



Universiteit Utrecht

[Faculty of Science
Information and Computing Sciences]

Optimizing Generics Is Easy!

José Pedro Magalhães

Joint work with

Stefan Holdermans Johan Jeuring Andres Löh

Dept. of Information and Computing Sciences, Universiteit Utrecht

P.O. Box 80.089, 3508 TB Utrecht, The Netherlands

Webpage: <http://www.dreixel.net>

PEPM'2010, Madrid, 18/01/2010

Overview

Generic programming

Optimizing generics through inlining

A benchmark suite for generics

Conclusions and future work



1. Generic programming



What is datatype-generic programming?

- ▶ Programming with the *structure* of types
- ▶ Conversion functions map user datatypes to/from *representation types*
- ▶ Generic functions are defined on representation types

Generic functions work for all types for which we can write conversion functions.



Performance of generics

Generic functions tend to be slow...

Generic function	Overhead from popular GP library	Overhead from fast GP library	Overhead after our optimization
<i>show</i>	3	0.5	0
<i>update</i>	7	1	0 (times slower)

... but we will see how we can change that.



Generic representation I

Because generic programming is all about types, we use Haskell: a lazy, *pure* functional programming language with strong static typing.

Haskell supports the definition of algebraic datatypes, like:

```
data Logic = Logic ∨ Logic -- disjunction
           | Var String    -- variables
           | Not Logic     -- negation
```

To represent these, we need to know how to handle:

- ▶ Different alternatives: disjoint sums.
- ▶ Arguments of a constructor: products.
- ▶ Constructors and field labels.
- ▶ Primitive types: *String*, *Int*, ...
- ▶ ...



Generic representation II

We need to translate every datatype to this set of constructs and apply the appropriate code in the right place.

Haskell's **data** construct combines several features: type abstraction, type recursion, (labeled) sums, and (possibly labeled) products, but they are essentially *sums of products*.



Generic representation II

We need to translate every datatype to this set of constructs and apply the appropriate code in the right place.

Haskell's **data** construct combines several features: type abstraction, type recursion, (labeled) sums, and (possibly labeled) products, but they are essentially *sums of products*.

To represent them we can use the following *representation datatypes*:

data $(f + g)$	$r = L (f r) \mid R (g r)$	-- Choice
data $(f \times g)$	$r = f r \times g r$	-- Multiple arguments
data $K a$	$r = K a$	-- Constants
data I	$r = I r$	-- Recursive occurrences
data U	$r = U$	-- No arguments
data $C f$	$r = C \text{String} (f r)$	-- Constructors



Generic representation III

We encapsulate conversion to and from the generic representation using a type class. The generic type is given using a type family:

```
type family PF a :: * → *
```

```
class Regular a where
```

```
  from :: a      → PF a a
```

```
  to   :: PF a a → a
```



Generic representation IV

Back to our *Logic* example:

type instance *PF Logic* = $C (I \times I)$ -- disjunction
+ $C (K \textit{String})$ -- variables
+ $C I$ -- negation



Generic representation IV

Back to our *Logic* example:

type instance *PF Logic* = $C (I \times I)$ -- disjunction
+ $C (K \text{ String})$ -- variables
+ $C I$ -- negation

instance *Regular Logic* where

from $(p \vee q) = L (C \text{ "V" } ((I p) \times (I q)))$

from $(\text{Var } x) = R (L (C \text{ "Var" } (K x)))$

from $(\text{Not } p) = R (R (C \text{ "Not" } (I p)))$

to $(L (C _ ((I p) \times (I q)))) = p \vee q$

to $(R (L (C _ (K x)))) = \text{Var } x$

to $(R (R (C _ (I p)))) = \text{Not } p$



Generic functions: *gmap*

Now we can write generic functions:

```
class GMap f where  
  gmap :: (a → b) → f a → f b
```



