedro Magalhães)

URL: <http://michaeldadams.org> (Michael D. Adams)

11 documented many times in the literature (Rodriguez Yakushev, 2009; Brown
12 and Sampson, 2009; Chakravarty et al., 2009; Magalhães et al., 2010; Adams
13 and DuBuisson, 2012; Sculthorpe et al., 2014).

14 While attempts have been made in the past to use general-purpose optimiza-
15 tions to improve the performance of SYB, they have met with only moderate
16 success. For example, while setting the compiler’s optimizer to be exceptionally
17 aggressive about unfolding and inlining can slightly improve the performance
18 of SYB, doing so can harm the performance of the program as a whole as code
19 may be inlined that should not be (Magalhães et al., 2010).

20 Nevertheless, SYB-style code exhibits a structure that we can take advantage
21 of in our optimizations. This paper presents a domain-specific optimization that
22 transforms SYB-style code to be as fast as handwritten code. This optimiza-
23 tion uses the types of expressions to direct where the inlining process should
24 be more aggressive. In essence, it is a specialized form of supercompilation
25 (Turchin, 1979) and partial evaluation (Jones et al., 1993) that uses type in-
26 formation to determine whether an expression should be computed statically
27 at compile time or dynamically at runtime. Using this technique and domain-
28 specific knowledge about the structure of the SYB library and the code that
29 uses it, we show that optimizing SYB-style code can be easily implemented
30 with standard transformations.

31 This optimization was first implemented using HERMIT (Farmer et al., 2012;
32 Sculthorpe et al., 2013), an interactive optimization system implemented as a
33 GHC plugin. HERMIT makes it easy to quickly develop these sorts of optimiza-
34 tions, as a prelude for improving the GHC optimizer with similar techniques.
35 Afterwards, we used the information gained while developing our HERMIT plu-
36 gin to improve the GHC optimizer, and obtained equally promising results,
37 but without the need for a HERMIT script. The code for our optimization
38 is available at <https://github.com/xich/hermit-syb> and as the `hermit-syb`
39 package on Hackage.

40 This paper is a rewritten and extended version of our earlier work (Adams
41 et al., 2014). The implementation of our optimization in GHC (together with the
42 description of the changes necessary to the optimizer and SYB) is entirely new
43 to this version. We have also revisited our HERMIT script and its benchmarks
44 to find and eliminate cases of poor performance.

45 The remainder of this paper is organized as follows. We start with an
46 overview of SYB in Section 2. In Section 3 we show a step-by-step “manual”
47 optimization of an SYB program. This is followed by a formal description of
48 our optimization in Section 4. In Section 5 we discuss an implementation of the
49 optimization for GHC and present benchmarks validating its effectiveness. This
50 is followed in Section 6 by a discussion of the limitations and future work for
51 our system. In Section 7, we consider the GHC specializer and how it can be
52 adapted to achieve the same optimization. Finally, we review related work in
53 Section 8, and conclude in Section 9.

54 2. Overview of SYB

55 In order to understand why SYB is slow, we must first understand how it
 56 works. SYB is a generic-programming system for concisely expressing traversals.
 57 For example, suppose we have a type of abstract syntax trees, `AST`, and wish
 58 to apply a name mangling function, `mangle`, to every identifier in a given `AST`.
 59 Writing this by hand requires a large amount of “boilerplate” code that merely
 60 recurs until we get to an identifier where we can apply `mangle`. With SYB,
 61 however, we can use the `everywhere` and `mkT` functions to write this traversal
 62 simply as `everywhere (mkT mangle)`.

63 SYB defines many traversals in addition to `everywhere`, and the optimiza-
 64 tion presented in this paper handles these, but for the sake of simplicity our
 65 examples focus on the `everywhere` traversal. In addition, since traversals over
 66 an `AST` type can be unwieldy, we use the following traversal over slightly simpler
 67 types as our running example.

```
68     inc :: Int -> Int
69     inc n = n + 1
70
71     incrementsyb :: [Int] -> [Int]
72     incrementsyb x = everywhere (mkT inc) x
```

73 This traversal applies `inc` to every object in `x` that has type `Int` and thus
 74 increments every integer in a list of integers.

75 We now turn to how `mkT` and `everywhere` work before considering the ques-
 76 tion of why this SYB-style traversal is slower than an equivalent handwritten
 77 traversal.

78 2.1. Transformations

79 The `mkT` function applies a transformation `f` to a term `x` if the types are
 80 compatible. Otherwise, it behaves as an identity function and simply returns
 81 `x`. Its definition relies on the type-safe casting function `cast`, which in turn is
 82 defined in terms of the `typeOf` method provided by the `Typeable` class. The
 83 implementation of these functions is equivalent to the following although the
 84 actual implementation of `mkT` goes through several intermediate helper functions
 85 that are not shown here.

```
86     mkT :: (Typeable a, Typeable b) => (b -> b) -> a -> a
87     mkT f = case cast f of
88               Nothing -> id
89               Just g  -> g
90
91     cast :: (Typeable a, Typeable b) => a -> Maybe b
92     cast x = r where
93       r = if typeOf x == typeOf (fromJust r)
94           then Just (unsafeCoerce x)
95           else Nothing
```

96 The `typeOf` function used in this code returns a value of type `TypeRep` repre-
 97 senting the type of its argument. The value of its argument is ignored. The
 98 `unsafeCoerce` function has type $\forall a\ b. a \rightarrow b$ and unconditionally coerces a
 99 value of one type to another type. Assuming the `Typeable` instances are correct,
 100 this use of `unsafeCoerce` is safe because of the check that the types `a` and `b`
 101 are indeed the same.

102 2.2. Traversals

103 The `everywhere` function traverses a structure in a bottom-up fashion and
 104 is implemented as follows.

```
105     everywhere :: (forall b. Data b => b -> b)
106                -> (forall a. Data a => a -> a)
107     everywhere f x = f (gmapT (everywhere f) x)
```

108 It uses `gmapT` to apply `everywhere f` to every subterm of `x`, and afterwards
 109 it applies `f` to the result. The `gmapT` function applies a transformation to all
 110 the immediate subterms of a given term, and we discuss its implementation
 111 in Section 2.3. It does not itself recurse past the first layer of children, but by
 112 calling it with `everywhere f` as an argument, the `everywhere` function recurses
 113 to all the descendants of `x` in a bottom-up fashion.

114 2.3. Mapping subterms

115 The type of `gmapT` is the same as that of `everywhere`. The important
 116 difference is that `gmapT` is not recursive, and transforms only the immediate
 117 subterms of a term. For any constructor `C` with n arguments, `gmapT` obeys the
 118 following equality.

```
119     gmapT f (C x1...xn) = C (f x1) ... (f xn)
```

120 The function `gmapT` is a method of the `Data` class, and has a default implemen-
 121 tation in terms of the SYB primitive `gfoldl`, which has the following type.

```
122     gfoldl :: (Data a)
123            => (forall d b. Data d => c (d -> b) -> d -> c b)
124            -> (forall g. g -> c g) -> a -> c a
```

125 This is a method of the `Data` class so its implementation is different for every
 126 type, but the general structure of such implementations can be seen in the
 127 following class instance for lists.

```
128     instance Data a => Data [a] where
129         gfoldl k z []      = z []
130         gfoldl k z (x:xs) = z (:) 'k' x 'k' xs
```

131 The `gfoldl` function takes three arguments. The first, `k`, combines an argu-
 132 ment with the constructor. The second, `z`, is applied to the constructor itself.
 133 Finally, the third is the value over which the `gfoldl` method traverses. The
 134 implementation always follows the same pattern. For any constructor `C` with n
 135 arguments, `gfoldl` obeys the following equality.

```
136   gfoldl k z (C x1...xn) = z C 'k' x1 ... 'k' xn
```

137 While extremely general, `gfoldl` is not easy to use directly. However, generic
 138 functions such as `gmapT` that are easier to use can be built in terms of it. Re-
 139 turning to `gmapT`, its default implementation is defined in terms of `gfoldl` as
 140 follows.

```
141   gmapT :: (∀b. Data b => b -> b)
142         -> (∀a. Data a => a -> a)
143   gmapT f x = unID (gfoldl k ID x) where
144     k (ID c) y = ID (c (f y))
145
146   newtype ID x = ID { unID :: x }
```

147 Since `gmapT` does not need to take advantage of the type changing ability pro-
 148 vided by the `c` type parameter to `gfoldl`, it instantiates `c` to the trivial type `ID`.
 149 Aside from wrapping and unwrapping `ID`, `gmapT` operates by using `k` to rebuild
 150 the constructor application after applying `f` to each constructor argument and
 151 thus obeys the previously given equality for `gmapT`.

152 2.4. Why SYB is slow

153 The slow performance of SYB is well documented. Rodriguez Yakushev
 154 (2009, Figure 4.9) benchmarked three SYB functions, and found them to be
 155 36, 52, and 69 times slower than handwritten code. Chakravarty et al. (2009)
 156 also benchmark SYB on three functions, finding them to be 45, 73, and 230
 157 times slower than handwritten code. Brown and Sampson (2009) developed a
 158 new generic-programming library because SYB was too slow and found SYB
 159 to be 4 to 23 times slower than their own approach. Magalhães et al. (2010)
 160 report SYB performing between 3 and 20 times slower than handwritten code.
 161 Adams and DuBuisson (2012) developed an optimized variant of SYB using
 162 Template Haskell and report SYB performing between 10 and nearly 100 times
 163 slower than handwritten code. Sculthorpe et al. (2014) benchmark SYB on two
 164 generic traversals, finding it to be around 5 times slower than handwritten code.
 165 All of these papers conclude that SYB is one of the slowest generic-programming
 166 libraries.

167 After analyzing how SYB works, these results should not be surprising. Con-
 168 sider for example, the runtime behavior of the `incrementSYB` function. When
 169 applied to a value of type `[Int]` such as `[0,1]`, it recurses down the structure
 170 while applying `mkT inc` to every subterm. In this case, there are five subterms.
 171 Three of them are the lists `[0,1]`, `[1]` and `[]`. The remaining two are the `Int`

172 values 0 and 1. For each subterm, `mkT` attempts to cast `inc` to have a type
 173 that is applicable to that subterm. On the lists, it fails to do so, and thus `mkT`
 174 returns them unchanged. On the `Int` values, however, the cast succeeds, and
 175 thus `mkT` applies `inc` to them. This process involves significant overhead as it
 176 uses five dynamic type checks in order to update only two values.

177 Existing techniques for optimizing other generic-programming libraries are
 178 unable to eliminate this overhead in SYB-style code. Since SYB relies heavily
 179 on runtime type comparison, the type specializer cannot guide the optimiza-
 180 tion as it does in the work of Magalhães (2013). Instead, in order to find out
 181 if `inc` can be applied to a term, we must inline `mkT`, `cast`, and the `Typeable`
 182 methods all the way to the comparison of the type representation computed
 183 for the type of a term. If all of those are appropriately inlined, `mkT inc` re-
 184 duces to either `inc` or `id` depending on whether the types match. However,
 185 the GHC inliner (Peyton Jones and Marlow, 2002), while often eager to inline
 186 small expressions, will not perform as aggressive an inlining as is required here.
 187 Coercing GHC to inline aggressively has the side-effect of inlining parts of the
 188 code that were not intended to be inlined (Magalhães et al., 2010). Furthermore,
 189 because `everywhere` is a recursive function, GHC avoids inlining it in order to
 190 ensure termination of the inlining process. Even if GHC would inline recur-
 191 sive definitions, it would have to do so in a way that avoids infinitely inlining
 192 nested recursive occurrences. Implementing these optimizations would require
 193 fundamental changes to the way the inliner behaves, and their applicability to
 194 non-SYB-style code is not clear.

