

## A Proofs

**Theorem 4** For any type  $\forall \alpha_1, \dots, \alpha_n : \tau$ ,  $n \geq 0$  we have:

$$\forall f_{\alpha_i}, g_{\alpha_i}, x, y : \bigwedge_i f_{\alpha_i} \circ g_{\alpha_i} = \text{id} \Rightarrow (\mathcal{P}[\tau](x, y) \Leftrightarrow x = \mathcal{F}[\tau](\bar{f}, \bar{g}) y)$$

where  $\bar{f} = f_{\alpha_1}, \dots, f_{\alpha_n}$  and  $\bar{g} = g_{\alpha_1}, \dots, g_{\alpha_n}$ .

*Proof* (by induction over the structure of the type  $\tau$ ):

- If  $\tau = \text{basic}$ , then  $\mathcal{F}[\tau](\bar{f}, \bar{g}) = \text{id}$  and  $\mathcal{P}[\tau](x, y) \Leftrightarrow (x = y)$ .
- If  $\tau = \alpha_i$ , then  $\mathcal{F}[\tau](\bar{f}, \bar{g}) = f_{\alpha_i}$  and  $\mathcal{P}[\tau](x, y) \Leftrightarrow (x = f_{\alpha_i}(y))$ .
- If  $\tau = \tau_1 \times \tau_2$ , then  $\mathcal{F}[\tau](\bar{f}, \bar{g}) = \mathcal{F}[\tau_1](\bar{f}, \bar{g}) \times \mathcal{F}[\tau_2](\bar{f}, \bar{g})$ , and

$$\begin{aligned} & \mathcal{P}[\tau](x, y) \\ \Leftrightarrow & \mathcal{P}[\tau_1](\pi_1(x), \pi_1(y)) \wedge \mathcal{P}[\tau_2](\pi_2(x), \pi_2(y)) && \text{from Th. 1} \\ \Leftrightarrow & (\pi_1(x) = \mathcal{F}[\tau_1](\bar{f}, \bar{g})(\pi_1(y))) \wedge (\pi_2(x) = \mathcal{F}[\tau_2](\bar{f}, \bar{g})(\pi_2(y))) && \text{induction hypothesis} \\ \Leftrightarrow & x = (\mathcal{F}[\tau_1](\bar{f}, \bar{g}) \times \mathcal{F}[\tau_2](\bar{f}, \bar{g})) y \\ \Leftrightarrow & x = \mathcal{F}[\tau](\bar{f}, \bar{g}) y \end{aligned}$$

- If  $\tau = \tau_1 \rightarrow \tau_2$ , then  $\mathcal{F}[\tau](\bar{f}, \bar{g}) = \lambda h. \mathcal{F}[\tau_2](\bar{f}, \bar{g}) \circ h \circ \mathcal{F}[\tau_1](\bar{g}, \bar{f})$ , and

$$\begin{aligned} & \mathcal{P}[\tau](k, h) \\ \Leftrightarrow & \mathcal{P}[\tau_1](x, y) \Rightarrow \mathcal{P}[\tau_2](k(x), h(y)) && \text{from Th. 1} \\ \Leftrightarrow & x = \mathcal{F}[\tau_1](\bar{f}, \bar{g}) y \Rightarrow k(x) = \mathcal{F}[\tau_2](\bar{f}, \bar{g})(h(y)) && \text{induction hypothesis} \\ \Leftrightarrow & k \circ \mathcal{F}[\tau_1](\bar{f}, \bar{g}) = \mathcal{F}[\tau_2](\bar{f}, \bar{g}) \circ h \\ \Leftrightarrow & k \circ \mathcal{F}[\tau_1](\bar{f}, \bar{g}) \circ \mathcal{F}[\tau_1](\bar{g}, \bar{f}) = \mathcal{F}[\tau_2](\bar{f}, \bar{g}) \circ h \circ \mathcal{F}[\tau_1](\bar{g}, \bar{f}) && (*) \\ \Leftrightarrow & k \circ \mathcal{F}[\tau_1](\overline{f \circ g}, \overline{f \circ g}) = \mathcal{F}[\tau_2](\bar{f}, \bar{g}) \circ h \circ \mathcal{F}[\tau_1](\bar{g}, \bar{f}) && \text{from Eq. 5} \\ \Leftrightarrow & k \circ \mathcal{F}[\tau_1](\overline{\text{id}}, \overline{\text{id}}) = \mathcal{F}[\tau_2](\bar{f}, \bar{g}) \circ h \circ \mathcal{F}[\tau_1](\bar{g}, \bar{f}) && \text{since } f_i \circ g_i = \text{id} \\ \Leftrightarrow & k \circ \text{id} = \mathcal{F}[\tau_2](\bar{f}, \bar{g}) \circ h \circ \mathcal{F}[\tau_1](\bar{g}, \bar{f}) && \text{from Eq. 4} \\ \Leftrightarrow & k = \mathcal{F}[\tau](\bar{f}, \bar{g}) h \end{aligned}$$

The equivalence (\*) is true because  $\text{range}(h) = \text{domain}(f) \Rightarrow (f = g \Leftrightarrow f \circ h = g \circ h)$ .

- If  $\tau = T(\tau')$ , then  $\mathcal{F}[\tau](\bar{f}, \bar{g}) = \text{map}^T(\mathcal{F}[\tau'](\bar{f}, \bar{g}))$ , and

$$\begin{aligned} & \mathcal{P}[\tau](x, y) \\ \Leftrightarrow & \forall h : (\forall z : \mathcal{P}[\tau'](h(z), z)) \Rightarrow x = \text{map}^T(h) y && \text{from Th. 1} \\ \Leftrightarrow & \forall h : (\forall z : h(z) = \mathcal{F}[\tau'](\bar{f}, \bar{g}) z) \Rightarrow x = \text{map}^T(h) y && \text{induction hypothesis} \\ \Leftrightarrow & x = \text{map}^T(\mathcal{F}[\tau'](\bar{f}, \bar{g})) y \\ \Leftrightarrow & x = \mathcal{F}[\tau](\bar{f}, \bar{g}) y \end{aligned} \quad \square$$

- [3] L. Fegaras and T. Sheard. Revisiting Catamorphisms over Datatypes with Embedded Functions. In *23rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, St. Petersburg Beach, Florida*, January 1996. To Appear.
- [4] L. Fegaras, T. Sheard, and T. Zhou. Improving Programs which Recurse over Multiple Inductive Structures. In *ACM SIGPLAN Workshop on Partial Evaluation and Semantics-Based Program Manipulation, Orlando, Florida*, pp 21–32, June 1994.
- [5] A. Gill, J. Launchbury, and S. Peyton Jones. A Short Cut to Deforestation. *Sixth Conference on Functional Programming Languages and Computer Architecture, Copenhagen, Denmark*, pp 223–232, June 1993.
- [6] Z. Hu, H. Iwasaki, and M. Takeichi. Deriving Structural Hylomorphisms from Recursive Definitions. In *ACM SIGPLAN International Conference on Functional Programming (ICFP'96), Philadelphia*, pp 73–82. ACM Press, May 1996.
- [7] Z. Hu, H. Iwasaki, M. Takeichi, and A. Takano. Tupling Calculation Eliminates Multiple Data Traversals. In *ACM SIGPLAN International Conference on Functional Programming*, pp 164–175, Amsterdam, The Netherlands, June 1997. ACM Press.
- [8] J. Launchbury and T. Sheard. Warm Fusion. *Seventh Conference on Functional Programming Languages and Computer Architecture, La Jolla, California*, June 1995.
- [9] E. Meijer, M. Fokkinga, and R. Paterson. Functional Programming with Bananas, Lenses, Envelopes and Barbed Wire. In *Proceedings of the 5th ACM Conference on Functional Programming Languages and Computer Architecture, Cambridge, Massachusetts*, pp 124–144, August 1991. LNCS 523.
- [10] S. Peyton Jones and J. Launchbury. Unboxed Values as First Class Citizens in a Non-strict Functional Language. *Fifth Conference on Functional Programming Languages and Computer Architecture, Cambridge, MA*, pp 636–665, August 1991.
- [11] T. Sheard and L. Fegaras. A Fold for All Seasons. *Sixth Conference on Functional Programming Languages and Computer Architecture, Copenhagen, Denmark*, pp 233–242, June 1993.
- [12] P. Wadler. Theorems for Free! *Fourth Conference on Functional Programming Languages and Computer Architecture, Imperial College, London*, September 1989.
- [13] P. Wadler. Deforestation: Transforming Programs to Eliminate Trees. *Proceedings of the 2nd European Symposium on Programming, Nancy, France*, pp 344–358, March 1988. LNCS 760.