Generic functions: *gmap*

Now we can write generic functions:

class *GMap* *f* **where**

gmap :: (*a* → *b*) → *f* *a* → *f* *b*

instance *GMap* *I* **where**

gmap *f* (*I* *r*) = *I* (*f* *r*)

instance *GMap* (*K* *a*) **where**

gmap _ (*K* *x*) = *K* *x*

instance *GMap* *U* **where**

gmap _ *U* = *U*

instance (*GMap* *f*, *GMap* *g*) ⇒ *GMap* (*f* + *g*) **where**

gmap *f* (*L* *x*) = *L* (*gmap* *f* *x*)

gmap *f* (*R* *x*) = *R* (*gmap* *f* *x*)

instance (*GMap* *f*, *GMap* *g*) ⇒ *GMap* (*f* × *g*) **where**

gmap *f* (*x* × *y*) = *gmap* *f* *x* × *gmap* *f* *y*



Generic functions: *gshow* I

Another function we can define is generic show. For that we need to use constructor information.

```
class GShow f where  
  gshow f :: (a → String) → f a → String
```



Generic functions: *gshow* I

Another function we can define is generic show. For that we need to use constructor information.

class *GShow* *f* **where**

gshowf :: (*a* → *String*) → *f a* → *String*

instance *GShow* *I* **where**

gshowf *f* (*I* *r*) = *f r*

instance (*Show* *a*) ⇒ *GShow* (*K a*) **where**

gshowf _ (*K* *x*) = *show* *x*

instance *GShow* *U* **where**

gshowf _ *U* = ""

instance (*GShow* *f*) ⇒ *GShow* (*C f*) **where**

gshowf *f* (*C* *conName* *x*) = "(" ++ *conName* ++ " "
++ *gshowf* *f* *x* ++ ")"



Generic functions: *gshow* II

instance (*GShow* *f*, *GShow* *g*) \Rightarrow *GShow* (*f* + *g*) **where**
 gshowf *f* (*L* *x*) = *gshowf* *f* *x*
 gshowf *f* (*R* *x*) = *gshowf* *f* *x*

instance (*GShow* *f*, *GShow* *g*) \Rightarrow *GShow* (*f* \times *g*) **where**
 gshowf *f* (*x* \times *y*) = *gshowf* *f* *x* ++ " " ++ *gshowf* *f* *y*



Generic functions: *gshow* II

instance (*GShow* *f*, *GShow* *g*) \Rightarrow *GShow* (*f* + *g*) **where**
gshowf *f* (*L* *x*) = *gshowf* *f* *x*
gshowf *f* (*R* *x*) = *gshowf* *f* *x*

instance (*GShow* *f*, *GShow* *g*) \Rightarrow *GShow* (*f* \times *g*) **where**
gshowf *f* (*x* \times *y*) = *gshowf* *f* *x* ++ " " ++ *gshowf* *f* *y*

This function works only on the generic representations. For normal datatypes we first have to convert them:

gshow :: (*Regular* *a*, *GShow* (*PF* *a*)) \Rightarrow *a* \rightarrow *String*
gshow *x* = *gshowf* *gshow* (*from* *x*)

At the recursive occurrences we apply *gshow* again.



2. Optimizing generics through inlining



Efficiency I

While representation types are useful, they incur a performance penalty:

- ▶ Generic functions keep converting back and forth
- ▶ Generic representation types are present in the final generated code
- ▶ Overhead from using a generic programming library can range from 0.5–16 times slower, compared to a handwritten variant.



Efficiency II

Generic representation types should not be present in the generated code. Generic functions can be specialized to particular types.

We can see that if we inline definitions and apply equational reasoning we can remove the generic representations.

As an example, let us see one-level generic identity on the *Logic* datatype:

$$gid_{Logic} :: Logic \rightarrow Logic$$
$$gid_{Logic} = to \circ gmap\ id \circ from$$


Efficiency III

to (gmap id (from l))



Efficiency III

to (gmap id (from l))

\Rightarrow { choose l to be $p \vee q$ (other constructors similar) }
to (gmap id (from (p \vee q)))



Efficiency III

$to (gmap\ id\ (from\ l))$

$\Rightarrow \{ \text{choose } l \text{ to be } p \vee q \text{ (other constructors similar)} \}$
 $to (gmap\ id\ (from\ (p \vee q)))$

$\equiv \{ \text{definition of } from_{Logic} \}$
 $to (gmap\ id\ (L\ (I\ p \times I\ q)))$