195 3. Optimizing SYB-style code

196 In order to gain an intuition for optimizing SYB-style code, we now consider
 197 the `incrementSYB` function from Section 2 and how we can manually transform
 198 it into non-generic code. Our goal is to reach the following more efficient non-
 199 generic implementation that avoids the runtime casts and dictionary dispatches
 200 that slow down the code as discussed in Section 2.4.

```
201     incrementHand :: [Int] -> [Int]
202     incrementHand []      = []
203     incrementHand (x : xs) = inc x : incrementHand xs
```

204 In order to optimize `incrementSYB`, we can exploit the fact that, due to the
 205 types of `incrementSYB` and `inc`, the concrete types and dictionaries needed by
 206 `everywhere` and `mkT` are known at compile time. These can be aggressively
 207 inlined, yielding code without any dynamic type checks or runtime casts. In
 208 Haskell, type and dictionary arguments are implicit. In order to make them
 209 explicit, we represent `incrementSYB` in terms of `Core`, which is the intermediate
 210 representation on which GHC does most of its optimizations. The result is the
 211 following.

```
incrementSYB :: [Int] -> [Int]
incrementSYB = λ x →
```

```

everywhere
  (λ b0 $dData0 →
    mkT Int b0 ($p1Data b0 $dData0)
      $fTypeableInt inc)
[Int] $dData x

```

212 Explicit type arguments are highlighted here in green, and we elide type co-
 213 coercions as they make the code difficult to read. In the following, we also skip
 214 many intermediate transformations as the full derivation requires several hun-
 215 dred steps.

216 In this code, the `$dData` and `$dTypeableInt` variables are `Data` and `Typeable`
 217 dictionaries that were previously implicit. The `$p1Data b0 $dData0` expression
 218 computes the `Typeable` dictionary corresponding to `$dData0`. We will see more
 219 such expressions as we proceed.

220 Since the dynamic type checks in `mkT` cause this code to be slow, we could
 221 try inlining `mkT` immediately. However, we would not have enough information
 222 to eliminate these checks if we did so as `b0` and `$dData0` do not yet have values
 223 and thus we do not know enough about the arguments to which `mkT` is applied.
 224 Instead, in order to get the λ -expression containing `mkT` to a fully applied po-
 225 sition, we inline `everywhere`, the function to which it is an argument. This
 226 results in the following.

```

incrementSYB :: [Int] -> [Int]
incrementSYB = λ x →
  mkT Int [Int] ($p1Data [Int] $dData) $fTypeableInt inc
  (gmapT [Int] $dData
    (λ b1 $dData1 →
      everywhere
        (λ b0 $dData0 →
          mkT Int b0 ($p1Data b0 $dData0)
            $fTypeableInt inc)
        b1 $dData1)
    x)

```

227 The call to `mkT` at the beginning of this code can now be inlined, and this exposes
 228 a call to `cast`.

```

incrementSYB :: [Int] -> [Int]
incrementSYB =
  let $dTypeable4 = ...
      $dTypeable5 = ...
  in λ x →
    (case cast (Int -> Int) ([Int] -> [Int])
      $dTypeable5 $dTypeable4 inc of
      Nothing → id [Int]
      Just g0 → g0)
    (gmapT [Int] $dData

```

```

      (λ b1 $dData1 →
        everywhere
          (λ b0 $dData0 →
            mkT Int b0 ($p1Data b0 $dData0)
              $fTypeableInt inc)
          b1 $dData1)
    x)

```

229 This code attempts to cast `inc` from type `Int -> Int` to type `[Int] ->`
 230 `[Int]` by using the `cast` function. Inlining `cast` exposes calls to `typeOf` that
 231 we can symbolically evaluate. After several more simplification steps, this call to
 232 `cast` reduces to `Nothing`, and in turn the `case` statement can be reduced to the
 233 identity function. Thus, we have removed one of the runtime type comparisons
 234 that slow down this code, and after simplification, the code now looks like the
 235 following.

```

incrementSYB :: [Int] -> [Int]
incrementSYB = λ x →
  gmapT [Int] $dData
    (λ b1 $dData1 →
      everywhere
        (λ b0 $dData0 →
          mkT Int b0 ($p1Data b0 $dData0)
            $fTypeableInt inc)
        b1 $dData1)
  x

```

236 At this point, the outer `mkT` has gone away completely. This is to be expected
 237 as `mkT` applied `inc` to only `Int` values, but at the outer level it is being applied
 238 to a `[Int]` value in which case `mkT` is an identity.

239 Here again we choose not to inline the λ -expression containing `mkT` since it
 240 is not fully applied. Instead we inline `gmapT` and get the following code.

```

incrementSYB :: [Int] -> [Int]
incrementSYB = λ x →
  case x of
    [] → [] Int
    (:) x0 xs0 →
      (:) Int
      (everywhere
        (λ b0 $dData0 →
          mkT Int b0 ($p1Data b0 $dData0)
            $fTypeableInt inc)
        Int $fDataInt x0)
      (everywhere
        (λ b0 $dData0 →
          mkT Int b0 ($p1Data b0 $dData0)

```



```

    $fTypeableInt inc)
  [Int] $dData xs0)

```

241 Since the eliminated `gmapT` is a class method, this inlining is particular to the
 242 type at which `gmapT` is applied. In this case it is over the list type, and `gmapT`
 243 inlines to a `case` expression over lists. As this `case` expression corresponds to the
 244 one in `incrementHand`, we can now recognize the structure of `incrementHand`
 245 becoming manifest in the code.

246 The code now contains two calls to `everywhere` that are inside the `(:)`
 247 branch of the `case` expression. One is on the head of the list and is at the type
 248 `Int`. The other is on the tail of the list and is at the type `[Int]`. We can inline
 249 the first of these which results in calls to `mkT` and `gmapT` just as before. This
 250 time, however, they are over the `Int` type. Thus, not only does the `cast` in `mkT`
 251 succeed and the `mkT` reduce to `inc`, but the call to `gmapT` reduces to the identity
 252 function. After a bit of simplification, the code now looks like the following.

```

incrementSYB :: [Int] -> [Int]
incrementSYB = λ x →
  case x of
    [] → [] Int
    (:) x0 xs0 →
      (:) Int
        (inc x0)
        (everywhere
          (λ b0 $dData0 →
            mkT Int b0 ($p1Data b0 $dData0)
              $fTypeableInt inc)
          [Int] $dData xs0)

```

253 Thus far we have eliminated several runtime costs merely by inlining and some
 254 basic simplifications, and this has brought us close to our goal of transforming
 255 `incrementSYB` into `incrementHand`. The only generic part of the code that
 256 remains is the call to `everywhere` on the tail of the list. While it is tempting to
 257 also inline this call, this expression is the same one that `incrementSYB` started
 258 with, and continuing to inline will thus lead us in a loop. Instead, we can take
 259 advantage of the fact that `incrementSYB` equals this expression and replace it
 260 with a reference to `incrementSYB`. Once we perform that replacement, we get
 261 the following code, which is identical to that of `incrementHand`.

```

incrementSYB :: [Int] -> [Int]
incrementSYB = λ x →
  case x of
    [] → [] Int
    (:) x0 xs0 → (:) Int (inc x0) (incrementSYB xs0)

```

$e, u := x$	Variables
l	Literals
$\Lambda a : \kappa. e \mid e \tau$	Type abstraction and application
$\lambda x : \sigma. e \mid e_1 e_2$	Term abstraction and application
$K \mid \mathbf{case} \overrightarrow{e_0} \mathbf{of} \overrightarrow{p_i \rightarrow e_i}$	Constructors and case matching
$\mathbf{let} \overrightarrow{x : \tau} = \overrightarrow{e} \mathbf{in} u$	Local variable binding
$e \triangleright \gamma$	Casts
$[\gamma]$	Coercions as expressions
$p := K \overrightarrow{x : \tau}$	Patterns
$\tau := a \mid \forall a : \kappa. \tau \mid \tau_1 \tau_2 \mid \dots$	Types
$\kappa := \star \mid \# \mid \kappa \rightarrow \kappa$	Kinds
$\gamma := \mathbf{sym} \gamma$	Symmetry rule for coercions
$\mathbf{nth} 1 \gamma$	Arg part of function coercion
$\mathbf{nth} 2 \gamma$	Result part of function coercion
$\gamma @ \tau$	Type application for coercions
\dots	

Figure 1: Syntax of System F_C (Excerpt)

BETA	$(\lambda x : \tau. e_1) e_2$	$\rightsquigarrow e_1 [e_2/x]$
TYBETA	$(\Lambda a : \kappa. e) \tau$	$\rightsquigarrow e [\tau/a]$
CASEBETA	$\mathbf{case} K \overrightarrow{e_i} \mathbf{of} \dots K \overrightarrow{x_i : \tau_i} \rightarrow e_j \dots$	$\rightsquigarrow e_j [\overrightarrow{e_i/x_i}]$
PUSH	$(e_1 \triangleright \gamma) e_2$	$\rightsquigarrow (e_1 (e_2 \triangleright \mathbf{sym} (\mathbf{nth} 1 \gamma)))$ $\qquad \qquad \qquad \triangleright (\mathbf{nth} 2 \gamma)$
TYPUSH	$(e \triangleright \gamma) \tau$	$\rightsquigarrow (e \tau) \triangleright (\gamma @ \tau)$

Figure 2: Reductions of System F_C (Excerpt)

262 4. A more principled attempt

263 The transformation in Section 3 is achieved by a simple combination of manu-
 264 ally selected inlining, memoization, simplification, and symbolic evaluation. In
 265 order to automate it, we must be precise about what we choose to inline, memo-
 266 ize, and evaluate. For a general-purpose optimization, designing such a heuristic
 267 is hard. However, because we are optimizing a particular style of code, namely
 268 SYB-style code, we can take advantage of domain-specific knowledge.

269 We express these transformations in terms of System F_C (Vytiniotis et al.,
 270 2012), the formal language corresponding to GHC’s `Core` language. Figure 1
 271 presents the relevant parts of the syntax of System F_C , and Figure 2 presents
 272 some of the core reduction rules of System F_C . For simplicity of presentation
 273 these figures omit aspects of System F_C that are not relevant to the optimization
 274 considered in this paper. In particular, System F_C contains additional types and
 275 coercions not listed in Figure 1, as well as additional reductions and machinery
 276 for specifying the evaluation contexts for the reduction rules in Figure 2. The
 277 judgments used by our optimization are listed in Figure 3 and defined in the
 278 following figures.

279 At a high level, the complete optimization can be summarized as follows.
 280 The details and rationale of the individual steps are explained in the remainder
 281 of this section.

282 **Algorithm 1.** *[SYB Optimization] Repeatedly loop until none of the following*
 283 *rules apply. On each loop choose the first rule that applies.*

- 284 1. *Replace any expression with a memoization that it matches as discussed*
 285 *in Section 4.2.*
- 286 2. *Simplify any expression using the rules from Figure 8 as discussed in Sec-*
 287 *tion 4.3.*
- 288 3. *Evaluate any primitive call using the rules from Figure 10 as discussed in*
 289 *Section 4.4.*
- 290 4. *(OPTIONAL) Eliminate any case expression over a manifest constructor*
 291 *as discussed in Section 4.5.1.*
- 292 5. *(OPTIONAL) Float memoization bindings if possible as discussed in Sec-*
 293 *tion 4.5.2.*
- 294 6. *Choose the outermost expression at which we can do either of the following*
 295 *as discussed in Section 4.1.*
 - 296 (a) *Memoize an expression having an undesirable type using the rules*
 297 *from Figure 6.*
 - 298 (b) *Eliminate an expression having an undesirable type using the rules*
 299 *from Figure 4.*

300 Note that our optimization relies on later optimizations already in GHC to fur-
 301 ther clean up the resulting code after our optimization completes. For example,
 302 it may leave behind unused memoization bindings that downstream optimiza-
 303 tions will eliminate. In addition, steps 4 and 5 of this algorithm are optional in

$e : \tau$	Expression typing
$\gamma : \tau_1 \sim \tau_2$	Coercion typing
$e \rightsquigarrow e'$	System F_C evaluation step (See Figure 2)
$e \mapsto e'$	Optimization step (See Figures 4, 6 and 8)
$e \xrightarrow{\gamma} e'$	Symmetric cast elimination (See Figure 9)
$e \rightsquigarrow e'$	Force step (See Figures 7 and 10)
$e \rightsquigarrow e'$	Deep force step (See Figure 10)
Und τ	Undesirable type
ElimUnd e	Elimination expression (See Figure 4)
Memo e	Memoizable expression (See Figure 6)

Figure 3: Judgments

304 that they reduce the work that the optimization has to do but are not essential
 305 for eliminating expressions that have undesirable types.