$$\begin{aligned}\text{reverse}'(f,g) [] w &= f(w) \\ \text{reverse}'(f,g) (a:x) w &= g(\text{reverse}'(f,g) x (a:w))\end{aligned}$$

If we use the regular fusion algorithm, we derive:

$$\text{len}(\text{reverse } x [] ) = \text{reverse}'(\text{len},\text{id}) x []$$

which uses the accumulator in the same way as **reverse** does: it starts with an empty list and at each step it conses an element to the accumulator. At the end, it returns the length of the accumulator. We can avoid building the list accumulator by using a simple integer accumulator instead:

$$\begin{aligned}h [] w &= w \\ h (a:x) w &= h x (1+w)\end{aligned}$$

Such a definition can be derived by abstracting the accumulator from each recursive call in **reverse'**:

$$\begin{aligned}\text{reverse}'(f,g) [] w &= f(w) \\ \text{reverse}'(f,g) (a:x) w &= g(a,w,\lambda z. \text{reverse}'(f,g) x z)\end{aligned}$$

As before, **reverse** can be expressed in terms of **reverse'**:

$$\text{reverse} = \text{reverse}'(\lambda z.z, \lambda(a,w,f). f(a:w))$$

The type of **reverse'** is:

$$\forall \alpha, \beta, \gamma. (\gamma \rightarrow \beta) \times (\alpha \times \gamma \times (\gamma \rightarrow \beta) \rightarrow \beta) \rightarrow \text{list}(\alpha) \rightarrow \gamma \rightarrow \beta$$

From Th. 4 we have  $\forall f_{\alpha_i}, x, y : \mathcal{P}[\tau](x, y) \Leftrightarrow x = \mathcal{F}[\tau](\overline{f}, \overline{\mathcal{INV}(f)}) y$ , since  $f_i \circ \mathcal{INV}(f_i) = \text{id}$ . The fusion algorithm can be extended accordingly to handle any function  $f : \tau_1 \rightarrow \dots \rightarrow \tau_n \rightarrow \alpha$ :

$$\begin{aligned}\forall h, x_i : h(f x_1 \dots x_n) &\rightarrow f(\mathcal{F}[\tau_1](h, \mathcal{INV}(h)) x_1) \dots (\mathcal{F}[\tau_n](h, \mathcal{INV}(h)) x_n) \\ \forall x : h(\mathcal{INV}(h) x) &\rightarrow x\end{aligned}$$

Therefore, for  $\alpha = ()$  and  $\beta = \gamma$ , the **reverse'** fusion rule is:

$$h(\text{reverse}'(f,g) x w) \rightarrow \text{reverse}'(h \circ f \circ \mathcal{INV}(h), h \circ g \circ (\text{id} \times \mathcal{INV}(h) \times (\lambda k. \mathcal{INV}(h) \circ k \circ h))) x (h(w))$$

Using this rule and after some simplifications, we get:

$$\text{len}(\text{reverse } x [] ) = \text{reverse}'(\lambda z.z, \lambda(a,w,f). f(1+w)) x 0$$

which is equivalent to the desired definition of **len(rev x)**.

## References

- [1] R. Burstall and J. Darlington. A Transformation System for Developing Recursive Programs. *Journal of the ACM*, 24(1):44–67, January 1977.
- [2] W. Chin. Safe Fusion of Functional Expressions. *Proceedings of the ACM Symposium on Lisp and Functional Programming, San Francisco, California*, pp 11–20, June 1992.

$$\begin{aligned} \text{append}'(f,g) [] y &= f(y) \\ \text{append}'(f,g) (a:x) y &= g(a,x,y,\text{append}'(f,g) x y) \end{aligned}$$

If we apply the fusion algorithm once again to fuse `len` with `append`, we get:

$$\begin{aligned} &\text{len}(\text{append} (\mathcal{I}\mathcal{N}\mathcal{V}(\text{len}) r) [a]) \\ &= \text{len}(\text{append}'(\lambda y. y, \lambda(b,x,y,s). b:s) (\mathcal{I}\mathcal{N}\mathcal{V}(\text{len}) r) [a]) \\ &= \text{append}'(\lambda y. \text{len}(y), \lambda(b,x,y,s). \text{len}(b:(\mathcal{I}\mathcal{N}\mathcal{V}(\text{len}) s))) (\mathcal{I}\mathcal{N}\mathcal{V}(\text{len}) r) [a] \\ &= \text{append}'(\lambda y. \text{len}(y), \lambda(b,x,y,s). 1+s) (\mathcal{I}\mathcal{N}\mathcal{V}(\text{len}) r) [a] \end{aligned}$$

Therefore, `len(rev(x))` is

$$\text{rev}'(\lambda(). 0, \lambda(a,x,r). \text{append}'(\lambda y. \text{len}(y), \lambda(b,x,y,s). 1+s) (\mathcal{I}\mathcal{N}\mathcal{V}(\text{len}) r) [a]) x$$

which contains a term  $\mathcal{I}\mathcal{N}\mathcal{V}(\text{len})$  that has not been cancelled out. This problem occurs whenever the result of a recursive call of a function is handled by another recursive function in a non-trivial way (i.e., when it is deconstructed and/or recursed upon it). In that case, the intermediate data structure has a true computational value in the program and cannot be eliminated. All the other cases can be handled effectively by the fusion algorithm.

There are two ways to handle cases like this: One is to actually synthesize  $\mathcal{I}\mathcal{N}\mathcal{V}(h)$  from  $h$ . We may have many solutions for  $\mathcal{I}\mathcal{N}\mathcal{V}(h)$  that satisfy  $f \circ \mathcal{I}\mathcal{N}\mathcal{V}(h) = \text{id}$ . One solution is

$$\mathcal{I}\mathcal{N}\mathcal{V}(\text{len}) = \text{int\_fold}(\lambda(). [], \lambda n. 0:n)$$

where `int_fold` is the fold over integers. That way, we can fuse `append'` with  $\mathcal{I}\mathcal{N}\mathcal{V}(\text{len})$  yielding:

$$\text{len}(\text{rev}(x)) = \text{rev}'(\lambda(). 0, \lambda(a,x,r). \text{int\_fold}(\lambda(). \lambda y. \text{len}(y), \lambda f. \lambda y. 1+f(y)) r [a]) x$$

The other way is to use additional laws. These laws, called *promotion laws*, instead of fusing  $g$  with  $f$  in  $g(f(x))$ , they promote  $g$  to the right of  $f$ . That is, these laws take the form:  $g \circ f = h \circ F(g)$ , for some function  $h$  that does not depend on  $g$  and some function  $F$  that depends on  $g$ . For example, we have the law:

$$\text{len}(\text{append } x y) = (\text{len}(x)) + (\text{len}(y))$$

which actually states that `len` is a homomorphism. That way, `len(rev(x))` becomes:

$$\begin{aligned} &\text{rev}'(\lambda(). 0, \lambda(a,x,r). \text{len}(\text{append} (\mathcal{I}\mathcal{N}\mathcal{V}(\text{len}) r) [a])) x \\ &= \text{rev}'(\lambda(). 0, \lambda(a,x,r). (\text{len}(\mathcal{I}\mathcal{N}\mathcal{V}(\text{len}) r) + (\text{len}([a]))) x \\ &= \text{rev}'(\lambda(). 0, \lambda(a,x,r). r+1) x \end{aligned}$$

#### 4.4 Accumulator Deforestation

The linear version of list reverse is `rev x = reverse x []`, where `reverse` uses an extra accumulator:

$$\begin{aligned} \text{reverse} [] w &= w \\ \text{reverse} (a:x) w &= \text{reverse } x (a:w) \end{aligned}$$

Its recursive skeleton is:

Using  $\mathcal{I}\mathcal{N}\mathcal{V}(h)$ , the premise of Eq. 6 becomes:

$$\begin{aligned}
& h \circ g_i = g'_i \circ \mathcal{F}[\tau_i](h, \text{id}) \\
\Leftrightarrow & h \circ g_i \circ \mathcal{F}[\tau_i](\mathcal{I}\mathcal{N}\mathcal{V}(h), \text{id}) = g'_i \circ \mathcal{F}[\tau_i](h, \text{id}) \circ \mathcal{F}[\tau_i](\mathcal{I}\mathcal{N}\mathcal{V}(h), \text{id}) \\
\Leftrightarrow & h \circ g_i \circ \mathcal{F}[\tau_i](\mathcal{I}\mathcal{N}\mathcal{V}(h), \text{id}) = g'_i \circ \mathcal{F}[\tau_i](h \circ \mathcal{I}\mathcal{N}\mathcal{V}(h), \text{id}) && \text{from Eq. 5} \\
\Leftrightarrow & h \circ g_i \circ \mathcal{F}[\tau_i](\mathcal{I}\mathcal{N}\mathcal{V}(h), \text{id}) = g'_i \circ \mathcal{F}[\tau_i](\text{id}, \text{id}) && \text{since } h \circ \mathcal{I}\mathcal{N}\mathcal{V}(h) = \text{id} \\
\Leftrightarrow & h \circ g_i \circ \mathcal{F}[\tau_i](\mathcal{I}\mathcal{N}\mathcal{V}(h), \text{id}) = g'_i \circ \text{id} && \text{from Eq. 4} \\
\Leftrightarrow & h \circ g_i \circ \mathcal{F}[\tau_i](\mathcal{I}\mathcal{N}\mathcal{V}(h), \text{id}) = g'_i
\end{aligned}$$

Thus, we derive the following algorithm from Eq. 6:

**Algorithm 1 (Fusion Algorithm)** *The fusion  $h(f(\bar{x}))$  is achieved by the following rewrite rules:*

$$\forall h, g_i, x_i : h(\mathcal{SK}_f(\bar{g}) \bar{x}) \rightarrow \mathcal{SK}_f(\bar{g}') \bar{x} \quad \text{where } g'_i = h \circ g_i \circ \mathcal{F}[\tau_i](\mathcal{I}\mathcal{N}\mathcal{V}(h), \text{id}) \quad (7)$$

$$\forall x : h(\mathcal{I}\mathcal{N}\mathcal{V}(h) x) \rightarrow x \quad (8)$$

Eq. 7 fuses  $h$  with  $f$  into a function that has the same recursive skeleton as  $f$ . The components  $g'_i$  of the resulting skeleton are derived by fusing  $h$ ,  $g_i$ , and  $\mathcal{F}[\tau_i](\mathcal{I}\mathcal{N}\mathcal{V}(h), \text{id})$ . Since  $\mathcal{I}\mathcal{N}\mathcal{V}(h)$  is unknown, we only fuse  $h$  with  $g_i$ . This fusion can be achieved by using Eq. 7, Eq. 8, as well some standard partial evaluation techniques (such as applying a function to a value construction). At the end, we hope to derive terms of the form  $h \circ \mathcal{I}\mathcal{N}\mathcal{V}(h)$  only, which cancel  $\mathcal{I}\mathcal{N}\mathcal{V}(h)$ . Otherwise, if there is a term  $\mathcal{I}\mathcal{N}\mathcal{V}(h)$  in the program that is not cancelled out, then the fusion algorithm fails and the fusion  $h(f(\bar{x}))$  is not performed.

As an example, we will fuse `nth (zip(x,y)) n`, where

$$\begin{aligned}
\text{nth } [ ] \text{ n} &= \text{error}() \\
\text{nth } (\text{a:x}) \text{ n} &= \text{if } n=0 \text{ then a else nth x (n-1)}
\end{aligned}$$

Using the fusion algorithm we get:

$$\begin{aligned}
\text{nth } (\text{zip}(x,y)) \text{ n} &= \text{nth } (\text{zip}'(\lambda(\text{a,b,r}). (\text{a,b}):r, \lambda(). [ ])(x,y)) \text{ n} \\
&= \text{zip}'(\lambda(\text{a,b,r}). \text{nth } ((\text{a,b}):(\mathcal{I}\mathcal{N}\mathcal{V}(\text{nth}) r)), \lambda(). \text{nth } [ ])(x,y) \text{ n} \\
&= \text{zip}'(\lambda(\text{a,b,r}). \lambda n. \text{if } n=0 \text{ then } (\text{a,b}) \text{ else nth}(\mathcal{I}\mathcal{N}\mathcal{V}(\text{nth}) r) (n-1), \\
&\quad \lambda(). \lambda n. \text{error}())(x,y) \text{ n} \\
&= \text{zip}'(\lambda(\text{a,b,r}). \lambda n. \text{if } n=0 \text{ then } (\text{a,b}) \text{ else } r(n-1), \lambda(). \lambda n. \text{error}())(x,y) \text{ n}
\end{aligned}$$

If we unfold `zip'`, we get:

$$\begin{aligned}
f(\text{a:x,b:y}) \text{ n} &= \text{if } n=0 \text{ then } (\text{a,b}) \text{ else } f(x,y) (n-1) \\
f \_ \text{ n} &= \text{error}()
\end{aligned}$$

But there are examples in which  $\mathcal{I}\mathcal{N}\mathcal{V}(h)$  cannot not cancelled out, such as in the case of the length of the quadratic reverse. Using the fusion algorithm we get:

$$\begin{aligned}
\text{len}(\text{rev}(x)) &= \text{len}(\text{rev}'(\lambda(). [ ], \lambda(\text{a,x,r}). \text{append } r [\text{a}] ) x) \\
&= \text{rev}'(\lambda(). \text{len}([ ]), \lambda(\text{a,x,r}). \text{len}(\text{append } (\mathcal{I}\mathcal{N}\mathcal{V}(\text{len}) r) [\text{a}]) ) x \\
&= \text{rev}'(\lambda(). 0, \lambda(\text{a,x,r}). \text{len}(\text{append } (\mathcal{I}\mathcal{N}\mathcal{V}(\text{len}) r) [\text{a}]) ) x
\end{aligned}$$

The skeleton of `append` is:

*Proof:* Since all type variables  $\alpha_i$  are positive in  $\tau$ , none of the  $g_i$ s in Th. 4 is used. Therefore, each  $g_i$  can be replaced by an arbitrary function, including the identity function itself.  $\square$

For example, the bifunctor for  $\tau = \text{int} \rightarrow (\alpha \times (\text{int} \rightarrow \text{list}(\alpha)))$  is:

$$\mathcal{F}[\tau](f_\alpha, g_\alpha) = \lambda h. (f_\alpha \times (\lambda k. \text{map}(f_\alpha) \circ k)) \circ h$$

We have  $\mathcal{P}[\text{int} \rightarrow \text{list}(\alpha)](g', g) \Leftrightarrow g' = \text{map}(f_\alpha) \circ g$  and

$$\begin{aligned} \mathcal{P}[\tau](h', h) &\Leftrightarrow x = y \Rightarrow (\pi_1(h'(x)) = f_\alpha(\pi_1(h(y)))) \wedge (\pi_2(h'(x)) = \text{map}(f_\alpha) \circ (\pi_2(h(y)))) \\ &\Leftrightarrow h'(x) = (f_\alpha \times (\lambda k. \text{map}(f_\alpha) \circ k))(h(x)) \end{aligned}$$

which is equivalent to  $h' = \mathcal{F}[\tau](f_\alpha, g_\alpha)h$ .

### 4.3 The Inverse Trick

The problem of fusing two recursive functions  $f$  and  $h$  in  $h(f(x))$  is to derive a new recursive function with the same functionality as  $h(f(x))$ . This is not always possible. Typically,  $f$  produces an intermediate data structure which is consumed by  $h$ . When these two functions are fused, this data structure is not generated. In our framework, the fusion  $h(f(x))$  is achieved by fusing  $h(\mathcal{SK}_f(\bar{g})x)$ , since  $f = \mathcal{SK}_f(\bar{g})x$ , for some functions  $g_i$ . The law for this fusion is derived directly from the parametricity theorem of  $\mathcal{SK}_f$  given below.