Efficiency III

$to (gmap\ id\ (from\ l))$

$\Rightarrow \{ \text{choose } l \text{ to be } p \vee q \text{ (other constructors similar)} \}$
 $to (gmap\ id\ (from\ (p \vee q)))$

$\equiv \{ \text{definition of } from_{Logic} \}$
 $to (gmap\ id\ (L\ (I\ p \times I\ q)))$

$\equiv \{ \text{definition of } gmap_{+}, gmap_{\times} \}$
 $to (L\ (gmap\ id\ (I\ p) \times gmap\ id\ (I\ q)))$



Efficiency III

to (*gmap id (from l)*)

\Rightarrow { choose l to be $p \vee q$ (other constructors similar) }
to (*gmap id (from (p \vee q))*)

\equiv { definition of *from_{Logic}* }
to (*gmap id (L (I p \times I q))*)

\equiv { definition of *gmap₊*, *gmap _{\times}* }
to (*L (gmap id (I p) \times gmap id (I q))*)

\equiv { definition of *gmap_I* }
to (*L (I (id p) \times (I (id q)))*)



Efficiency III

to (*gmap id (from l)*)

\Rightarrow { choose l to be $p \vee q$ (other constructors similar) }
to (*gmap id (from (p \vee q))*)

\equiv { definition of *from_{Logic}* }
to (*gmap id (L (I p \times I q))*)

\equiv { definition of *gmap₊*, *gmap _{\times}* }
to (*L (gmap id (I p) \times gmap id (I q))*)

\equiv { definition of *gmap_I* }
to (*L (I (id p) \times (I (id q)))*)

\equiv { definition of *id*, *to_{Logic}* }
 $p \vee q$



Core code I

Can we not get the compiler to do the same for us? The core code GHC generates for our example

$$\begin{aligned}gid_{Logic} &:: Logic \rightarrow Logic \\gid_{Logic} &= to \circ gmap\ id \circ from\end{aligned}$$

is

$$\begin{aligned}gid_{Logic}^{01} &:: Logic \rightarrow Logic \\gid_{Logic}^{01} &= \lambda(x :: Logic) \rightarrow to\ (from\ x)\end{aligned}$$

This is good, but not ideal. We also know that $to_{Logic} \circ from_{Logic} \equiv id$.



Core code II

The problem is that the compiler is conservative with *inlining*—replacing function calls with their body. We can force inlining by tweaking some flags:

Flag	Default	Abbr.
<code>-funfolding-creation-threshold</code>	45	CT
<code>-funfolding-use-threshold</code>	6	UT

Compiling with `-O2 -funfolding-use-threshold=60` produces the wanted result:

$$\begin{aligned} \text{gid}_{\text{Logic}}^{\text{02UT60}} &:: \text{Logic} \rightarrow \text{Logic} \\ \text{gid}_{\text{Logic}}^{\text{02UT60}} &= \lambda(x :: \text{Logic}) \rightarrow x \end{aligned}$$


Core code III

For *gshow*, with standard optimizations we get:

$$\begin{aligned} gshow_{Logic}^{01} &:: Logic \rightarrow String \\ gshow_{Logic}^{01} &= \lambda(x :: Logic) \rightarrow \\ &\quad \text{case (from } x \text{) 'cast' (sym (trans ...)) of } w \{ \\ &\quad \quad L \ y \rightarrow \dots \\ &\quad \quad R \ y \rightarrow \dots \} \end{aligned}$$


Core code III

For *gshow*, with standard optimizations we get:

$$\begin{aligned} gshow_{Logic}^{01} &:: Logic \rightarrow String \\ gshow_{Logic}^{01} &= \lambda(x :: Logic) \rightarrow \\ &\quad \text{case (from } x) \text{ 'cast' (sym (trans ...)) of } w \{ \\ &\quad \quad L \ y \rightarrow \dots \\ &\quad \quad R \ y \rightarrow \dots \} \end{aligned}$$