306 With the benchmarks in Section 5 we show that this algorithm successfully
 307 optimizes typical SYB-style code to be as fast as handwritten code. Remarkably,
 308 this optimization algorithm requires no changes to the standard SYB library
 309 other than what is necessary to ensure inlining information is available for the
 310 appropriate methods, operators, and traversals defined by SYB.

311 4.1. Elimination of expressions with undesirable types

312 In Section 2.4, we identified the presence of expressions with certain types
 313 as a source of performance problems in SYB-style code. However, the transfor-
 314 mations performed in Section 3 allowed us to eliminate expressions with those
 315 types from the code for `incrementSYB`. One of the primary goals of our opti-
 316 mization then is eliminating these occurrences. In particular, objects of type
 317 `TypeRep`, as well as the `TyCon` objects used to construct them, slow down the
 318 code when they are used by `mkT` and similar functions. In addition, the `Data`
 319 and `Typeable` dictionaries contain functions that may generate and manipulate
 320 `TypeRep` and `TyCon` objects. Finally, the default implementations of several of
 321 the methods in the `Data` class use `newtype` wrappers such as `ID` that interfere
 322 with the optimization process and should also be eliminated.

323 In Section 3, we were able to eliminate expressions that have these unde-
 324 sirable types by a combination of inlining and simplification. Moreover, the
 325 only inlining operations necessary were ones that eliminated such expressions.
 326 For example, we inlined calls to `everywhere`, `mkT`, and `gmapT`, which all take
 327 `Data` dictionaries as argument. We also inlined and simplified the call to `cast`

$$\begin{array}{c}
\frac{\mathbf{ElimUnd} \ e \quad e \rightsquigarrow e'}{e \rightsquigarrow e'} \text{ELIMUND} \\
\\
\frac{e_1 : \tau_1 \rightarrow \tau_2 \quad \mathbf{Und} \ \tau_1}{\mathbf{ElimUnd} \ (e_1 \ e_2)} \text{ELIMUNDAPP} \\
\\
\frac{e_0 : \tau \quad \mathbf{Und} \ \tau}{\mathbf{ElimUnd} \ (\mathbf{case} \ e_0 \ \mathbf{of} \ p \rightarrow e_i)} \text{ELIMUNDCASE} \\
\\
\frac{e : \tau \quad \mathbf{Und} \ \tau}{\mathbf{ElimUnd} \ (e \triangleright \gamma)} \text{ELIMUNDCAST}
\end{array}$$

Figure 4: Undesirably Typed Expression Elimination

$$\begin{array}{c}
\tau \in \left\{ \begin{array}{l} \mathbf{Data}, \mathbf{Typeable}, \mathbf{TypeRep}, \mathbf{TyCon}, \\ \mathbf{Fingerprint}, \mathbf{ID}, \mathbf{CONST}, \mathbf{Qi}, \mathbf{Qr}, \mathbf{Mp} \end{array} \right\} \\
\hline
\mathbf{Und} \ \tau \\
\\
\frac{\mathbf{Und} \ \tau_1}{\mathbf{Und} \ (\tau_1 \ \tau_2)} \quad \frac{\mathbf{Und} \ \tau_2}{\mathbf{Und} \ (\tau_1 \ \tau_2)} \quad \frac{\mathbf{Und} \ \tau_1}{\mathbf{Und} \ (\tau_1 \rightarrow \tau_2)} \\
\\
\frac{\mathbf{Und} \ \tau_2}{\mathbf{Und} \ (\tau_1 \rightarrow \tau_2)} \quad \frac{\mathbf{Und} \ \tau}{\mathbf{Und} \ (\forall x : \kappa. \tau)}
\end{array}$$

Figure 5: Undesirable Types

328 which had `Typeable` dictionaries as arguments. This exposed `TypeRep` objects
329 in the scrutinee of a `case` that we then also symbolically evaluated. Thus we
330 can design a heuristic that focuses on expressions that both have these types
331 and are in elimination positions. Expressions in elimination positions are those
332 that are arguments to function applications, scrutinees of `case` expressions, and
333 the bodies of casts. If we can simplify the expression far enough to be able to
334 apply the `BETA` or the `CASEBETA` rules in Figure 2 or expose nested casts that
335 cancel each other out, we can eliminate those occurrences and thus remove the
336 expressions with undesirable types from our code.

337 Essentially what we need to do is symbolically evaluate these expressions until
338 they are values and then apply the appropriate reduction rules to the elimination
339 forms. Formally this is specified by the `ELIMUND` rule in Figure 4. If e is an
340 elimination form for an expression with an undesirable type and we can symbolically
341 evaluate e to e' , then the optimization simplifies e to e' . The elimination
342 forms are specified in the `ELIMUNDAPP`, `ELIMUNDCASE`, and `ELIMUNDCAST`
343 rules, and the rules for forcing a step of evaluation are specified in Figure 7.
344 These rules use the `Und` τ judgment, which holds if and only if the type τ
345 syntactically contains an occurrence of an undesirable type and is defined in

Figure 5. In addition, we use typing judgments for expressions, $e : \tau$, and coercions, $\gamma : \tau_1 \sim \tau_2$. These judgments respectively assert that expression e has type τ and that the coercion γ casts type τ_1 to type τ_2 . The inference rules for these typing judgments are omitted as they are standard in System F_C . In these and other rules, we elide details about the environment as it is not relevant to the optimization other than to support the typing judgments.

Finally, Figure 7 gives the `FORCEBETA`, `FORCETYBETA`, `FORCECASEBETA`, `FORCEPUSH`, and `FORCETYPUSH` rules, which implement symbolic evaluation for the `BETA`, `TYBETA`, `CASEBETA`, `PUSH`, and `TYPUSH` reduction rules respectively. The `FORCEBETA`, `FORCETYBETA`, and `FORCECASEBETA` rules avoid code duplication by introducing `let` bindings instead of substituting. It is then up to `FORCEVAR` to inline forced variables at their use sites. In order to ensure that the `let` forms in the code do not interfere with the optimization process, we also introduce the rules `FORCELETFLOATAPP` and `FORCELETFLOATSCR` which float `let` bindings out of the way so that other rules can fire. The rules `FORCEAPPFUN`, `FORCEAPPTYFUN`, `FORCESCR`, `FORCELETBODY`, and `FORCECAST` implement structural congruences that allow the forcing process to recur down the expression. The guiding principle in all these rules is to make the smallest transformation necessary to expose an expression form that can be eliminated.

4.2. Memoization

In Section 3, we needed to recognize the repeated occurrence of `everywhere` (`mkT inc`) and replace it with a variable reference bound to an equivalent expression. Essentially this is a memoization of the inlining process. Without such memoization, the recursive structure of `everywhere` makes the optimization diverge.

In the example in Section 3, we already have a binding for such an expression, namely `incrementsgyb`. In general, however, we cannot rely on such a binding already being in scope. The original call to `everywhere` might be embedded in another expression, or `everywhere` might be called over a non-uniform type for abstract syntax that contains mutually recursive types for expressions and statements. Even if there is a binding for the type at which `everywhere` is originally called, we need bindings for the other types. Since we cannot rely on the existence of these bindings, we introduce them when we first start simplifying an expression for which we might later need a binding.

Rather than performing a deep analysis of what inlinings and expansions should be memoized, we adopt the very simple strategy of memoizing when the expression e in `ELIMUND` is the application of a variable to one or more arguments. Thus we have `MEMOUND` in Figure 6. This rule has higher priority than `ELIMUND` and should be used instead of that rule whenever possible. In Section 3, the applications of `everywhere`, `mkT`, `gmapT`, and `cast` would all be memoized under this rule. This strategy may lead to unnecessary extra memoization bindings. However, this heuristic is easy to implement, and the extra bindings do not get in the way of the rest of the optimization.

$$\begin{array}{c}
\frac{\mathbf{ElimUnd} \ e \quad e \rightsquigarrow e' \quad \mathbf{Memo} \ e \quad x \notin fv(e')}{e \mapsto \mathbf{let} \ x : \tau = e' \ \mathbf{in} \ x} \text{MEMOUND} \\
\\
\frac{\mathbf{Memo} \ e_1}{\mathbf{Memo} \ (e_1 \ e_2)} \text{MEMOUNDAPP} \quad \frac{\mathbf{Memo} \ e_1}{\mathbf{Memo} \ (e_1 \ \tau)} \text{MEMOUNDTYAPP} \\
\\
\frac{}{\mathbf{Memo} \ x} \text{MEMOUNDVAR}
\end{array}$$

Figure 6: Undesirably Typed Expression Memoization

390 Note that we memoize inlinings only when they eliminate an expression with
391 an undesirable type. The reason for this is that we want to memoize only code
392 that would have triggered `ELIMUND` and not necessarily every intermediate
393 expression.

394 When `MEMOUND` fires we also add e to a memoization table and if e ever
395 occurs again, we replace it with x . We detect reoccurrences only when an
396 expression is manifestly equal to e as we use a simple, syntactic comparison
397 modulo alpha equivalence. For example if e is the expression `mkT f`, then we
398 do not consider `mkT f'` to be a reoccurrence of e even if f' is bound to f .
399 While in theory the optimization could as a result miss opportunities to take
400 advantage of the memoization, in practice there are only a few ways that this
401 happens in SYB-style code, and they are automatically eliminated by the other
402 simplifications in the optimization.

403 4.3. Simplification

404 As we symbolically evaluate the code, detritus can build up in the form of
405 dead and trivial `let` bindings and unnecessary casts. Though in some cases we
406 can leave the elimination of these for later optimization passes in the compiler,
407 some of these `let` bindings and casts get in the way of the core optimization
408 rules from Figure 4 and Figure 6. In the example in Section 3, many of the
409 intermediate simplifications were omitted in order to focus on the core aspects
410 of the optimization, but now we formally specify these by applying the simplifi-
411 cations from Figure 8 to the code as we are optimizing it. These simplifications
412 are chosen based on an empirical observation of the sort of code generated when
413 optimizing SYB-style code and what forms need to be simplified in that process.
414 While there are a number of other simplifications that could be used, we restrict
415 ourselves to a minimal number of conservative simplifications that never make
416 the code worse while still being sufficient to enable the core optimization rules.