We have seen that even with all the extensions described in Section 4.1, the type of  $\mathcal{SK}_f$  has the following form:

$$\forall \alpha. (\tau_1 \rightarrow \alpha) \times \cdots \times (\tau_m \rightarrow \alpha) \rightarrow t_1 \rightarrow \cdots \rightarrow t_n \rightarrow \alpha$$

Note that, type  $\tau_i$  does not contain any negative instances of  $\alpha$ . As we have seen from Eq. 3, the parametricity theorem for this type is:

$$\forall h, g_i, g'_i, x_i : \bigwedge_i h \circ g_i = g'_i \circ \mathcal{F}[\tau_i](h, \text{id}) \Rightarrow h(\mathcal{SK}_f(\bar{g})\bar{x}) = \mathcal{SK}_f(\bar{g}')\bar{x} \quad (6)$$

This gives us a law for fusing any function  $h$  with  $f$ :  $h$  is given, each  $g_i$  is derived directly from the definition of  $f$ , and each  $\mathcal{F}[\tau_i]$  is derived from Definition 1. The only functions that need to be computed are the  $g'_i$  functions.

In a previous work [11], we describe a method for solving a similar set of equations for fusing folds. It relies on the fact that for each function  $h$  there exists a function  $\mathcal{INV}(h)$  such that  $\forall x \in \text{range}(h) : h(\mathcal{INV}(h)x) = x$ , i.e.,  $\mathcal{INV}(h)$  is a right inverse of  $h$ . To see why this is true, consider all the values  $a_1, \dots, a_n, n > 0$  that satisfy  $h(a_1) = \cdots = h(a_n) = b$  for some value  $b \in \text{range}(h)$ . Then,  $\mathcal{INV}(h)b$  is defined to be one of these  $a_i$ . This inverse trick, which is described in [11] and is also used in this paper, does not actually derive a function  $\mathcal{INV}(h)$  from  $h$ . It only assumes that such function exists and uses the property  $h \circ \mathcal{INV}(h) = \text{id}$  to eliminate it.

where  $\neg(+)= -$  and  $\neg(-)= +$ . For example,

$$((\alpha \rightarrow \beta) \rightarrow \gamma \rightarrow \delta)^+ = (\alpha^+ \rightarrow \beta^-) \rightarrow \gamma^- \rightarrow \delta^+$$

A type variable  $\alpha$  is positive (resp., negative) in a type  $\tau$  if all occurrences of  $\alpha$  in  $\tau^+$  are  $\alpha^+$  (resp.,  $\alpha^-$ ). For example,  $\alpha$  and  $\delta$  are positive in the type  $(\alpha \rightarrow \beta) \rightarrow \gamma \rightarrow \delta$ , while  $\beta$  and  $\gamma$  are negative. The following theorem can be proved using structural induction over types:

**Theorem 3** *If  $f_i : \alpha_i^- \rightarrow \alpha_i^+$  and  $g_i : \alpha_i^+ \rightarrow \alpha_i^-$ , then  $\mathcal{F}[\tau](\bar{f}, \bar{g}) : \tau^- \rightarrow \tau^+$ .*

The following theorem expresses the parametricity theorem of a type in terms of the type's bifunctor (the proof is given in the appendix):

**Theorem 4** *For any type  $\forall \alpha_1, \dots, \alpha_n : \tau$ ,  $n \geq 0$  we have:*

$$\forall f_{\alpha_i}, g_{\alpha_i}, x, y : \bigwedge_i f_{\alpha_i} \circ g_{\alpha_i} = \text{id} \Rightarrow (\mathcal{P}[\tau](x, y) \Leftrightarrow x = \mathcal{F}[\tau](\bar{f}, \bar{g}) y)$$

where  $\bar{f} = f_{\alpha_1}, \dots, f_{\alpha_n}$  and  $\bar{g} = g_{\alpha_1}, \dots, g_{\alpha_n}$ .

For example, the bifunctor for  $\tau = (\alpha \rightarrow \beta) \rightarrow \gamma$  is:

$$\mathcal{F}[\tau](f_\alpha, f_\beta, f_\gamma, g_\alpha, g_\beta, g_\gamma) = \lambda h. f_\gamma \circ h \circ (\lambda k. g_\beta \circ k \circ f_\alpha)$$

and the parametricity theorem is:

$$\begin{aligned} \mathcal{P}[\tau](x, y) &\Leftrightarrow \forall m, n : m \circ f_\alpha = f_\beta \circ n \Rightarrow x(m) = f_\gamma(y(n)) \\ &\Leftrightarrow \forall m, n : m \circ f_\alpha \circ g_\alpha = f_\beta \circ n \circ g_\alpha \Rightarrow x(m) = f_\gamma(y(n)) & (\epsilon 1) \\ &\Leftrightarrow \forall m, n : m = f_\beta \circ n \circ g_\alpha \Rightarrow x(m) = f_\gamma(y(n)) \\ &\Leftrightarrow \forall n : x(f_\beta \circ n \circ g_\alpha) = f_\gamma(y(n)) \\ &\Leftrightarrow x \circ (\lambda n. f_\beta \circ n \circ g_\alpha) = f_\gamma \circ y \\ &\Leftrightarrow x \circ (\lambda n. f_\beta \circ n \circ g_\alpha) \circ (\lambda n. g_\beta \circ n \circ f_\alpha) = f_\gamma \circ y \circ (\lambda n. g_\beta \circ n \circ f_\alpha) & (\epsilon 2) \\ &\Leftrightarrow x \circ (\lambda n. f_\beta \circ g_\beta \circ n \circ f_\alpha \circ g_\alpha) = f_\gamma \circ y \circ (\lambda n. g_\beta \circ n \circ f_\alpha) \\ &\Leftrightarrow x \circ (\lambda n. n) = f_\gamma \circ y \circ (\lambda n. g_\beta \circ n \circ f_\alpha) \\ &\Leftrightarrow x = \mathcal{F}[\tau](f_\alpha, f_\beta, f_\gamma, g_\alpha, g_\beta, g_\gamma) y \end{aligned}$$

The equivalences  $(\epsilon 1)$  and  $(\epsilon 2)$  are based on the fact that  $\text{range}(h) = \text{domain}(f) \Rightarrow (f = g \Leftrightarrow f \circ h = g \circ h)$ .

**Corollary 1** *Let  $\forall \alpha_1, \dots, \alpha_n : \tau$ ,  $n \geq 0$  be a type whose type variables are positive in  $\tau$ , then:*

$$\forall f_{\alpha_i}, x, y : \mathcal{P}[\tau](x, y) \Leftrightarrow x = \mathcal{F}[\tau](\bar{f}, \bar{\text{id}}) y$$

where  $\bar{f} = f_{\alpha_1}, \dots, f_{\alpha_n}$ .

## 4.2 Expressing the Parametricity Theorem using Bifunctors

Most useful parametricity theorems are quite complex to use directly for program fusion since they typically include both a premise and a consequence. Thus, to fuse programs we need to synthesize a number of functions that satisfy the premise, which is hard in general. In this subsection we describe a theory for solving parametricity theorems effectively. A parametricity theorem relates two values of the same type. Here we prove that one value can be calculated in terms of the other under certain conditions. This calculation is based on the notion of the *bifunctor*. We will show in the next subsection that this expression of the parametricity theorem using bifunctors, even though it too has preconditions to be satisfied, these preconditions are eliminated if we use the *inverse trick* (to be described shortly), thus, calculating a solution directly.

A bifunctor is a generalization of a functor. In contrast to regular functors, bifunctors can capture types with contravariant type variables [3]. Here we define bifunctors only on a type whose universal quantifications appear in the outermost level, i.e., for types of the form  $\forall \alpha_1, \dots, \alpha_n : \tau$  for  $n \geq 0$ .