But we can force inlining to obtain a better result:

$$\begin{aligned} gshow_{Logic}^{CT90UT30} &:: Logic \rightarrow String \\ gshow_{Logic}^{CT90UT30} &= \lambda(x :: Logic) \rightarrow \text{case } x \text{ of } w \{ \\ &\quad (\vee) \ p \ q \rightarrow (++) \dots gshow_{Logic}^{CT90UT30} \ p \ \dots gshow_{Logic}^{CT90UT30} \ q \ \dots \\ &\quad \text{Var } v \rightarrow (++) \dots show \ v \ \dots \\ &\quad \text{Not } p \rightarrow (++) \dots gshow_{Logic}^{CT90UT30} \ p \ \dots \\ &\quad \text{Const } b \rightarrow (++) \dots show \ b \ \dots \} \end{aligned}$$


3. A benchmark suite for generics



A benchmark suite for generics: functions

To visualize the impact of increased inlining we designed a benchmark suite of generic programs. We will show two functions:

- show* Requires constructor information, such as name and fixity.
- update* Transform all odd *Int* values by adding one to them, or prepend all non-empty *String* values with a "y".

In our paper we present also the results for generic equality, map and read.



A benchmark suite for generics: datatypes

We use two datatypes. The *Tree* datatype is a simple labeled binary leaf tree:

```
data Tree a = Bin a (Tree a) (Tree a) | Leaf
```

The *Logic* type is similar to the one we introduced before, only with more constructors:

```
data Logic = Impl Logic Logic | Equiv Logic Logic  
           | Conj Logic Logic | Disj Logic Logic  
           | Not Logic | Var String | T | F
```



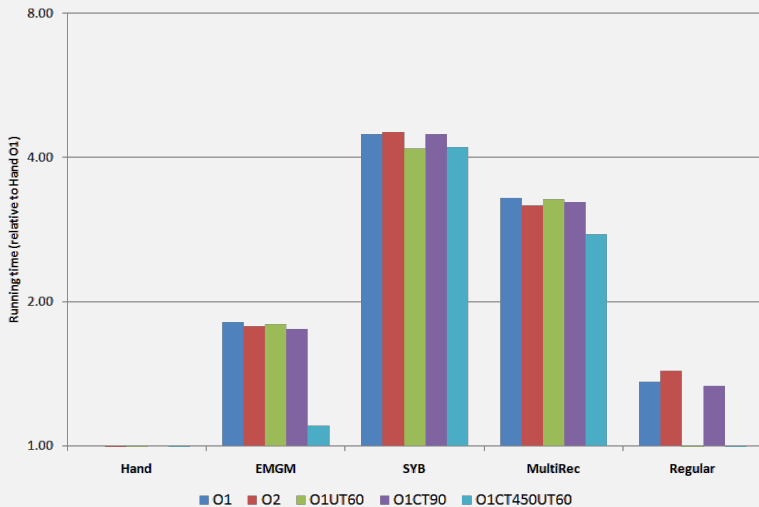
A benchmark suite for generics: libraries

We have chosen a few representative, mainstream, and maintained libraries to benchmark:

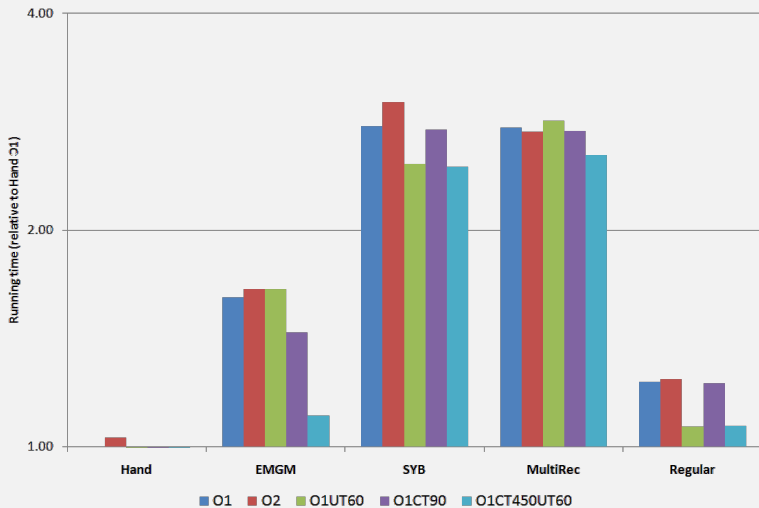
- emgm** Extensible and Modular Generics for the Masses. Its fundamental characteristic is to encode datatype representations through a type class.
- syb** Scrap Your Boilerplate is a very popular library based on generic combinators and type-safe cast. It comes with GHC.
- regular** The library described in the introduction.
- multirec** The first approach able to express mutually recursive datatypes. Structurally similar to **regular**, but makes use of a few more advanced concepts to deal with mutual recursion.