417 4.3.1. Cast elimination

418 GHC's implementation of `newtype` definitions and some class dictionaries
419 makes use of casts. For example, a call to the `typeOf` method of the `Typeable`
420 class is implemented as a cast. In addition, many of the SYB functions use

FORCEBETA		
	$(\lambda x : \tau. e_1) e_2$	$\rightsquigarrow \mathbf{let } x : \tau = e_2 \mathbf{ in } e_1$
FORCETYBETA		
	$(\Lambda a : \kappa. e) \tau$	$\rightsquigarrow \mathbf{let } a : \kappa = \tau \mathbf{ in } e$
FORCECASEBETA		
	$\mathbf{case } K \overrightarrow{e_i} \mathbf{ of } \dots K \overrightarrow{x_i : \tau_i} \rightarrow e_j \dots$	$\rightsquigarrow \mathbf{let } \overrightarrow{x_i : \tau_i} = \overrightarrow{e_i} \mathbf{ in } e_j$
FORCEPUSH		
	$(e_1 \triangleright \gamma) e_2$	$\rightsquigarrow (e_1 (e_2 \triangleright \mathbf{sym} (\mathbf{nth } 1 \gamma))) \triangleright (\mathbf{nth } 2 \gamma)$
FORCETYPUSH		
	$(e \triangleright \gamma) \tau$	$\rightsquigarrow (e \tau) \triangleright (\gamma @ \tau)$
FORCEVAR		
	x	$\rightsquigarrow e$ if e is the inlining of x
FORCELETFLOATAPP		
	$(\mathbf{let } \overrightarrow{x : \tau} = \overrightarrow{e_i} \mathbf{ in } e_0) u$	$\rightsquigarrow \mathbf{let } \overrightarrow{x : \tau} = \overrightarrow{e_i} \mathbf{ in } e_0 u$
FORCELETFLOATSCR		
	$\mathbf{case } (\mathbf{let } \overrightarrow{x : \tau} = \overrightarrow{u} \mathbf{ in } e_0) \mathbf{ of } \overrightarrow{p_i} \rightarrow \overrightarrow{e_i}$	$\rightsquigarrow \mathbf{let } \overrightarrow{x : \tau} = \overrightarrow{u} \mathbf{ in } (\mathbf{case } e_0 \mathbf{ of } \overrightarrow{p_i} \rightarrow \overrightarrow{e_i})$
FORCEAPPFUN		
	$e_1 e_2$	$\rightsquigarrow e'_1 e_2$ if $e_1 \rightsquigarrow e'_1$
FORCEAPPTYFUN		
	$e_1 \tau$	$\rightsquigarrow e'_1 \tau$ if $e_1 \rightsquigarrow e'_1$
FORCESCR		
	$\mathbf{case } e_0 \mathbf{ of } \overrightarrow{p_i} \rightarrow \overrightarrow{e_i}$	$\rightsquigarrow \mathbf{case } e'_0 \mathbf{ of } \overrightarrow{p_i} \rightarrow \overrightarrow{e_i}$ if $e_0 \rightsquigarrow e'_0$
FORCELETBODY		
	$\mathbf{let } \overrightarrow{x_i : \tau_i} = \overrightarrow{u_i} \mathbf{ in } e$	$\rightsquigarrow \mathbf{let } \overrightarrow{x_i : \tau_i} = \overrightarrow{u_i} \mathbf{ in } e'$ if $e_0 \rightsquigarrow e'_0$
FORCECAST		
	$e \triangleright \gamma$	$\rightsquigarrow e' \triangleright \gamma$ if $e \rightsquigarrow e'$

Figure 7: Forcing Rules

CASTREFL	$e \triangleright \gamma$	$\mapsto e$ if $\gamma : \tau \sim \tau$
CASTSYM	$e \triangleright \gamma$	$\mapsto e'$ if $e \xrightarrow{\gamma} e'$
DEADLET	$\mathbf{let} x : \tau = u \mathbf{ in} e$	$\mapsto e$ if $x \notin fv(e)$ and x is not a memoization
SUBSTSTAR	$\mathbf{let} x : \star = \tau \mathbf{ in} e$	$\mapsto e[\tau / x]$
SUBSTHASH	$\mathbf{let} x : \# = \tau \mathbf{ in} e$	$\mapsto e[\tau / x]$
SUBSTVAR	$\mathbf{let} x : \tau = x' \mathbf{ in} e$	$\mapsto e[x' / x]$
SUBSTLIT	$\mathbf{let} x : \tau = l \mathbf{ in} e$	$\mapsto e[l / x]$
SUBSTDFUN	$\mathbf{let} x : \tau = v \vec{u} \mathbf{ in} e$	$\mapsto e[v \vec{u} / x]$ if v is a dictionary constructor

Figure 8: Simplifications

421 **newtype** definitions to define higher level operations in terms of lower level
 422 operations. For example, the default implementation of `gmapT` is in terms of
 423 `gfoldl` with `ID` for the `c` type parameter. Since `ID` is a **newtype**, calls to the `ID`
 424 constructor and `unID` destructor get translated to casts.

425 As these higher-level operations project into and out of these types, these
 426 constructors and destructors may be directly or indirectly nested on each other
 427 as pairs of symmetric of casts that could be eliminated. In addition, the types
 428 in these casts may be refined by `FORCE TYBETA` or shuffled around by the
 429 `FORCE PUSH` or `FORCE TY PUSH` rules to result in casts that are reflexive.

430 These casts can quickly build up and get in the way of the core optimization
 431 rules. For example, it often happens that the scrutinee of a `case` contains a
 432 reflexive cast wrapped around a constructor. Until we eliminate the cast, we
 433 cannot use the `FORCE CASE BETA` rule even though the constructor involved is
 434 already manifest.

435 Reflexive casts from a type to itself are directly eliminated with the `CASTREFL`
 436 rule, which just checks the type of the cast. Symmetric casts, however, could
 437 be separated from each other by intermediate forms as in the following example
 438 where $\gamma_1 : \tau_1 \sim \tau_2$, $\gamma_2 : \tau_2 \sim \tau_1$, and $\gamma_3 : \tau_2 \sim \tau_1$.

$$(\mathbf{case} x \mathbf{ of} \{C_1 \rightarrow e_1 \triangleright \gamma_2; C_2 \rightarrow e_2 \triangleright \gamma_3\}) \triangleright \gamma_1$$

439 Simplifying this expression is accomplished by the `CASTSYM` rule. This rule uses
 440 the $e \xrightarrow{\gamma} e'$ judgment in Figure 9 to check whether all paths down e contain a cast
 441 symmetric to γ . That judgment returns the expression with those symmetric
 442 casts removed as e' , and thus `CASTSYM` reduces our example to to the following.

$$\mathbf{case} x \mathbf{ of} \{C_1 \rightarrow e_1; C_2 \rightarrow e_2\}$$

$$\begin{array}{c}
\frac{\gamma : \tau \sim \tau' \quad \gamma' : \tau' \sim \tau}{e \triangleright \gamma' \xrightarrow{\gamma} e} \text{CASTSYMCAST} \\
\\
\frac{\gamma : (\tau_1 \rightarrow \tau_2) \sim (\tau_1 \rightarrow \tau_2') \quad e \xrightarrow{\text{nth } 2 \gamma} e'}{\lambda x : \tau. e \xrightarrow{\gamma} \lambda x : \tau. e'} \text{CASTSYMFUN} \\
\\
\frac{e \xrightarrow{\gamma} e'}{\text{let } \bar{x} : \tau = e_i \text{ in } e \xrightarrow{\gamma} \text{let } \bar{x} : \tau = e_i \text{ in } e'} \text{CASTSYMLET} \\
\\
\frac{\begin{array}{c} \xrightarrow{\gamma} \\ e_i \xrightarrow{\gamma} e_i' \end{array}}{\text{case } e \text{ of } \bar{p} \rightarrow e_i \xrightarrow{\gamma} \text{case } e \text{ of } \bar{p} \rightarrow e_i'} \text{CASTSYMCASE}
\end{array}$$

Figure 9: Cast Symmetry Rules

443 4.3.2. Let elimination

444 We also eliminate `let` bindings that are either trivial, dead, or bind a type as
445 they may interfere with our ability to apply the core optimization rules. These
446 are implemented by the remaining rules in Figure 8. Note that when doing
447 this we are careful to not eliminate bindings introduced by memoization. In
448 particular, due to the way that GHC implements class dictionaries, it is quite
449 common for a memoized call to expand to another memoized call in a way that
450 results in the memoized binding for the original call becoming trivial. We must
451 avoid eliminating these as the memoization process may add new references to
452 such bindings.

453 4.4. Primitives

454 Recall that the `cast` function is implemented by testing the equality of two
455 `TypeRep` objects returned by calls to `typeOf`. This `typeOf` operator is imple-
456 mented in terms of `fingerprintFingerprints`, which computes unique hashes
457 for `TypeRep` objects. Furthermore, equality over these objects is implemented
458 in terms of the `eqWord#` and `tagToEnum#` primitives. As we are attempting to
459 eliminate the dynamic dispatches implemented by `cast`, it is important that we
460 eliminate calls to these primitives. In order to do so, our optimization fully eval-
461 uates the arguments to these functions when attempting to force an expression.
462 Once those arguments are fully evaluated, the calls themselves are statically
463 evaluated. The rules that implement this are specified in Figure 10 where we
464 use double brackets (`[[` and `]]`) for compile time evaluation. These rules effectively
465 implement constant folding for these operators.

PRIMFF		
	<code>fingerprintFingerprints e</code>	$\rightsquigarrow \llbracket \text{fingerprintFingerprints } e \rrbracket$ if e is a value
PRIMFFARG		
	<code>fingerprintFingerprints e</code>	$\rightsquigarrow \text{fingerprintFingerprints } e'$ if $e \rightsquigarrow e'$
PRIMEQWORD		
	<code>eqWord# e₁ e₂</code>	$\rightsquigarrow \llbracket \text{eqWord# } e_1 e_2 \rrbracket$ if e_1 and e_2 are values
PRIMEQWORDARG1		
	<code>eqWord# e₁ e₂</code>	$\rightsquigarrow \text{eqWord# } e'_1 e_2$ if $e_1 \rightsquigarrow e'_1$
PRIMEQWORDARG2		
	<code>eqWord# e₁ e₂</code>	$\rightsquigarrow \text{eqWord# } e_1 e'_2$ if $e_2 \rightsquigarrow e'_2$
TAGTOENUM		
	<code>tagToEnum# e₁</code>	$\rightsquigarrow \llbracket \text{tagToEnum# } e_1 \rrbracket$ if e is a value
TAGTOENUMARG		
	<code>tagToEnum# e₁</code>	$\rightsquigarrow \text{tagToEnum# } e'_1$ if $e \rightsquigarrow e'$
FORCEDEEP		
	<code>e</code>	$\rightsquigarrow e'$ if $e \rightsquigarrow e'$
FORCEDEEPARG		
	<code>e₁ e₂</code>	$\rightsquigarrow e_1 e'_2$ if $e_2 \rightsquigarrow e'_2$

Figure 10: Rules for Primitives

466 *4.5. Optional optimizations*

467 While not essential to the core optimization and the elimination of expres-
 468 sions with undesirable types, there are certain transformations that help keep
 469 the generated code compact and reduce the amount of work to be done by the
 470 optimization.

471 *4.5.1. Case reduction*

472 SYB-style traversals are based on the idea of dispatching to different code
 473 depending on the current type being traversed. At its core, this is the purpose of
 474 `mkT`. When optimizing SYB-style code, this often results in intermediate residual
 475 code with a structure similar to the following.

```
476     case typeOf t1 == typeOf t2 of
477       True  -> ...
478       False -> ...
```

479 The equality operator in this code is over the undesirable type `TypeRep`, so the
 480 optimization will reduce it to either `True` or `False`. After that, the scrutinee no
 481 longer contains an expression with an undesirable type, so the core optimiza-
 482 tion does not then simplify the `case` expression even though it has a known
 483 constructor in its scrutinee. In most cases this is not a problem as the code to
 484 be optimized under each branch of the `case` expression tends to be small, and
 485 we can simply rely on downstream optimizations to simplify the `case` expres-
 486 sion. However, when these branches are large, they can represent a significant
 487 amount of extra work to be done by the optimization. It would be better to
 488 detect the dead branch and skip the extra work in that branch. To do this we
 489 apply the rewrite in `FORCECASEBETA` whenever possible. This rewrite never
 490 makes the code worse or worsens the optimization result. Note that our use of
 491 this rewrite differs from the usual use of the rules in Figure 7 since we apply it at
 492 any position in the expression regardless of whether it eliminates an expression
 493 with an undesirable type.

494 *4.5.2. Memoization floating*

495 Duplicate memoizations of the same expression may arise if the first memo-
 496 ization is not in scope at the other occurrences of the same expression. For
 497 example, when traversing an abstract syntax tree, memoizations of the traversal
 498 at the identifier type may occur inside both the part of the code for λ -expressions
 499 and the part of the code for `let` expressions. If neither of these is within the
 500 scope of the other, the memoization rule will result in creating fresh memoiza-
 501 tions of the traversal on identifiers for each expression form even though the
 502 code for these memoizations are identical to each other.

503 As a consequence of this, it is relatively easy to get code that is exponen-
 504 tially large in the size of the types being traversed because the inlining process
 505 may not terminate until every path down the expanded expression contains a
 506 memoization for every type being traversed. Even in cases when the code does

507 not blow up to be exponentially large, these duplicated memoizations represent
 508 extra work for the optimizer and inflate the size of the resulting code.