**Definition 1 (Bifunctor)** *Let  $\forall \alpha_1, \dots, \alpha_n : \tau$ ,  $n \geq 0$  be a polymorphic type and let  $\vec{f} = f_1, \dots, f_n$  and  $\vec{g} = g_1, \dots, g_n$ , where each  $f_i/g_i$  pair is associated with the type variable  $\alpha_i$ . The bifunctor  $\mathcal{F}[\tau](\vec{f}, \vec{g})$  is defined as follows:*

$$\begin{aligned} \mathcal{F}[\text{basic}](\vec{f}, \vec{g}) &\rightarrow \text{id} \\ \mathcal{F}[\alpha_i](\vec{f}, \vec{g}) &\rightarrow f_i \\ \mathcal{F}[\tau_1 \times \tau_2](\vec{f}, \vec{g}) &\rightarrow \mathcal{F}[\tau_1](\vec{f}, \vec{g}) \times \mathcal{F}[\tau_2](\vec{f}, \vec{g}) \\ \mathcal{F}[\tau_1 \rightarrow \tau_2](\vec{f}, \vec{g}) &\rightarrow \lambda h. \mathcal{F}[\tau_2](\vec{f}, \vec{g}) \circ h \circ \mathcal{F}[\tau_1](\vec{g}, \vec{f}) \\ \mathcal{F}[T(\tau)](\vec{f}, \vec{g}) &\rightarrow \text{map}^T(\mathcal{F}[\tau](\vec{f}, \vec{g})) \end{aligned}$$

where the product of functions is defined by  $(f \times g)(x, y) = (f x, g y)$ . For example,

$$\mathcal{F}[(\alpha \rightarrow \alpha) \rightarrow \alpha](f, g) = \lambda h. f \circ h \circ (\lambda k. g \circ k \circ f)$$

It is easy to prove the following theorem:

**Theorem 2** *A bifunctor  $\mathcal{F}[\tau](\vec{f}, \vec{g})$  is a functor that is covariant over  $f_i$  and contravariant over  $g_i$ . That is,*

$$\mathcal{F}[\tau](\vec{\text{id}}, \vec{\text{id}}) = \text{id} \tag{4}$$

$$\mathcal{F}[\tau](\vec{f}, \vec{g}) \circ \mathcal{F}[\tau](\vec{f}', \vec{g}') = \mathcal{F}[\tau](\vec{f} \circ \vec{f}', \vec{g}' \circ \vec{g}) \tag{5}$$

We annotate type variables by a sign  $s \in \{+, -\}$  as follows:

$$\begin{aligned} \text{basic}^s &= \text{basic} \\ \alpha^s &= \alpha^s \\ (\tau_1 \times \tau_2)^s &= \tau_1^s \times \tau_2^s \\ (\tau_1 \rightarrow \tau_2)^s &= \tau_1^{-s} \rightarrow \tau_2^s \\ (T(\tau))^s &= T(\tau^s) \end{aligned}$$



The type of **rev'** is:

$$\forall \alpha, \beta. ((\rightarrow \alpha) \times (\beta \times \text{list}(\beta) \times \alpha \rightarrow \alpha) \rightarrow \text{list}(\beta) \rightarrow \alpha)$$

which satisfies the following parametricity theorem (for  $\beta = ()$ ):

$$\forall f_\alpha, n, c, n', c', x : f_\alpha \circ n = n' \circ \text{id} \wedge f_\alpha \circ c = c' \circ (\text{id} \times \text{id} \times f_\alpha) \Rightarrow f_\alpha(\text{rev}'(n, c)x) = \text{rev}'(n', c')x$$

The above algorithm can be easily extended to allow recursive calls  $f a_1 \cdots a_n$  in which some variables in  $a_i$  are bound in an outer case statement. In that case, we would need more extra parameters for  $\mathcal{SK}_f$ ; one for each case branch. A more substantial extension can be achieved by permitting  $a_i$  to contain variables that do not appear in the function patterns or in the outer case statements. In that case we do not abstract the function call but instead we construct a lambda expression that captures all these free variables. If all arguments to the recursive calls are free variables, our method deteriorates to the unfold-simplify-fold law (Eq. 2), because the recursive skeleton becomes the  $Y$  combinator. Since this is undesirable, we try to abstract as many arguments of the recursive calls as possible.

For example, consider the following program that computes the map over bushes:

$$\begin{aligned} \text{mapB}(f)(\text{Leaf } x) &= \text{Leaf}(f x) \\ \text{mapB}(f)(\text{Branch } r) &= \text{Branch}(\lambda z. \text{mapB}(f) z) r \end{aligned}$$

where, **Leaf** and **Branch** are the value constructors of **Bush**:

$$\text{data Bush}(\alpha) = \text{Leaf } \alpha \mid \text{Branch list}(\text{Bush}(\alpha))$$

Notice that the variable **z** in the second equation is not bound in the pattern of the equation. Thus, in this case we do not abstract the recursive call alone, but a lambda abstraction that contains the recursive call:

$$\begin{aligned} \text{mapB}'(l, b)(f)(\text{Leaf } x) &= l(f, x) \\ \text{mapB}'(l, b)(f)(\text{Branch } r) &= b(r, \lambda z. \text{mapB}'(l, b)(f) z) \end{aligned}$$

Function **mapB** is defined in terms of **mapB'**:

$$\text{mapB} = \text{mapB}'(\lambda(f, x). \text{Leaf}(f x), \lambda(r, g). \text{Branch}(\text{map}(g) r) )$$

The type of **mapB'** is:

$$\forall \alpha, \beta, \gamma. (\alpha \times \beta \rightarrow \gamma) \times ((\text{list}(\text{Bush}(\beta)) \times (\text{Bush}(\beta) \rightarrow \gamma)) \rightarrow \gamma) \rightarrow \alpha \rightarrow \text{Bush}(\beta) \rightarrow \gamma$$

The parametricity theorem for  $\alpha = \beta = ()$  is:

$$\begin{aligned} \forall f_\gamma, l, b, b', f, x : f_\gamma \circ b = b' \circ (\text{id} \times (\lambda g. f_\gamma \circ g)) \\ \Rightarrow f_\gamma(\text{mapB}'(l, b)(f) x) = \text{mapB}'(f_\gamma \circ l, b')(f) x \end{aligned}$$

## 4.1 Extracting the Recursive Skeleton of a Function

This section presents an algorithm for extracting the recursive skeleton of a function  $f$  of type  $t_1 \rightarrow t_2 \rightarrow \dots \rightarrow t_n$ . It can handle any function  $f$  defined in terms of  $m$  recursive equations:

$$\begin{aligned} f p_{1,1} \dots p_{1,n} &= \epsilon_1(\overline{v_1}, \overline{f \epsilon_{1,1} \dots \epsilon_{1,n}}) \\ &\vdots \\ f p_{m,1} \dots p_{m,n} &= \epsilon_m(\overline{v_m}, \overline{f \epsilon_{m,1} \dots \epsilon_{m,n}}) \end{aligned}$$

where  $p_{i,j}$  is a pattern and  $\epsilon_i$  is a function in which all variables  $\overline{v_i}$  that appear in the patterns and all the recursive calls  $f \epsilon_{i,1} \dots \epsilon_{i,n}$  have been abstracted out of the body of  $\epsilon_i$ . That is, the body of  $\epsilon_i$  does not contain any free variable or any recursive calls to  $f$ . In that case, the skeleton of  $f$  is  $\mathcal{SK}_f$ , whose  $i$ th inductive equation is defined as follows:

$$\mathcal{SK}_f(g_1, \dots, g_m) p_{i,1} \dots p_{i,n} = g_i(\overline{v_i}, \overline{\mathcal{SK}_f(g_1, \dots, g_m) \epsilon_{i,1} \dots \epsilon_{i,n}})$$

Therefore,  $f = \mathcal{SK}_f(\epsilon_1, \dots, \epsilon_m)$ .