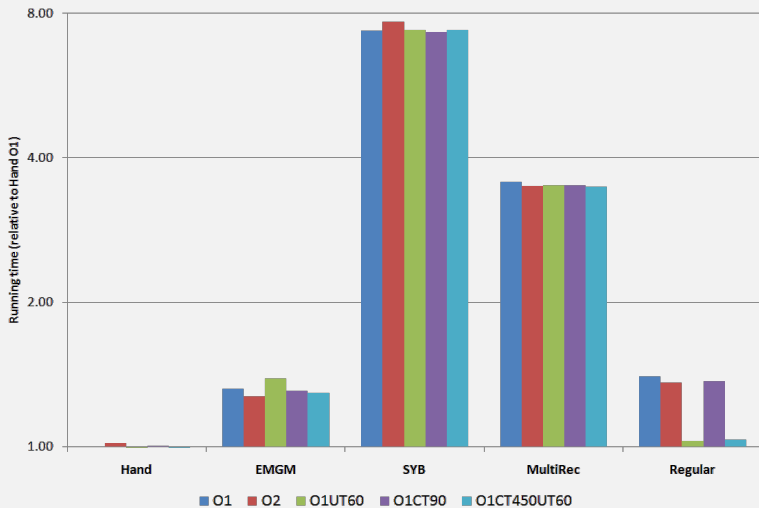
Results: *show for Tree*



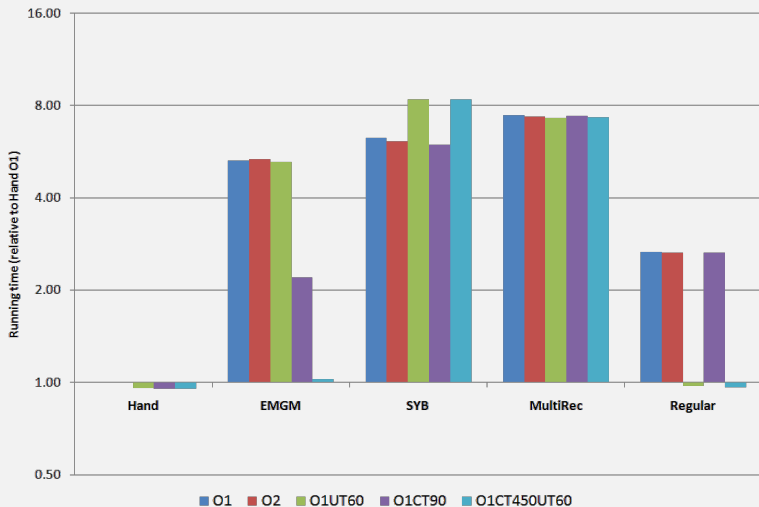
Results: *show for Logic*



Results: *update for Tree*



Results: *update for Logic*



4. Conclusions and future work



Conclusions

- ▶ Generic programs do not have to be slow
- ▶ Inlining is the way to go
- ▶ Facilities for inlining are already present in the compiler and can be reused for optimizing generics
- ▶ Both `emgm` and `regular` are fast and can be optimized to handwritten code speed with inlining
- ▶ The slowest (but most popular) generic programming library is `syb`
- ▶ `multirec` is not benefiting much from increased inlining, as opposed to the similar `regular` library



Future work

- ▶ Specifying the behavior of the inliner should be more localized: use the `INLINE` pragmas of the upcoming version of GHC
- ▶ Not all libraries benefit equally from increased inlining: why?
 - ▶ Are GADTs preventing inlining in `multirec`?
 - ▶ What can we do about `syb`?
- ▶ Investigate generic producers more thoroughly