509 To avoid this size explosion, we `let`-float memoized bindings as far outward
 510 as possible. By floating the memoized bindings outwards, we maximize their
 511 scope and thus avoid creating duplicate memoizations due to already created
 512 memoizations being out of scope. For example, once the memoization created
 513 for the identifier in a λ -expression floats outwards, the traversal for the identifier
 514 in a `let` expression can use the existing memoization instead of creating a new
 515 one. We also consolidate memoization bindings into a common recursive `let`
 516 binding when possible as, while they may not initially refer to each other, the
 517 process of replacing expressions with their memoized bindings may make them
 518 refer to each other at some later point.

519 5. Implementation

520 We implemented the custom optimization pass described in Section 4 using
 521 HERMIT, a recently developed GHC plugin for applying transformations to
 522 `Core` (Farmer et al., 2012; Sculthorpe et al., 2013). HERMIT was used inter-
 523 actively to gain an intuition about the transformations necessary and was then
 524 extended with new primitive transformations implementing the rules given in
 525 Section 4. The overall optimization in Algorithm 1 was implemented as a HER-
 526 MIT plugin. After the optimization completes, we use HERMIT’s `simplify`
 527 command to perform basic simplification like dead `let`-binding elimination.

528 HERMIT provides several facilities to ease the implementation of `Core`-to-
 529 `Core` transformations such as our optimization. This includes KURE, a strategic
 530 rewriting library allowing transformations to be expressed in a high-level, declar-
 531 ative style (Gill, 2009; Farmer et al., 2012; Sculthorpe et al., 2014), a versioning
 532 kernel which manages the application of rewrites, congruence combinators for
 533 `Core` which automatically update the rewriting context, error reporting facili-
 534 ties, and a large set of existing primitive rewrites and queries. Not including
 535 primitive transformations already available in HERMIT, the entire optimization
 536 was implemented in approximately 450 lines of Haskell and did not require any
 537 modifications to GHC itself.

538 5.1. Benchmarks

539 We applied the optimization to a selection of benchmarks taken from the
 540 Haskell generic-programming literature. The resulting programs were bench-
 541 marked using a version of the framework from Magalhães et al. (2010) that
 542 was adapted to support compilation with HERMIT. The benchmarks were as
 543 follows.

544 **RmWeights** Taken from `GPBench` (Rodriguez et al., 2008), the `RmWeights`
 545 benchmark traverses a weighted binary tree while removing the weight
 546 annotations. It is implemented in SYB using the `everywhere` and `mKT`
 547 combinators.

548 **SelectInt** Also from `GPBench`, `SelectInt` traverses a weighted binary tree
 549 while collecting all the `Ints` into a single list. It is naively implemented
 550 in SYB using the `everything` and `mkQ` combinators, but as we discuss in
 551 Section 5.2, it had to be modified to ensure a fair comparison.

552 **Map** Found in Magalhães et al. (2010), `Map` performs a mapping over a struc-
 553 ture. It is implemented in SYB using `everywhere` and `mkT`. This traversal
 554 is performed on two¹ data types. The first is a binary tree of integers. The
 555 second is a logic formula. For the binary tree, all integers are incremented.
 556 For the other type, all characters are replaced with the character ‘y’.

557 **RenumberInt** Taken from Adams and DuBuisson (2012), the `RenumberInt`
 558 benchmark replaces each integer in a structure with a new, unique integer
 559 that is drawn from a state monad. This traversal is also performed on
 560 both binary tree and logic formula data types. It is implemented in SYB
 561 using `everywhereM` and `mkM`.

562 5.2. Benchmark setup

563 Each benchmark was implemented both non-generically (`Hand`) and using
 564 SYB combinators (`SYB`). The SYB implementation was also benchmarked with
 565 our optimization (`SYB/Hermit`). The benchmarking framework used in Maga-
 566 lhães et al. (2010) was used to run each program 10 times and take the average
 567 running time. We compiled the benchmarks with GHC 7.8.2 using the `-O2`
 568 compiler option and ran them with the `-K1g` RTS option on a 1.7 GHz, 64-bit
 569 Intel i7 with 8 GB of RAM running Darwin 13.1.0.

570 The implementation of `SelectInt` in `GPBench` uses two different algorithms
 571 for the `Hand` and `SYB` implementations. The `Hand` implementation uses a
 572 linear-time, accumulating-style traversal, while the `SYB` implementation uses a
 573 quadratic-time, non-accumulating traversal. To ensure a fair comparison, we
 574 modified the `SYB` implementation to use an accumulating traversal. Similiary,
 575 the `Hand` implementation of `RenumberInt` in `GPBench` did not descend into
 576 strings, whereas the `SYB` implementation does. We modified the `Hand` imple-
 577 mentation to match the `SYB` traversal.

578 5.3. Performance results

579 Figure 11 summarizes the resulting execution times of the benchmarks. The
 580 results are normalized relative to the `Hand` version and are displayed on a log-
 581 arithmic scale in order to accommodate the large differences between execution
 582 times. These benchmarks confirm previous results about the poor performance
 583 of `SYB` as it performed on average an order of magnitude slower than the hand-
 584 written code.

¹Note to reviewers: Due to a combination of lost files and newly discovered bugs in GHC and/or HERMIT, we are currently unable to present the `MapSrc` benchmark results that were in the previous version of this paper. We are working on these problems and hope to add these results back into the paper once these issues are solved.

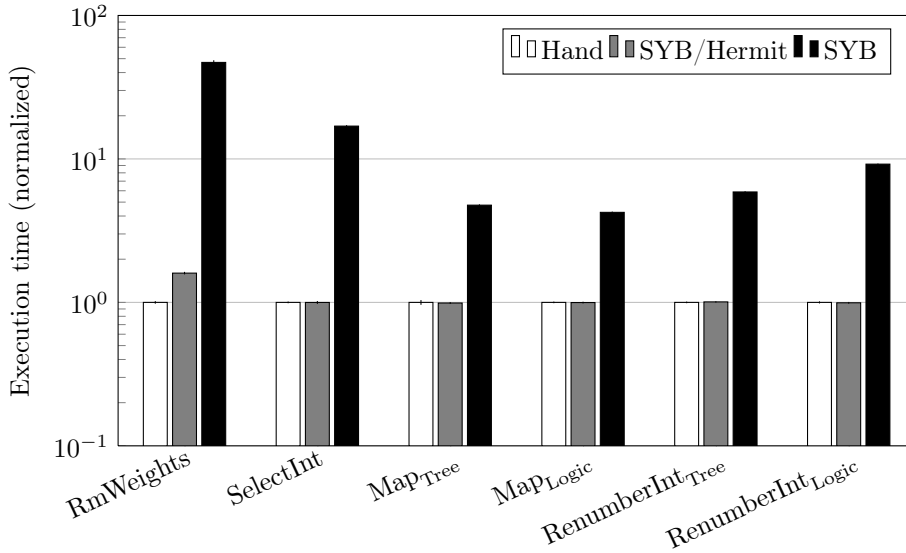


Figure 11: Benchmarks Results

585 For all of the benchmarks except `RmWeights`, the optimization completely
 586 eliminates the runtime costs associated with SYB. A manual inspection of the
 587 generated `Core` confirms that the optimization does indeed eliminate all run-
 588 time type checks and dictionary dispatches in the SYB-style code and that the
 589 resulting code is equivalent to the handwritten code.

590 When initially running these benchmarks, the SYB/Hermit versions of `MapTree`
 591 and `MapLogic` actually ran *faster* than the Hand versions by about 20%. Analy-
 592 sis of the resulting `Core` revealed that, as a side effect of our optimization, the
 593 traversal was being specialized to the particular function being mapped over
 594 the structure. The Hand version did not do this. Rewriting the Hand ver-
 595 sion by applying a static-argument transformation (Santos, 1995) improved its
 596 performance to match that of the SYB/Hermit version.

597 On the `RmWeights` benchmark, SYB/Hermit fails to achieve parity with the
 598 Hand version. This contrasts with a previous version of this paper (Adams
 599 et al., 2014) in which `RmWeights` fully optimized. Inspecting the `Core` reveals
 600 why. The optimization successfully eliminates all runtime type checks and dic-
 601 tionary dispatches as expected. After several of GHC’s own optimization passes
 602 run, including two rounds of the simplifier, we are left with the following two
 603 mutually-recursive functions. (Casts have been omitted for clarity.)

```

memo_everywhere :: WTree Int Int → WTree Int Int
memo_everywhere = λ x →
  case memo_gfoldl x of
    WithWeight t w → t
    wild → wild

```

```

memo_gfoldl :: WTree Int Int → WTree Int Int
memo_gfoldl = λ ds →
  case ds of
    Leaf a1 → Leaf Int Int a1
    Fork a1 a2 →
      Fork Int Int (memo_everywhere a1) (memo_everywhere a2)
    WithWeight a1 a2 →
      WithWeight Int Int (memo_everywhere a1) a2

```

604 In order to achieve the same performance as Hand, the `memo_gfoldl` function
605 needs to be inlined into `memo_everywhere`. This causes a subsequent `case-of-`
606 `case` transformation and `case-reduction` by the simplifier, resulting in a single
607 self-recursive traversal function. When we tested this by forcing the inlining with
608 HERMIT, we observe the desired speedup. However, GHC marks `memo_gfoldl`
609 as a loop breaker, an annotation it uses to ensure that inlinings in mutually
610 recursive binding groups terminate. This prevents the full optimization of this
611 code. We speculate that `RmWeights` fully optimized in the previous version of
612 this paper because a different function was chosen by GHC as the loop breaker.
613 However, we have no way to test this as we no longer have the particular devel-
614 opment version of GHC used in that paper.

615 In the previous version of this paper, `ReNumberIntLogic` performed 2.2 times
616 slower than the Hand version. Subsequent investigation has revealed the cause.
617 Part of this slowdown was due to the selective traversal issue we mention in
618 Section 5.2, namely, the SYB version descended into strings while the Hand
619 version did not. Fixing this brought SYB/Hermit to 1.6 times slower than
620 Hand. Investigating the resulting Core showed that the remaining slowdown
621 was caused by poor interactions with GHC’s unboxing optimizations.

622 Recall that the `ReNumberInt` benchmark uses a `State` monad to generate
623 fresh integers during the traversal. The `State` monad in Haskell is implemented
624 using a function which returns a tuple of value and state. Combining two `State`
625 computations with `>>=` results in the allocation of a tuple for the result of the
626 first computation followed by a `case` expression to extract the value and state
627 from the tuple for use by the second computation. This intermediate allocation
628 of tuples is wasteful so, when possible, GHC’s Constructed Product Result
629 (CPR) analysis pass (Baker-Finch et al., 2004) eliminates tuples by unboxing.

630 The code resulting from our optimization prevents this unboxing. We spec-
631 ulate that residual casts are interfering with the CPR analysis. We can improve
632 the situation by switching to the strict `State` monad, which immediately scru-
633 tinizes the result of the first `State` computation rather than allocating it with a
634 `let` binding. This makes the code resulting from our optimization amenable
635 to CPR, which then successfully unboxes the tuples. Switching to a strict
636 `State` monad for `ReNumberIntTree` and `ReNumberIntLogic` improves the run-
637 ning time of Hand by a factor of 1.1, the unoptimized SYB by a factor of 1.2,
638 and SYB/Hermit by a factor of 1.8 at which point SYB/Hermit matches the
639 performance of Hand. The results for `ReNumberIntTree` and `ReNumberIntLogic`
640 in Figure 11 are for the strict `State` monad.

641 6. Limitations and Future Work

642 While the algorithm described in Section 4 is effective for most instances of
 643 SYB-style code, it does have limitations and areas that future work can improve.
 644 Many of these problems will be familiar to the partial-evaluation community.
 645 As these are active research topics in their own right, we do not attempt a
 646 general solution to them but where possible note how they can be mitigated for
 647 our particular optimization. As it is domain-specific, this optimization may not
 648 be appropriate for all code, and the compiler may require assistance from the
 649 programmer in the form of pragmas or annotations to determine when to use
 650 or not use this optimization.