The type of  $\mathcal{SK}_f$  is (when we set all but the output type variable to  $()$ ):

$$\forall \alpha. (\tau_1 \rightarrow \alpha) \times \dots \times (\tau_m \rightarrow \alpha) \rightarrow t_1 \rightarrow \dots \rightarrow t_n \rightarrow \alpha$$

The parametricity theorem for this type is:

$$\forall f_\alpha, g_i, g'_i, x_i: \bigwedge_i f_\alpha \circ g_i = g'_i \circ \mathcal{F}[\tau_i](f_\alpha, \text{id}) \Rightarrow f_\alpha(\mathcal{SK}_f(\overline{g}) \overline{x}) = \mathcal{SK}_f(\overline{g}') \overline{x} \quad (3)$$

*Proof:* The parametricity theorem for the type of  $\mathcal{SK}_f$  is  $\bigwedge_i \mathcal{P}[\tau_i \rightarrow \alpha](g'_i, g_i) \Rightarrow f_\alpha(\mathcal{SK}_f(\overline{g}) \overline{x}) = \mathcal{SK}_f(\overline{g}') \overline{x}$ . We have:

$$\begin{aligned} \mathcal{P}[\tau_i \rightarrow \alpha](g'_i, g_i) &\Leftrightarrow \mathcal{P}[\tau_i](x, y) \Rightarrow g'_i(x) = f_\alpha(g_i(y)) && \text{from Th. 1} \\ &\Leftrightarrow x = \mathcal{F}[\tau_i](f_\alpha, \text{id}) y \Rightarrow g'_i(x) = f_\alpha(g_i(y)) && \text{from Th. 1} \\ &\Leftrightarrow g'_i(\mathcal{F}[\tau_i](f_\alpha, \text{id}) y) = f_\alpha(g_i(y)) && \square \end{aligned}$$

For example, consider the list reverse function:

$$\begin{aligned} \text{rev } [] &= [] \\ \text{rev } (a:x) &= \text{append } (\text{rev } x) [a] \end{aligned}$$

Its recursive skeleton  $\mathcal{SK}_{\text{rev}} = \text{rev}'$  is straightforward:

$$\begin{aligned} \text{rev}'(n,c) [] &= n() \\ \text{rev}'(n,c) (a:x) &= c(a,x, \text{rev}'(n,c) x) \end{aligned}$$

which is actually equivalent to the list primitive recursion. Function  $\text{rev}$  can be expressed in terms of  $\text{rev}'$ :

$$\text{rev} = \text{rev}'(\lambda(). [], \lambda(a,x,r). \text{append } r [a])$$

Following the terminology in [10],  $\mathbf{zip}'$  can be viewed as a worker and the above definition of  $\mathbf{zip}$  as a wrapper. The parametricity theorem for  $\mathbf{zip}'$  with  $\alpha = \beta = ()$  is:

$$\forall f_\gamma, c, n, c' : f_\gamma \circ c = c' \circ (\text{id} \times \text{id} \times f_\gamma) \Rightarrow f_\gamma \circ \mathbf{zip}'(c, n) = \mathbf{zip}'(c', f_\gamma \circ n)$$

Suppose now that we want to perform the program fusion  $\mathbf{len}(\mathbf{zip}(x, y))$ . We can achieve this fusion by unwrapping  $\mathbf{zip}$  and by using the  $\mathbf{zip}'$  fusion law for  $f_\gamma = \mathbf{len}$  and  $\mathbf{c}(a, b, r) = (a, b) : r$ :

$$\begin{aligned} \mathbf{len}(\mathbf{zip}(x, y)) &= \mathbf{len}(\mathbf{zip}'(\lambda(a, b, r). (a, b) : r, \lambda(). [])(x, y)) \\ &= \mathbf{zip}'(\mathbf{c}', \lambda(). \mathbf{len} [ ])(x, y) \\ &= \mathbf{zip}'(\mathbf{c}', \lambda(). 0)(x, y) \end{aligned}$$

The premise of the  $\mathbf{zip}'$  fusion law gives us a value for  $\mathbf{c}'$ :

$$\begin{aligned} \mathbf{c}'(x, y, \mathbf{len} z) &= \mathbf{len}(c(x, y, z)) \\ &= \mathbf{len}((x, z) : z) \\ &= 1 + (\mathbf{len} z) \end{aligned}$$

Therefore, if we generalize the term  $\mathbf{len} z$  to a variable  $w$ , we get  $\mathbf{c}'(x, y, w) = 1 + w$ . That is,

$$\mathbf{len}(\mathbf{zip}(x, y)) = \mathbf{zip}'(\lambda(x, y, w). 1 + w, \lambda(). 0)(x, y)$$

Finally, if we unfold  $\mathbf{zip}'$  and set  $\mathbf{f} = \mathbf{zip}'(\lambda(x, y, w). 1 + w, \lambda(). 0)$ , we get:

$$\begin{aligned} \mathbf{f}(a : x, b : y) &= 1 + (\mathbf{f}(x, y)) \\ \mathbf{f} \_ &= 0 \end{aligned}$$

## 4 The Fusion Algorithm

Program fusion in our framework is performed in four steps:

- Given a function  $f$ , generate a sufficiently polymorphic function  $\mathcal{SK}_f$ , called the *recursive skeleton* of  $f$ , which captures the recursion scheme of  $f$ ;
- Redefine  $f$  as a wrapper of  $\mathcal{SK}_f$ , i.e.,  $f = \mathcal{SK}_f(e_1, \dots, e_n)$ , for some expressions  $e_i$ ;
- Generate the parametricity theorem for the type of  $\mathcal{SK}_f$ , with  $\alpha = ()$  for any type variable  $\alpha$  other than the type variable of the output;
- Whenever there is an application  $g(f e)$  in a program, unwrap  $f$  into an  $\mathcal{SK}_f$  call and use the parametricity theorem to fuse the application.

As it will be explained shortly, our fusion algorithm is based on a simple term rewriting system and does not involve theorem proving or program synthesis.

The above steps are described in greater detail below.

If  $t_n = \beta$ , then the parametricity theorem for  $f$  becomes:

$$\begin{aligned} \mathcal{P}[[t_1]](x_1, y_1) &\Rightarrow (\mathcal{P}[[t_2]](x_2, y_2) \Rightarrow \cdots (\mathcal{P}[[t_{n-1}]](x_{n-1}, y_{n-1}) \\ &\Rightarrow f_\beta(f x_1 x_2 \dots x_n) = f y_1 y_2 \dots y_n)) \end{aligned}$$

which gives us a fusion law for fusing any function  $f_\beta$  with  $f$ .

Our previous analysis indicates that we should transform a function into a traversal scheme in such a way that the scheme's output type be completely parametric (i.e., a type variable). Having done this, we can easily generate the parametricity theorem for the scheme and use it to perform program fusion in the same way we use the fold fusion law to fuse a function with a fold.

The following is an example of a function that does not have a direct representation as a fold. We will transform it to get a completely polymorphic output. This function is **zip**<sup>1</sup>:

$$\begin{aligned} \mathbf{zip}(a:x, b:y) &= (a, b):(\mathbf{zip}(x, y)) \\ \mathbf{zip} \_ &= [ ] \end{aligned}$$

Even though, its type is polymorphic:

$$\forall \alpha, \beta. (\text{list}(\alpha) \times \text{list}(\beta)) \rightarrow \text{list}(\alpha \times \beta)$$

its parametricity theorem is a simple natural transformation:

$$\forall f_\alpha, f_\beta : \text{map}(f_\alpha \times f_\beta) \circ \mathbf{zip} = \mathbf{zip} \circ (\text{map}(f_\alpha) \times \text{map}(f_\beta))$$

Notice that the output type of **zip** is not a type variable. Thus, the above equation cannot be used as is for fusing any function  $g$  with **zip**. To make the output type of **zip** a type variable, we need to generalize both the inductive equations of **zip**. First observe that the second equation returns `[ ]`; this should be replaced by an extra parameter, **n**, of **zip**. Finally, the first equation returns a list construction; this too should be abstracted into another extra parameter, **c**. The transformed function **zip'** is now:

$$\begin{aligned} \mathbf{zip}'(c, n)(a:x, b:y) &= c(a, b, \mathbf{zip}'(c, n)(x, y)) \\ \mathbf{zip}'(c, n) \_ &= n() \end{aligned}$$

which has a sufficiently polymorphic type:

$$\forall \alpha, \beta, \gamma. ((\alpha \times \beta \times \gamma \rightarrow \gamma) \times ((\ ) \rightarrow \gamma)) \rightarrow (\text{list}(\alpha) \times \text{list}(\beta)) \rightarrow \gamma$$

since its output type is the type variable  $\gamma$ . Function **zip** can be computed in terms of **zip'**:

$$\mathbf{zip} = \mathbf{zip}'(\lambda(a, b, r). (a, b):r, \lambda(). [ ])$$

---

<sup>1</sup>Function **zip** can be expressed as a second-order fold that traverses one of the **zip** arguments and deconstructs the other argument during the traversal. This results into an asymmetry: the fold fusion law can only be used for fusing one argument only. An alternative, symmetric, definition of **zip** is given elsewhere [4] but it requires a more general traversal scheme than fold.

where  $\otimes$  can be calculated from the premise of Eq. 1:

$$\begin{aligned} \mathbf{a} \otimes (\mathbf{len} \mathbf{r}) &= \mathbf{len}((1+\mathbf{a}):r) \\ &= 1+(\mathbf{len} \mathbf{r}) \end{aligned}$$

If we substitute  $(\mathbf{len} \mathbf{r})$  for  $\mathbf{s}$ , we get  $\mathbf{a} \otimes \mathbf{s}=1+\mathbf{s}$  and, finally,

$$\mathbf{len}(\mathbf{inc} \mathbf{x}) = \mathbf{fold}(\lambda \mathbf{a}. \lambda \mathbf{s}. 1+\mathbf{s}) \mathbf{0} \mathbf{x}$$

The resulting fold does not create the intermediate list of the original program.