651 6.1. Missing inlining information

652 The first and most obvious limitation is that this optimization relies heavily
 653 on inlining and thus depends on having the appropriate inlining information
 654 available. If that information is not available, then the optimization may fail to
 655 complete its task of eliminating expressions with undesirable types. Fortunately,
 656 this is an easily detected situation, and the optimization can abort while leaving
 657 the original code intact and issue a warning so the user can make appropriate
 658 adjustments to expose the necessary inlining information.

659 Missing inlining information can be caused by using functions from imported
 660 modules for which GHC has not recorded inlining information. For example, by
 661 default, the inlining information for several operations in SYB were not available
 662 so we had to use `-fexpose-all-unfoldings` or add `INLINABLE` pragmas to
 663 expose these. Missing inlining information may also be caused by running the
 664 optimization over code in which the types over which `Data` or `Typeable` are
 665 quantified are underspecified. For example, consider the following code that
 666 one might write as a helper function.

```
667     mapSYB :: (Data a) => (a -> a) -> [a] -> [a]
668     mapSYB f x = everywhere (mkT f) x
```

669 Since this function is polymorphic in `a`, there is no concrete dictionary available
 670 for the class constraint `Data a`, and we cannot fully optimize this function.

671 There are, however, two important points to consider about this limitation.
 672 First, as it is obviously impossible to specialize a generic traversal when we do
 673 not yet know the type at which to specialize, this limitation is inherent in the
 674 optimization task and not merely a failure of the optimization algorithm. For
 675 example, if `a` is instantiated with `[Char]`, then `f` must be applied not only to the
 676 elements of the list passed to `mapSYB` but also to the sub-lists of those elements.
 677 Until we know `a`, it is impossible to know how to traverse those elements.

678 Second and more importantly, this limitation is not a problem in practice.
 679 It simply means that the optimization must be deferred to uses of the function
 680 that specify types at which to specialize. For example, instead of optimizing
 681 `mapSYB`, we optimize uses of `mapSYB` such as the following.

```

682     incrementSYB/Int :: [Int] -> [Int]
683     incrementSYB/Int = mapSYB inc

```

684 Because this definition completely determines the type of `a` in `mapSYB` and thus
685 calls `mapSYB` with a concrete `Data` dictionary for `a`, the optimization will suc-
686 cessfully complete on `incrementSYB/Int` even though it would fail on `mapSYB`.

687 Finally, note that specialized versions of `mapSYB` that are successfully opti-
688 mized by our optimization can be explicitly generated by specifying their types
689 as in the following.

```

690     mapSYB/Int :: (Int -> Int) -> [Int] -> [Int]
691     mapSYB/Int = mapSYB

```

692 6.2. Essential occurrences of undesirable types

693 Since the primary design heuristic behind this optimization is the elimination
694 of expressions that have undesirable types, it will fail if there are expressions that
695 have undesirable types but should not be eliminated. An obvious example is
696 when the type being traversed itself contains undesirable types such as `TypeRep`
697 or `TyCon`, but less obvious examples of this include types like the following from
698 Hinze et al. (2006).

```

699     data Spine b
700       = Unit b
701       | ∀a. (Data a) => App (Spine (a -> b)) a

```

702 Here the existential² type `a` is qualified by the `Data` class and thus the `App`
703 constructor contains a dictionary for the `Data` class.

704 Along similar lines, it may be possible for a particular traversal to contain
705 essential uses of undesirable types. For example, `SYB` allows code to arbitrarily
706 synthesize `TypeRep` and `TyCon` objects. This may result in occurrences of unde-
707 sirable types that are essential to the traversal and either should not or cannot
708 be eliminated. Note that though such a traversal is possible, it is exceedingly
709 rare in `SYB`-style code. None of the standard traversals exhibit such a structure.

710 This limitation may be mitigated by annotating the code with information
711 about which occurrences of undesirable types are genuinely undesirable and
712 which are not. Then as the optimization transforms the code, we can keep
713 careful account of each occurrence and whether it is genuinely undesirable.

714 6.3. Polymorphic recursion in types

715 As with other forms of partial evaluation, polymorphic recursion is a concern
716 with this optimization. Most types in Haskell programs are regular, but non-
717 regular, polymorphically recursive types do occasionally occur. Consider, for
718 example, the following polymorphically recursive, non-regular type.

²GHC uses the `∀` keyword for both existential and universal types. The distinction between the two is where the keyword is placed.

```

719     data T a
720         = Base a
721         | Double (T (a, a))

```

722 If we attempt to traverse over the type `T Int`, then the traversal will initially
723 be memoized at `T Int`. Since at this type the argument to the `Double` con-
724 structor is of type `T (Int, Int)`, the traversal will also have to be memoized
725 at type `T (Int, Int)`. In turn, at that type, the argument to the `Double` con-
726 structor has type `T ((Int, Int), (Int, Int))` and so on. Naively running
727 the optimization on this type would thus continue forever as the memoization
728 process depends on the assumption that there are a finite number of types to be
729 traversed, but the `T Int` type effectively contains an infinite number of types.

730 In order to successfully handle this, we would need to account for the fact
731 that in many cases a non-generic traversal over a polymorphic type must be
732 structured differently from a generic traversal. In these cases it is impossible
733 to generate non-generic code that naively mirrors the structure of the generic
734 code. For example, consider a traversal that increments all values of type `Int`
735 inside an object of type `T Int`. The generic code for this is the following.

```

736     incrementT :: T Int -> T Int
737     incrementT x = everywhere (mkT inc) x

```

738 Now consider how one would write this with non-generic code. The recursion
739 over the elements of `T` cannot have type `T Int -> T Int` since the `Double`
740 constructor changes the type argument of `T`. On the other hand, the recursion
741 cannot have type $\forall a. T a \rightarrow T a$ since being polymorphic in `a` prevents the
742 function from manipulating the `Int` that occur in `a`. Instead, a more sophisti-
743 cated implementation such as the following is necessary.

```

744     incrementT :: T Int -> T Int
745     incrementT x = go inc x where
746         go :: (a -> a) -> T a -> T a
747         go f (Base x) = f x
748         go f (Double t) = Double (go (f' f) t)
749         f' :: (a -> a) -> (a, a) -> (a, a)
750         f' f (x1, x2) = (f x1, f x2)

```

751 Since the optimization presented in this paper preserves the structure of the
752 generic traversal and `incrementT` does not follow that structure, it is unsur-
753 prising that our optimization fails on such a traversal. However, note that the
754 `f` argument to `go` serves essentially the same role as the `Data` dictionary in the
755 generic traversal in that it provides the necessary information for implementing
756 the parts of the traversal that operate over the type `a`. Thus an interesting
757 direction for future work would be deriving such a non-generic implementation
758 from the generic traversal by appropriately specializing and simplifying the `Data`
759 dictionary.

760 *6.4. Polymorphic recursion in terms*

761 In addition to types being polymorphically recursive, the traversal itself may
 762 be polymorphically recursive in an argument whose type contains undesirable
 763 types. Traversals like this are rare in SYB-style code, but one could imagine an
 764 example like the following.

```
765     poly :: (∀b. Data b => b -> b)
766           -> (∀a. Data a => a -> a)
767     poly f x = f (gmapT (poly (f `extT` g)) x)
768     where g = ...
```

769 Note how the `f` argument to the traversal is extended each time through the
 770 traversal. As a result, the previously memoized instances of `poly` cannot be
 771 used and the optimization algorithm will never be able to completely eliminate
 772 all expressions with undesirable types.

773 Of course, this is a concern only because the type of `f` contains an undesirable
 774 type. Parameters such as `x` that do not have a type containing an undesirable
 775 type can freely vary from call to call as the memoization does not care about
 776 them.

777 As with polymorphically recursive types, this limitation is not unique to
 778 optimizing SYB-style code. Polymorphic recursion is an area of active research
 779 in the partial evaluation community for which we do not have a solution in the
 780 general case.

781 *6.5. Selective traversal*

782 An instance where the optimization does not fail but the results could be
 783 improved is when parts of the generic traversal expand to trivial traversals that
 784 do no useful work. For example, a traversal that modifies only integers can safely
 785 skip over any strings that it finds and avoid processing the individual characters
 786 in the string. Adams and DuBuisson (2012) call this selective traversal and
 787 document the significant performance improvements this can achieve. SYB does
 788 not do selective traversal unless it is explicitly told what expressions to skip. In
 789 the code produced by our optimization, these skippable parts of the traversal
 790 are manifest as functions that do a trivial deconstruction and reconstruction.
 791 For example, in a traversal that effects only integers, we might find code for
 792 traversing strings similar to the following.

```
793     memoChar  c      = c
794     memoString []    = []
795     memoString (c : cs) = memoChar c : memoString cs
```

796 Here `memoString` is equivalent to the identity function and can thus be more
 797 efficiently implemented by not doing the traversal and simply returning its ar-
 798 gument. Depending on the structure of the data being traversed, this can lead
 799 to significant speedups.

800 Similar situations arise for queries and monadic traversals. For queries,
 801 some parts of the traversal may produce trivial query results, and for monadic

802 traversals, some parts of the traversal may be equivalent to simply applying
 803 `return` to the tree being traversed.

804 Identifying and optimizing these trivial functions is fairly easy and can be
 805 done by a post-processing pass after our optimization. We plan to add this in
 806 future versions of our implementation.

807 7. The GHC Specializer

808 Given that the core rules of our optimization specialize functions to partic-
 809 ular arguments, a natural question is whether the existing specializer in GHC
 810 can achieve the same effect. However, the GHC specializer focuses on class-
 811 dictionary specialization (Jones, 1995) and does not specialize non-dictionary
 812 arguments. This is a problem in `incrementSYB` where we need to specialize
 813 `everything` over the non-dictionary argument `mkT inc`. As a consequence, the
 814 default optimization pipeline in GHC does not do the specialization needed to
 815 effectively optimize SYB-style code.

816 The situation is not a total loss, however. Under appropriate conditions
 817 the GHC specializer will specialize some parts of `incrementSYB` over the `[Int]`
 818 type and produce the following code.

```
819     incrementSYB :: [Int] -> [Int]
820     incrementSYB x = everywhere[Int] (mkT inc) x
821
822     everywhere[Int] :: (forall b. Data b => b -> b)
823                    -> [Int] -> [Int]
824     everywhere[Int] f x = f (gmapT[Int] (everything f) x)
825
826     gmapT[Int] :: (forall b. Data b => b -> b)
827                -> [Int] -> ID [Int]
828     gmapT[Int] f [] = []
829     gmapT[Int] f (x : xs) = f x : f xs
```

830 Unfortunately, while `everything` and `gmapT` are specialized to particular types
 831 in this code, the `f` arguments to these functions are not. They are still polymor-
 832 phic and take class dictionaries as arguments. This is because the techniques
 833 used by the GHC specializer do not handle the rank-2 polymorphism of these
 834 arguments. As a consequence, when `incrementSYB` is invoked, the outermost
 835 call to `everything` and `gmapT` use the specialized version, but the inner calls
 836 to `everything` that are made by `gmapT` use the unspecialized version. As a re-
 837 sult, the bulk of the computation runs slowly and does not use these specialized
 838 versions of `everything` and `gmapT`.

839 Even though the default GHC optimization pipeline does not do well on this
 840 code, there are some things we can do to help it. First, we can manually perform
 841 a static argument transformation on `everything` and define it as follows.

```
842     everywhere :: (forall b. Data b => b -> b)
```

```

843         -> (∀a. Data a => a -> a)
844     everywhere f x = go x where
845         go :: ∀c. Data c => c -> c
846         go x = f (gmapT go x)

```

847 With this definition, inlining `everywhere` produces a version of `go` that imple-
848 ments the work of `everywhere` but specialized to one particular value of `f`. For
849 example, inlining `everywhere` into `incrementSYB` results in the following.

```

850     incrementSYB :: [Int] -> [Int]
851     incrementSYB x = go where
852         go :: ∀c. Data c => c -> c
853         go x = mkT inc (gmapT go x)

```

854 The resulting `go` function implements `everywhere` but specialized to `mkT inc`
855 for `f`. More importantly though, `go` does not involve any higher-rank polymor-
856 phism. Thus if we run the specializer on this code, we get the following which
857 contains a version of `go` specialized to the `[Int]` type.

```

858     incrementSYB :: [Int] -> [Int]
859     incrementSYB x = go[Int] where
860         go :: ∀c. Data c => c -> c
861         go x = mkT inc (gmapT go x)
862         go[Int] :: [Int] -> [Int]
863         go[Int] x = mkT inc (gmapT go x)

```

864 In the default GHC optimization pipeline, the specializer runs before the inlining
865 process in the simplifier. Thus, in order to get this code, we have to modify
866 GHC to run the specializer after inlining.

867 This version is not yet fully optimized, however, as `go[Int]` still contains
868 calls to the polymorphic functions `mkT` and `gmapT`. Since `mkT` is not recursive,
869 inlining and symbolically evaluating should expose the `cast` in `mkT` and then
870 allow us to evaluate the comparison of `TypeRep` objects in the `cast`. That would
871 transform `go[Int]` to the following.

```

872     go[Int] :: [Int] -> [Int]
873     go[Int] x = gmapT go x

```

874 Unfortunately, in our experiments with GHC the optimization process often
875 simplified `mkT` and the contained `cast` but did not do the final step of removing
876 the comparison over `TypeRep` objects. This seems to be due to the simplifier not
877 knowing how to symbolically evaluate the `fingerprintFingerprints` function
878 that sometimes arises when simplifying such code. Adding primitive simplifica-
879 tion rules such as those in Figure 10 is a relatively trivial extension to the GHC
880 optimizer and allows us to eliminate that part of the code.

881 Next, consider the call to `gmapT` inside `go[Int]`. It is also over the concrete
882 type `[Int]` so the simplifier statically computes the dictionary dispatch and
883 changes the code to use `gmapT[Int]`, the `gmapT` implementation in the `Data`
884 instance for `[Int]`. This results in the following code.

```

885     go[Int] :: [Int] -> [Int]
886     go[Int] x = gmapT[Int] go x

```

887 Unfortunately, the simplifier does not inline this invocation of `gmapT[Int]`. This
888 is because there is a cycle between the dictionary for `Data` at the `[Int]` type
889 and `gmapT[Int]`, so GHC marks one of them as a loop breaker in order to avoid
890 infinite inlinings. As GHC avoids making class dictionaries be loop breakers,
891 `gmapT[Int]` is marked. As a result, the simplifier does not inline `gmapT[Int]`.
892 If we overlook this obstacle and force `gmapT[Int]` to inline, then we get the
893 following code.

```

894     go[Int] :: [Int] -> [Int]
895     go[Int] [] = []
896     go[Int] (x : xs) = go x : go xs

```

897 This exposes two calls to `go`. One is on `x` and is over the `Int` type. The
898 other is on `xs` and is over the `[Int]` type. When the GHC specializer adds
899 specializations, it also adds rewrite rules for those specializations. The simplifier
900 uses these rewrite rules to convert the call to `go` over the `[Int]` type to `go[Int]`.

901 At this point, all of the dictionaries involving the `[Int]` type have been
902 removed by the GHC specializer and simplifier. However, note that we still
903 have the call `go x` which is on the `Int` type. In this particular case, the `Data`
904 and `Typeable` instances for `Int` are simple enough that later passes of the GHC
905 optimization pipeline do transform this into `inc x`. However, this is not always
906 the case. Consider, for example, what happens if `incrementSyB` is over `[[Int]]`
907 instead of `[Int]`. After the initial inlining of `everywhere` we end up with the
908 following code.

```

909     incrementSyB :: [[Int]] -> [[Int]]
910     incrementSyB x = go where
911         go :: ∀c. Data c => c -> c
912         go x = mkT inc (gmapT go x)

```

913 This code contains a manifest call to `go` at the `[[Int]]` type, but there is no
914 such manifest call to `go` at the `[Int]` type. This is because `go` is passed as a
915 polymorphic argument to `gmapT`. This is exactly the sort of argument that the
916 GHC specializer does not know how to handle. With `everywhere` we were able
917 get around this by doing a static argument transformation, but `gmapT` is a class
918 method so we cannot do the same. Thus when the GHC specializer runs, `go` is
919 specialized only at the `[[Int]]` type. This results in the following code.

```

920     incrementSyB :: [[Int]] -> [[Int]]
921     incrementSyB x = go[[Int]] where
922         go :: ∀c. Data c => c -> c
923         go x = mkT inc (gmapT go x)
924         go[[Int]] :: [[Int]] -> [[Int]]
925         go[[Int]] x = mkT inc (gmapT go x)

```

926 As before, we simplify away the `mkT` in `go[[Int]]`. Since the `gmapT` in `go[[Int]]`
 927 is over the concrete type `[[Int]]`, we also simplify that, which results in the
 928 following for `go[[Int]]`.

```
929     go[[Int]] :: [[Int]] -> [[Int]]
930     go[[Int]] [] = []
931     go[[Int]] (x : xs) = go x : go xs
```

932 The call `go xs` is over the type `[[Int]]` for which we have a specialization so
 933 this is then turned into the following.

```
934     go[[Int]] :: [[Int]] -> [[Int]]
935     go[[Int]] [] = []
936     go[[Int]] (x : xs) = go x : go[[Int]] xs
```

937 As before, we have the call `go x` over a type which does not have a specialization.
 938 This time, however, it is over the type `[Int]`, which is complicated enough that
 939 it will not be optimized by the later stages of the pipeline.

940 Thus one pass of the specializer is not sufficient to optimize this SYB-style
 941 code. However, note that after specialization and simplification, we now have
 942 a manifest call of `go` on the `[Int]` type. The GHC specializer knows how to
 943 handle this sort of call. Thus if we run the specializer over this code, we get a
 944 `go[Int]` function. After another round of inlining and simplification this results
 945 in the following.

```
946     incrementsyB :: [[Int]] -> [[Int]]
947     incrementsyB x = go where
948         go :: ∀c. Data c => c -> c
949         go x = mkT inc (gmapT go x)
950         go[[Int]] :: [[Int]] -> [[Int]]
951         go[[Int]] [] = []
952         go[[Int]] (x : xs) = go[Int] x : go[[Int]] xs
953         go[Int] :: [Int] -> [Int]
954         go[Int] [] = []
955         go[Int] (x : xs) = go x : go[Int] xs
```

956 This in turn has exposed a call to `go` on the `Int` type inside `go[Int]`. This is
 957 simple enough that we can either leave this for later passes in the optimization
 958 pipeline or we can invoke the specializer and simplifier again, which results in
 959 the following.

```
960     incrementsyB :: [[Int]] -> [[Int]]
961     incrementsyB x = go where
962         go[[Int]] :: [[Int]] -> [[Int]]
963         go[[Int]] [] = []
964         go[[Int]] (x : xs) = go[Int] x : go[[Int]] xs
965         go[Int] :: [Int] -> [Int]
```



```

966     go[Int] [] = []
967     go[Int] (x : xs) = goInt x : go[Int] xs
968     goInt :: Int -> Int
969     goInt x = inc x

```

Each time we invoke the specializer and then the simplifier, we potentially uncover more types at which to specialize `go`. Thus for complex types we may need to perform this iteration multiple times.

In summary, while the default GHC optimization pipeline does not effectively optimize SYB-style code, a few modifications are sufficient to do so. First, we static-argument transform traversals like `everywhere`. This allows them to be inlined, which effectively specializes them to their first argument. Second, we run the specializer after the simplifier so that the `go` function resulting from the inlining of `everywhere` is specialized to particular types. Third, we run the simplifier again with two modifications. We symbolically evaluate `TypeRep`, `TyCon`, and `Fingerprint` calculations using rules such as those in Figure 10. We also inline the type-specific instances of `gfoldl`, `gmapT`, `gmapQ`, and `gmapM` even though the cycles in these would normally prevent it. Fourth, we iterate this specialization and simplification process with new iterations each time it reveals a new type at which to specialize `go`. The end result of all this is a fully optimized version of the code that contains no `Data` or `Typeable` class dispatches and no `TypeRep`, `TyCon`, or `Fingerprint` computations.

Note that using the specializer in this way is not as general as the optimization described in Section 4. It relies on the code having a structure similar to `everywhere` that we can static-argument transform. Thus, this technique will not be as effective when the code is not or cannot be of such a form.

7.1. Benchmarks

In order to test the efficacy of this technique, we created a plugin for GHC that iterates between specialization and simplification. This plugin also inlined of any instance specific implementations of `gfoldl`, `gmapT`, `gmapQ`, or `gmapM` that the GHC simplifier did not and used the simplification rules in Figure 10 to eliminate `TypeRep`, `TyCon`, and `Fingerprint` computations. We then re-ran the benchmarks in Section 5 with this plugin and a modified version of the SYB library with the static argument transformation applied to all of the traversal schemes.

We ran the benchmarks both with (`Spec+`) and without (`Spec`) the extra simplification rules for `TypeRep`, `TyCon`, and `Fingerprint` computations. Also, since these tests involve a different version of SYB than in Section 5, we reran the benchmarks without our optimization or plugin (SYB). The results are plotted in Figure 12 and are normalized relative to the handwritten code (`Hand`) just as in Figure 11.

Overall, both `Spec` and `Spec+` performed well. `Spec+` ran on par with the handwritten code across the board. `Spec` also ran on par with `Hand` for most benchmarks, but failed to significantly improve `RmWeights` or `SelectInt`. An inspection of the resulting core from each of these benchmarks reveals why. In

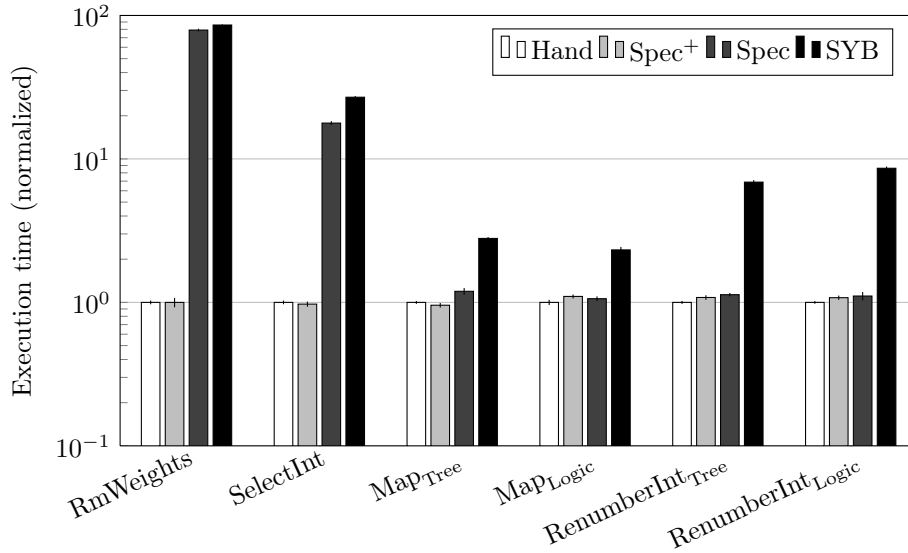


Figure 12: Benchmarks Results

1010 most of the benchmarks, the `TypeRep`, `TyCon`, and `Fingerprint` computations
 1011 that are left over after our plugin iteratively runs the specializer are simple
 1012 enough that they are eliminated by later passes of the compiler. In `RmWeights`
 1013 and `SelectInt`, however, the computations are more complex and are not elim-
 1014 inated.