Unfortunately, not all functions can be expressed as folds. Even though there are methods for translating a number of recursive functions into folds [8], these methods usually fail for complex functions. One solution to this problem is to use a list traversal scheme that is more flexible and maybe more expressive than **fold**. Some researchers have suggested hylomorphisms as a possible solution [9]. Recent work [6, 7] has been shown that most functional programs can be directly translated into hylomorphisms. It has yet to be shown the effectiveness of program fusion under this representation.

In this paper we propose an alternative solution to the above problem: instead of trying to make some recursive function fit the recursion pattern of a particular fixed traversal scheme, such as **fold**, we generate a traversal scheme that is individually tailored to this particular function. This scheme may not be useful for any other function. This traversal scheme is ‘polymorphic enough’ to satisfy a useful parametricity theorem.

We can make a function more polymorphic (i.e., with more type variables) by abstracting pieces of its code into some extra function arguments. But when a function becomes ‘polymorphic enough’? To answer this question we consider the parametricity theorem for **fold** (Equation 1). This theorem is useful because it has the conclusion:

$$f_\beta(\mathbf{fold}(\oplus) x y) = \mathbf{fold}(\otimes)(f_\beta x) y$$

The left part is the composition of any function  $f_\beta$  (which corresponds to the type variable  $\beta$ ) with a **fold**. The right part is another fold whose arguments can be calculated from the arguments of the first fold by using the equalities in the premise of the theorem. That way we can fuse any function  $f_\beta$  with a fold yielding another fold.

Given how the parametricity theorem should look like to be useful for fusion, we can easily guess how the type of a traversal scheme should look like to generate such a theorem: if the type of a traversal scheme  $f$  has the form  $t_1 \rightarrow t_2 \rightarrow \dots \rightarrow t_{n-1} \rightarrow t_n$ , then  $t_n$  should be a type variable, say  $\beta$ . To see why, we derive the parametricity theorem for  $f$  from Theorem 1 (universal quantifications are omitted here):

$$\begin{aligned} &\mathcal{P}[[t_1 \rightarrow t_2 \rightarrow \dots \rightarrow t_n]](f, f) \\ &= \mathcal{P}[[t_1]](x_1, y_1) \Rightarrow \mathcal{P}[[t_2 \rightarrow \dots \rightarrow t_n]](f x_1, f y_1) \\ &= \mathcal{P}[[t_1]](x_1, y_1) \Rightarrow (\mathcal{P}[[t_2]](x_2, y_2) \Rightarrow \mathcal{P}[[t_3 \rightarrow \dots \rightarrow t_n]](f x_1 x_2, f y_1 y_2)) \\ &= \mathcal{P}[[t_1]](x_1, y_1) \Rightarrow (\mathcal{P}[[t_2]](x_2, y_2) \Rightarrow \dots (\mathcal{P}[[t_{n-1}]](x_{n-1}, y_{n-1}) \\ &\quad \Rightarrow \mathcal{P}[[t_n]](f x_1 x_2 \dots x_n, f y_1 y_2 \dots y_n))) \end{aligned}$$

theorem for  $t(\alpha, \beta, \dots)$  and then setting  $f_\alpha = \text{id}$  is the same as finding the parametricity theorem for the type  $t(\alpha, \beta, \dots)$ , where  $()$  is the unit type.

Another example is the monad extension operator:

$$\text{ext}^T : \forall \alpha, \beta. (\alpha \rightarrow T(\beta)) \rightarrow T(\alpha) \rightarrow T(\beta)$$

that satisfies for  $\beta = ()$ :

$$\forall f_\alpha, g : \text{ext}^T(g) \circ \text{map}^T(f_\alpha) = \text{ext}^T(g \circ f_\alpha)$$

The  $Y$  combinator for functions is defined as  $Y(f) x = f(Y(f)) x$  and has type:

$$\forall \alpha, \beta. ((\alpha \rightarrow \beta) \rightarrow \alpha \rightarrow \beta) \rightarrow \alpha \rightarrow \beta$$

The parametricity theorem for this type with  $\alpha = ()$  is:

$$\forall f_\beta, f, f', g : f_\beta \circ (f g) = f'(f_\beta \circ g) \Rightarrow f_\beta \circ (Y f) = Y f' \quad (2)$$

which is actually the unfold-simplify-fold law [1] since it fuses any function  $f_\beta$  with a recursive function  $Y f$  yielding another recursive function  $Y f'$  (under certain conditions).

### 3 Using the Parametricity Theorem for Program Fusion

Consider the following non-polymorphic function of type  $\text{list}(\text{int}) \rightarrow \text{list}(\text{int})$  expressed in Haskell:

$$\begin{aligned} \text{inc } [] &= [] \\ \text{inc } (a:x) &= (1+a):(\text{inc } x) \end{aligned}$$

The parametricity theorem of a non-polymorphic type is always a tautology, which is of no use. Luckily, **inc** happens to be a **fold**, since a **fold** has a similar pattern of recursion:

$$\begin{aligned} \text{fold } f \ b \ [] &= b \\ \text{fold } f \ b \ (a:x) &= f \ a \ (\text{fold } f \ b \ x) \end{aligned}$$

In particular,  $\text{inc } x = \text{fold}(\lambda a. \lambda r. (1+a):r) [] x$ . This is quite useful because we know that **fold** satisfies a powerful parametricity theorem (Eq. 1). In fact, we have shown elsewhere [11] that there is an automated method for fusing a function composed with a fold. Suppose, for example that we want to fuse  $\text{len}(\text{inc } x)$ , where **len** computes the length of a list, so that the intermediate list produced by **inc** and consumed by **len** is eliminated. Function **len** is defined as follows:

$$\begin{aligned} \text{len } [] &= 0 \\ \text{len } (a:x) &= 1+(\text{len } x) \end{aligned}$$

Thus,  $\text{len}(\text{inc } x)$  can be calculated from Eq. 1, where  $f_\beta = \text{len}$ ,  $a \oplus r = (1+a):r$ ,  $x = []$ , and  $y = x$ . From the conclusion of Eq. 1 we have:

$$\begin{aligned} \text{len}(\text{inc } x) &= \text{len}(\text{fold}(\lambda a. \lambda r. (1+a):r) [] x) \\ &= \text{fold}(\otimes) (\text{len } []) x \\ &= \text{fold}(\otimes) 0 x \end{aligned}$$

type  $T(\alpha)$  is a functor that maps a function  $f$  into  $\text{map}^T(f)$ , where  $\text{map}^T : (\alpha \rightarrow \beta) \rightarrow T(\alpha) \rightarrow T(\beta)$  is the map function for type  $T$ . Since the composition of functors is also a functor, the type  $\text{list}(\text{list}(\alpha))$  is a functor that maps  $f$  into  $\text{map}(\text{map}(f))$ . The parametricity theorem for *flat* is the commuting diagram in Figure 1. It can be expressed as follows:

$$\forall f : \text{flat} \circ \text{map}(\text{map}(f)) = \text{map}(f) \circ \text{flat}$$

This commuting diagram represents a natural transformation between the functors  $\text{list} \circ \text{list}$  and  $\text{list}$ . It indicates that applying  $f$  to every element of a nested list and then flattening the resulting nested list is equivalent to flattening the list and then applying  $f$  to the flat list. This theorem is always true regardless of the actual definition of *flat* because, if  $f$  were changing in some way the elements of the nested list, then the type of *flat* would not be the polymorphic type given above.