1015 8. Related work

1016 Generic-programming systems in Haskell are often slow relative to hand-
 1017 written code. There has been a significant amount of work on designing more
 1018 efficient generic-programming systems (Mitchell and Runciman, 2007; Brown
 1019 and Sampson, 2009; Chakravarty et al., 2009; Augustsson, 2011; Adams and
 1020 DuBuisson, 2012), but there is little work on optimizing a pre-existing gener-
 1021 ic-programming system as we do here. Magalhães (2013) shows how to optimize
 1022 the `generic-deriving` system by using standard compiler optimizations, but
 1023 notes that his techniques are not sufficient to optimize SYB-style code. Ali-
 1024 marine and Smetsers (2004) have developed a similar optimization system for
 1025 generics in the Clean language.

1026 Our optimization is related to class dictionary specialization (Jones, 1995)
 1027 and call-pattern specialization (Peyton Jones, 2007). However, our optimization
 1028 specializes and memoizes over any expression with an undesirable type whereas
 1029 Jones (1995) specializes over only class dictionaries, and Peyton Jones (2007)
 1030 specializes over only manifest constructors. As discussed in Section 7, dictionary
 1031 specialization is not sufficient to optimize SYB-style code, but using the lessons

1032 and experience from our work we were able to find modifications of the GHC
 1033 specializer to effectively optimize SYB-style code.

1034 In a broader sense, our optimization is a form of partial evaluation (Jones
 1035 et al., 1993) with a binding-time analysis that uses type information to determine
 1036 whether code should be statically computed at compile time or dynamically
 1037 evaluated at runtime. However, because we use domain-specific knowledge, our
 1038 algorithm can be simpler and more direct than traditional partial evaluation.

1039 Our optimization can also be seen as a limited form of supercompilation
 1040 (Turchin, 1979). Like Bolingbroke and Peyton Jones (2010), we implement a
 1041 memoization scheme to ensure terms are optimized only once. We can draw
 1042 direct connections to many of the rules in Jonsson and Nordlander (2011). For
 1043 example, rules R1, R5, R6, R12, and R13 in that work correspond to several
 1044 of the forcing rules in our Figure 7. Rules R2, R11, and R15 correspond to
 1045 the primitive simplification rules in Figure 8. Rule R8 and R9 respectively
 1046 correspond to SUBSTLIT and SUBSTVAR in Figure 8.

1047 However, unlike general partial evaluation, we take advantage of domain
 1048 knowledge about SYB-style code. In particular, we use the types of expressions
 1049 to direct the optimization and start symbolically evaluating an expression only
 1050 when it is a form that eliminates an expression with an undesirable type. In
 1051 theory, we face the same problem of code explosion that supercompilers do, but
 1052 as we operate in the more limited setting of SYB-style code, this problem is
 1053 easier to handle.

1054 9. Conclusion

1055 SYB is widely used in the Haskell community. Its poor performance, how-
 1056 ever, can be a serious drawback in practical systems. Nevertheless, by using
 1057 domain specific knowledge about SYB-style code, we can design an optimization
 1058 that transforms this code to be as fast as equivalent handwritten, non-generic
 1059 code.

1060 The essential task of this optimization is the elimination of certain types by a
 1061 compile-time symbolic evaluation of the appropriate parts of the code. We have
 1062 first implemented this optimization in the HERMIT plugin for GHC. The inter-
 1063 active manipulation that HERMIT supports made it easy to rapidly prototype
 1064 such an optimization and trace how it transforms the code. This interactive ap-
 1065 proach was instrumental in empirically discovering the appropriate optimization
 1066 steps for optimizing SYB-style code. For example, a number of auxiliary code
 1067 simplifications had to be introduced in order to make it possible for the core
 1068 rules to run. We have then explored how to integrate this optimization directly
 1069 into GHC in order to obtain similar results without depending on HERMIT. In
 1070 the future, we plan to investigate how to make our optimization applicable to
 1071 other domains where expressions of certain types need to be eliminated.

1072 Benchmarks show that this optimization significantly improves the perform-
 1073 ance of several typical SYB-style traversals to closely match that of handwrit-
 1074 ten, non-generic code. In so doing, this optimization changes SYB from being

1075 one of the slowest generic-programming systems in the Haskell community to
 1076 being one of the fastest.

1077 **Acknowledgements** Our thanks go to Simon Peyton Jones for his help
 1078 with testing the GHC optimizer. The work presented in this paper was sup-
 1079 ported in part by the National Science Foundation (NSF) grants CCF-1218605
 1080 and 1117569, the Rockwell Collins contract 4504813093, and the Engineering
 1081 and Physical Sciences Research Council (EPSRC) grant EP/J010995/1.

1082 References

1083 Adams, M. D., DuBuisson, T. M., 2012. Template your boilerplate: using Tem-
 1084 plate Haskell for efficient generic programming. In: Proceedings of the 2012
 1085 symposium on Haskell symposium. Haskell '12. ACM, New York, NY, USA,
 1086 pp. 13–24.

1087 Adams, M. D., Farmer, A., Magalhães, J. P., 2014. Optimizing SYB is easy! In:
 1088 Proceedings of the ACM SIGPLAN 2014 Workshop on Partial Evaluation and
 1089 Program Manipulation. PEPM '14. ACM, New York, NY, USA, pp. 71–82.

1090 Alimarine, A., Smetsers, S., 2004. Optimizing generic functions. In: Kozen, D.
 1091 (Ed.), Mathematics of Program Construction. Vol. 3125 of Lecture Notes in
 1092 Computer Science. Springer Berlin Heidelberg, pp. 16–31.

1093 Augustsson, L., Nov. 2011. Geniplate version 0.6.0.0.
 1094 URL <http://hackage.haskell.org/package/geniplate/>

1095 Baker-Finch, C., Glynn, K., Peyton Jones, S., Mar. 2004. Constructed product
 1096 result analysis for Haskell. Journal of Functional Programming 14 (02), 211–
 1097 245.

1098 Bolingbroke, M., Peyton Jones, S., 2010. Supercompilation by evaluation. In:
 1099 Proceedings of the third ACM Haskell symposium on Haskell. Haskell '10.
 1100 ACM, New York, NY, USA, pp. 135–146.

1101 Brown, N. C. C., Sampson, A. T., 2009. Alloy: fast generic transformations for
 1102 Haskell. In: Proceedings of the 2nd ACM SIGPLAN symposium on Haskell.
 1103 Haskell '09. ACM, New York, NY, USA, pp. 105–116.

1104 Chakravarty, M. M. T., Ditu, G. C., Leshchinskiy, R., 2009. Instant generics:
 1105 Fast and easy, available at [http://www.cse.unsw.edu.au/~chak/papers/
 1106 instant-generics.pdf](http://www.cse.unsw.edu.au/~chak/papers/instant-generics.pdf).

1107 Farmer, A., Gill, A., Komp, E., Sculthorpe, N., 2012. The HERMIT in the
 1108 machine: a plugin for the interactive transformation of GHC core language
 1109 programs. In: Proceedings of the 2012 Haskell Symposium. Haskell '12. ACM,
 1110 New York, NY, USA, pp. 1–12.

- 1111 GHC Team, 2013. The Glorious Glasgow Haskell Compilation System User’s
1112 Guide, Version 7.6.2.
1113 URL <http://www.haskell.org/ghc>
- 1114 Gill, A., 2009. A Haskell hosted DSL for writing transformation systems. In:
1115 Taha, W. M. (Ed.), *Domain-Specific Languages*. Vol. 5658 of *Lecture Notes*
1116 *in Computer Science*. Springer Berlin Heidelberg, pp. 285–309.
- 1117 Hinze, R., Löh, A., Oliveira, B. C. d. S., 2006. “scrap your boilerplate” reloaded.
1118 In: Hagiya, M., Wadler, P. (Eds.), *Functional and Logic Programming*. Vol.
1119 3945 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, pp.
1120 13–29.
- 1121 Industrial Haskell Group, 2013. Hackage: Total downloads. Accessed on October
1122 8, 2013.
1123 URL <http://hackage.haskell.org/packages/top>
- 1124 Jones, M. P., Sep. 1995. Dictionary-free overloading by partial evaluation. *LISP*
1125 *and Symbolic Computation* 8 (3), 229–248.
- 1126 Jones, N. D., Gomard, C. K., Sestof, P., 1993. *Partial Evaluation and Automatic*
1127 *Program Generation*. Prentice-Hall International Series in Computer Science.
1128 Prentice Hall.
- 1129 Jonsson, P. A., Nordlander, J., 2011. Taming code explosion in supercompila-
1130 tion. In: *Proceedings of the 20th ACM SIGPLAN workshop on Partial evalua-*
1131 *tion and program manipulation*. PEPM ’11. ACM, New York, NY, USA,
1132 pp. 33–42.
- 1133 Lämmel, R., Peyton Jones, S., 2003. Scrap your boilerplate: a practical design
1134 pattern for generic programming. In: *Proceedings of the 2003 ACM SIGPLAN*
1135 *international workshop on Types in languages design and implementation*.
1136 *TLDI ’03*. ACM, New York, NY, USA, pp. 26–37.
- 1137 Lämmel, R., Peyton Jones, S., 2004. Scrap more boilerplate: reflection, zips, and
1138 generalised casts. In: *Proceedings of the ninth ACM SIGPLAN international*
1139 *conference on Functional programming*. ICFP ’04. ACM, New York, NY,
1140 USA, pp. 244–255.
- 1141 Magalhães, J. P., 2013. Optimisation of generic programs through inlining. In:
1142 *Implementation and Application of Functional Languages*.
- 1143 Magalhães, J. P., Holdermans, S., Jeurings, J., Löh, A., 2010. Optimizing generics
1144 is easy! In: *Proceedings of the 2010 ACM SIGPLAN workshop on Partial*
1145 *evaluation and program manipulation*. PEPM ’10. ACM, New York, NY,
1146 USA, pp. 33–42.
- 1147 Mitchell, N., Runciman, C., 2007. Uniform boilerplate and list processing. In:
1148 *Proceedings of the ACM SIGPLAN workshop on Haskell workshop*. Haskell
1149 ’07. ACM, New York, NY, USA, pp. 49–60.

- 1150 Peyton Jones, S., 2007. Call-pattern specialisation for Haskell programs. In:
1151 Proceedings of the 12th ACM SIGPLAN international conference on Func-
1152 tional programming. ICFP '07. ACM, New York, NY, USA, pp. 327–337.
- 1153 Peyton Jones, S., Marlow, S., Jul. 2002. Secrets of the Glasgow Haskell Compiler
1154 inliner. *Journal of Functional Programming* 12 (4–5), 393–434.
- 1155 Rodriguez, A., Jeuring, J., Jansson, P., Gerdes, A., Kiselyov, O., Oliveira, B.
1156 C. d. S., 2008. Comparing libraries for generic programming in Haskell. Tech.
1157 Rep. UU-CS-2008-010, Utrecht University.
- 1158 Rodriguez Yakushev, A., 2009. Towards getting generic programming ready for
1159 prime time. Ph.D. thesis, Utrecht University.
- 1160 Santos, A., 1995. Compilation by transformation in non-strict functional lan-
1161 guages. Ph.D. thesis, University of Glasgow.
- 1162 Sculthorpe, N., Farmer, A., Gill, A., 2013. The HERMIT in the tree: Mechaniz-
1163 ing program transformations in the GHC core language. In: *Implementation*
1164 *and Application of Functional Languages*.
- 1165 Sculthorpe, N., Frisby, N., Gill, A., 2014. The Kansas University Rewrite
1166 Engine: A Haskell-embedded strategic programming language with custom
1167 closed universes. *Journal of Functional Programming*.
- 1168 Turchin, V. F., Feb. 1979. A supercompiler system based on the language RE-
1169 FAL. *ACM SIGPLAN Notices* 14 (2), 46–54.
- 1170 Vytiniotis, D., Peyton Jones, S., Magalhães, J. P., 2012. Equality proofs and
1171 deferred type errors: a compiler pearl. In: *Proceedings of the 17th ACM*
1172 *SIGPLAN international conference on Functional programming. ICFP '12.*
1173 *ACM, New York, NY, USA, pp. 341–352.*