The proof of the following theorem comes directly from [12] (here binary relations between values are restricted to value equalities):

**Theorem 1 (Parametricity Theorem)** *Any expression  $e : \tau$  satisfies  $\mathcal{P}[\tau](e, e)$ , where:*

$$\begin{aligned} \mathcal{P}[\text{basic}](r, s) &\rightarrow r = s \\ \mathcal{P}[\alpha](r, s) &\rightarrow r = f_\alpha(s) \\ \mathcal{P}[\forall \alpha. \tau](r, s) &\rightarrow \forall f_\alpha : \mathcal{P}[\tau](r, s) \\ \mathcal{P}[\tau_1 \times \tau_2](r, s) &\rightarrow \mathcal{P}[\tau_1](\pi_1(r), \pi_1(s)) \wedge \mathcal{P}[\tau_2](\pi_2(r), \pi_2(s)) \\ \mathcal{P}[\tau_1 \rightarrow \tau_2](r, s) &\rightarrow \forall x, y : \mathcal{P}[\tau_1](x, y) \Rightarrow \mathcal{P}[\tau_2](r(x), s(y)) \\ \mathcal{P}[T(\tau)](r, s) &\rightarrow \forall f : (\forall x : \mathcal{P}[\tau](f(x), x)) \Rightarrow r = \text{map}^T(f) s \end{aligned}$$

We assume that each type variable,  $\alpha$ , has a distinct name and is associated with a function  $f_\alpha$  of type  $\alpha_1 \rightarrow \alpha_2$ , where  $\alpha_1$  and  $\alpha_2$  are instances of  $\alpha$ .

For example, the list fold has type:

$$\text{fold} : \forall \alpha, \beta. (\alpha \rightarrow \beta \rightarrow \beta) \rightarrow \beta \rightarrow \text{list}(\alpha) \rightarrow \beta$$

and satisfies the following parametricity theorem:

$$\forall f_\alpha, f_\beta, x, y : (f_\alpha x) \otimes (f_\beta y) = f_\beta(x \oplus y) \Rightarrow \text{fold}(\otimes)(f_\beta x)(\text{map}(f_\alpha) y) = f_\beta(\text{fold}(\oplus) x y)$$

If we set  $f_\alpha = \text{id}$ , where  $\text{id} = \lambda x. x$ , we get:

$$\forall f_\beta, x, y : x \otimes (f_\beta y) = f_\beta(x \oplus y) \Rightarrow \text{fold}(\otimes)(f_\beta x) y = f_\beta(\text{fold}(\oplus) x y) \quad (1)$$

which is the fusion law for list fold [11]. Binding the function  $f_\alpha$ , which corresponds to the type variable  $\alpha$ , to  $\text{id}$  will be done very often in this paper. In general, if we have a type  $t$  that depends on some type variables  $\alpha, \beta$ , etc., i.e.,  $t$  has the form  $t(\alpha, \beta, \dots)$ , then finding the parametricity

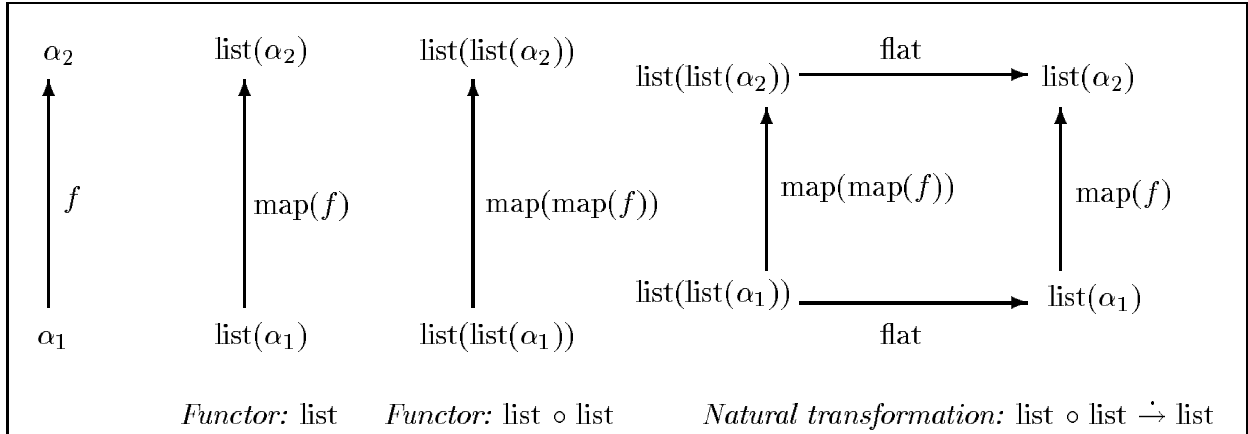


Figure 1: The Parametricity Theorem for  $\text{flat} : \forall \alpha. \text{list}(\text{list}(\alpha)) \rightarrow \text{list}(\alpha)$

theorem, which is very similar to the one for fold. In fact, if a function resembles a fold, then its recursive skeleton is exactly the fold operator. For each such recursive skeleton, we generate the parametricity theorem. Program fusion is achieved by using these theorems alone. In a way, our method generalizes all other methods that use the parametricity theorems of a fixed set of higher-order operators. It is also directly applicable to conventional functional programs. The drawback though is that our method uses many different theorems for program fusion, instead of just a fixed number. But, as we will see in this paper, using these theorems is actually no harder than using the parametricity theorems for folds.

For a functional type, the parametricity theorem becomes a proposition of the form  $\text{premise} \Rightarrow \text{conclusion}$ . The conclusion can only be useful if the premise is proved first, a process that typically involves program synthesis. In this paper we report on a method of using parametricity theorems effectively for program fusion by expressing these theorems in terms of bifunctors. That way, program fusion is performed by calculating the fused program directly, without the need of theorem proving or program synthesis.

We believe that our approach may well turned out to have practical uses for optimizing real functional languages. We also believe that it can be useful for proving equational theorems about functions.

## 2 Background: The Parametricity Theorem

Any function  $f$  of type  $\tau$  satisfies a parametricity theorem (also called *theorem for free* [12]), which is derived directly from the type  $\tau$ . For first-order functions, this theorem states that any polymorphic function is a natural transformation. For example, Figure 1 gives the parametricity theorem for any function  $\text{flat}$  of type:

$$\forall \alpha. \text{list}(\text{list}(\alpha)) \rightarrow \text{list}(\alpha)$$

The parametric type  $\text{list}(\alpha)$  is a functor that maps any type  $\alpha_i$  into the type  $\text{list}(\alpha_i)$  and any function  $f$  of type  $\alpha_1 \rightarrow \alpha_2$  into the function  $\text{map}(f)$  of type  $\text{list}(\alpha_1) \rightarrow \text{list}(\alpha_2)$ . In general, any parametric



# Using the Parametricity Theorem for Program Fusion

Leonidas Fegaras

Department of Computer Science and Engineering  
The University of Texas at Arlington  
416 Yates Street, P.O. Box 19015  
Arlington, TX 76019-19015  
email: *fegaras@cse.uta.edu*

## Abstract

Program fusion techniques have long been proposed as an effective means of improving program performance and of eliminating unnecessary intermediate data structures. This paper proposes a new approach on program fusion that is based entirely on the type signatures of programs. First, for each function, a recursive skeleton is extracted that captures its pattern of recursion. Then, the parametricity theorem of this skeleton is derived, which provides a rule for fusing this function with any function. This method generalizes other approaches that use fixed parametricity theorems to fuse programs.

## 1 Introduction

There is much work recently on using higher-order operators, such as *fold* [11] and *build* [8, 5], to automate program fusion [2] and deforestation [13]. Even though these methods do a good job on fusing programs, they are only effective if programs are expressed in terms of these operators. This limits their applicability to conventional functional languages. To ameliorate this problem, some researchers proposed methods to translate regular functional programs into folds [8]. These methods had a moderate success so far, mostly because they can apply to simple functions only.

The main reason for using these higher-order operators is that they satisfy some powerful theorems, which facilitate program optimization. But there is nothing special about these theorems. They are parametricity theorems [12] that are derived exclusively from the types of these operators. Any function satisfies a parametricity theorem. The difference is that most functions are not *sufficiently polymorphic* and, thus, their parametricity theorems are usually trivial.

This paper proposes a new approach on fusing programs. Instead of trying to express a function in terms of a particular higher-order operator, such as *fold*, we generate an individualized higher-order operator for this function. This operator, called the *recursive skeleton* of this function, captures its pattern of recursion. It is specific to this function only and may not be suitable for any other function. This operator is polymorphic enough to satisfy a useful parametricity